

Carlos Ivorra Castillo

ANÁLISIS MATEMÁTICO

Si una cantidad no negativa fuera tan pequeña que resultara menor que cualquier otra dada, ciertamente no podría ser sino cero. A quienes preguntan qué es una cantidad infinitamente pequeña en matemáticas, nosotros respondemos que es, de hecho, cero. Así pues, no hay tantos misterios ocultos en este concepto como se suele creer. Esos supuestos misterios han convertido el cálculo de lo infinitamente pequeño en algo sospechoso para mucha gente. Las dudas que puedan quedar las resolveremos por completo en las páginas siguientes, donde explicaremos este cálculo.

LEONHARD EULER

Índice General

Preámbulo	ix
Introducción	xiii
Capítulo I: Los números reales	1
1.1 Cuerpos métricos	3
1.2 Convergencia de sucesiones	13
1.3 Desarrollos decimales	19
1.4 Sucesiones de Cauchy	25
1.5 Cuerpos ordenados completos	29
1.6 El cardinal del continuo	35
Capítulo II: Topología	39
2.1 Espacios topológicos	39
2.2 Bases y subbases	45
2.3 Productos y subespacios	48
2.4 Algunos conceptos topológicos	52
2.5 Continuidad	58
2.6 Límites de funciones	70
2.7 Convergencia de sucesiones	80
2.8 Sucesiones y series numéricas	83
Capítulo III: Compacidad, conexión y completitud	97
3.1 Espacios compactos	97
3.2 Espacios localmente compactos	103
3.3 Espacios conexos	107
3.4 Espacios completos	113
3.5 Espacios de Hilbert	117
3.6 Espacios de funciones	127
3.7 El teorema de Baire	134
Capítulo IV: Cálculo diferencial de una variable	139
4.1 Derivación	139
4.2 Cálculo de derivadas	142
4.3 Propiedades de las funciones derivables	146

4.4	La diferencial de una función	155
4.5	El teorema de Taylor	158
4.6	Series de potencias	163
4.7	La función exponencial	168
4.8	Las funciones trigonométricas	175
4.9	Las funciones hiperbólicas	183
4.10	Primitivas	186
Capítulo V: Cálculo diferencial de varias variables		195
5.1	Diferenciación	195
5.2	Propiedades de las funciones diferenciables	204
5.3	Curvas parametrizables	216
Capítulo VI: Introducción a las variedades diferenciables		239
6.1	Variedades	240
6.2	Espacios tangentes, diferenciales	248
6.3	La métrica de una variedad	255
6.4	Geodésicas	260
6.5	Superficies	264
6.6	La curvatura de Gauss	267
Capítulo VII: Ecuaciones diferenciales ordinarias		275
7.1	Ecuaciones diferenciales de primer orden	278
7.2	Ecuaciones diferenciales de orden superior	286
7.3	Aplicaciones	291
Capítulo VIII: Teoría de la medida I		301
8.1	La medida de Jordan	302
8.2	Medidas	318
8.3	La medida de Lebesgue	328
8.4	Funciones medibles	334
8.5	La integral de Lebesgue	337
8.6	La integral de Lebesgue en \mathbb{R}	345
Capítulo IX: Teoría de la medida II		351
9.1	Producto de medidas	351
9.2	El teorema de Riesz	360
9.3	Espacios L^p	367
9.4	Medidas signadas	372
9.5	Derivación de medidas	382
9.6	El teorema de cambio de variable	386
9.7	Integración en variedades	395
Apéndice A: La completación de un espacio métrico		405

Apéndice B: Fracciones continuas	413
B.1 Propiedades básicas	413
B.2 Desarrollos de irracionales cuadráticos	418
B.3 Transformaciones modulares	421
B.4 El espacio de Baire	423
B.5 La fracción continua de e	425
Apéndice C: Resumen de dinámica clásica	433
C.1 El espacio y el tiempo	433
C.2 Cinemática de una partícula puntual	437
C.3 Fuerzas	444
C.4 Trabajo y energía	450
C.5 Dinámica de un sistema de partículas	454
C.6 Distribuciones continuas de materia	462
Bibliografía	465
Índice de Materias	466

Preámbulo

Este libro es una versión extendida de mi libro original del mismo título. Además de haber añadido bastante material nuevo, otra diferencia sustancial respecto de la versión anterior es que ahora las nuevas versiones de los libros *Álgebra* [Al], *Geometría* [G] y *Análisis matemático* [An] están concebidas para ser leídas simultáneamente, mientras que las anteriores podían leerse sucesivamente en el orden indicado. Esto ha hecho que parte del material original haya pasado de un libro a otro en las versiones nuevas.

Más concretamente, si los capítulos de los tres libros se leen de arriba hacia abajo y de izquierda a derecha según la disposición de la tabla de la página siguiente, cada uno sólo requiere el conocimiento de los capítulos precedentes, salvo por unas pocas excepciones indicadas después de la tabla. Consideramos que este orden es mucho más natural en la medida en que imita lo que sería el estudio simultáneo de tres asignaturas y se aprovecha dicha simultaneidad para enfatizar las relaciones entre ellas.

El libro [An] es el único que no cubre todos los contenidos de la versión anterior, sino que los últimos capítulos de dicha versión serán tratados en nuevas versiones de mi serie de libros, empezando por el de *Geometría diferencial*.

El primer capítulo de [Al] es una introducción a la teoría de conjuntos, cuyos aspectos más técnicos (los relacionados con el axioma de elección y la teoría de cardinales infinitos) se han relegado a dos apéndices. La teoría descrita es la teoría de Zermelo, que resulta más que suficiente para formalizar los contenidos de los tres libros. El único inconveniente es que “se queda corta” para desarrollar plenamente la teoría de cardinales infinitos, pero hemos preferido reducirla a lo imprescindible, aun al precio de no poder enunciar con total precisión algunos resultados sobre rango y dimensión de módulos y espacios vectoriales de dimensión infinita que, aunque resulta natural presentarlos al tratar estos conceptos, no son realmente necesarios en ningún momento.

Los contenidos de [Al I] (y sus apéndices) sirven de base a los tres libros. A su vez, los capítulos [G II] y [An I] se apoyan en las estructuras algebraicas introducidas en [Al II]. En estos dos capítulos se presentan dos construcciones alternativas de los números reales, mediante secciones de Dedekind en [G II] y mediante sucesiones de Cauchy en [An I]. Por otra parte, el capítulo [An II] se apoya en la geometría analítica expuesta en [G IV].

ÁLGEBRA	GEOMETRÍA	ANÁLISIS
Al I Preliminares		
Al II Anillos	G I Geometría absoluta	An I Números reales
Al III Aritmética	G II Geometría arquimediana	
Al IV Aplicaciones	G III Geometría euclídea	
Al V Módulos	G IV Geometría analítica	An II Topología
Al VI Grupos	G V Complejos y cuaternios	An III Compacidad, ...
Al VII Cuerpos	G VI Regla y compás	An IV Cálculo una variable
Al VIII Álgebra lineal	G VII Bijecciones afines	An V Cálculo varias variables
Al IX Ecuaciones	G VIII Geometría afín	An VI Variedades
Al X Enteros algebraicos	G IX Geometría proyectiva	An VII Ecuaciones diferenciales
Al XI Enteros cuadráticos	G X Cónicas	An VIII Medida I
Al XII Factorización ideal	G XI Geometría parabólica	An IX Medida II
Al XIII Complementos	G XII Geometría hiperbólica	An Ap A Compleción de un e.m.
Al Ap A Ax. de elección	G XIII Geometría elíptica	An Ap B Fracciones continuas
Al Ap B Conjuntos infinitos	G Ap A Geometría inversiva	An Ap C Dinámica clásica

Las únicas dependencias que no respetan el orden indicado en la tabla precedente son las siguientes:

- La sección [G 7.7] usa la medida de Jordan definida en [An 9.1], pero dicha sección puede leerse antes de [G 7.7], pues no requiere nada intermedio.
- El capítulo [An VII] usa en un ejemplo las cónicas definidas en la sección [G 10.1], pero esta sección (no ya el resto del capítulo) puede leerse tras [G IV] y resulta incluso recomendable hacerlo, como ilustración de las posibilidades de la geometría analítica.

Las líneas horizontales en la tabla separan bloques temáticos. El segundo bloque de [Al], después del capítulo de preliminares conjuntistas al que ya hemos hecho referencia, contiene una presentación de la aritmética básica desde el punto de vista del álgebra abstracta, junto con aplicaciones que conectan este

enfoque abstracto con resultados clásicos. En el tercer bloque se introducen nuevas estructuras abstractas con resultados que se aplican principalmente a la geometría (en [G]) y a la aritmética en el cuarto bloque, que contiene una introducción a la teoría algebraica de números. El libro termina con un capítulo en el que se recopilan algunos resultados que no han sido necesarios en los capítulos precedentes pero que son relevantes de cara a estudios más avanzados.

El primer bloque de [G] contiene un tratamiento axiomático de la geometría euclídea, el segundo desarrolla los elementos básicos de la geometría analítica, el tercero está dedicado a la geometría proyectiva y el cuarto a las geometrías no euclídeas.

Por último, [An] está dividido en tres bloques, dedicados respectivamente a la topología, al cálculo diferencial y al cálculo integral. El libro termina con tres apéndices, el primero de los cuales es una prolongación técnica del capítulo [An I] con material que no es necesario para los capítulos posteriores, el segundo expone la teoría de las fracciones continuas, en la que se combinan aspectos aritméticos con aspectos topológicos y, por último, hemos considerado oportuno incluir un resumen de la dinámica clásica que puede servir al lector para asimilar mejor las numerosas aplicaciones a la física presentadas en el libro. En realidad, todos los conceptos físicos involucrados se van explicando en los propios ejemplos a medida que van siendo necesarios, pero tal vez el lector no familiarizado con la física prefiera una exposición concentrada en unas pocas páginas que le sirva de referencia.

Introducción

En el siglo XVII Newton y Leibniz descubren independientemente el *análisis matemático* o *cálculo infinitesimal*, una potentísima herramienta que revolucionó el tratamiento matemático de la física y la geometría, y que más tarde impregnaría las más diversas ramas de la matemática, como la estadística o la teoría de números.

Esencialmente, el cálculo infinitesimal consistía por una parte en *analizar* o descomponer la dependencia entre varias magnitudes estudiando el comportamiento de unas al variar o *diferenciar* levemente otras (lo que constituía el *cálculo diferencial*) y por otra parte en *integrar* los resultados diferenciales para obtener de nuevo resultados globales sobre las magnitudes en consideración (el llamado *cálculo integral*).

Es difícil que un lector que no tenga ya algunas nociones de cálculo pueda entender el párrafo anterior, pero las nuevas ideas eran aún más difíciles de entender de la pluma de sus descubridores. El primer libro de texto que se publicó con el fin de explicarlas sistemáticamente fue el “Análisis” del marqués de L’Hôpital. Veamos algunos pasajes:

La parte infinitamente pequeña en que una cantidad variable es aumentada o disminuida de manera continua, se llama la diferencial de esta cantidad.

Siguiendo la notación leibniziana, L’Hôpital explica que la letra d se usa para representar uno de estos incrementos infinitamente pequeños de una magnitud, de modo que dx representa un incremento diferencial de la variable x , etc.

En ningún momento se precisa qué debemos entender por un aumento infinitamente pequeño de una cantidad, pero en compensación se presentan varias reglas para tratar con diferenciales. Por ejemplo:

Postúlese que dos cantidades cuya diferencia es una cantidad infinitamente pequeña pueden intercambiarse una por la otra; o bien (lo que es lo mismo) que una cantidad que está incrementada o disminuida solamente en una cantidad infinitamente menor, puede considerarse que permanece constante.

Así, por ejemplo, si analizamos el incremento infinitesimal que experimenta un producto xy cuando incrementamos sus factores, obtenemos

$$d(xy) = (x + dx)(y + dy) - xy = x dy + y dx + dx dy = x dy + y dx,$$

donde hemos despreciado el infinitésimo doble $dx dy$ porque es infinitamente menor que los infinitésimos simples $x dy$ e $y dx$.

Es fácil imaginar que estos razonamientos infinitesimales despertaron sospechas y polémicas. Baste citar el título del panfleto que en 1734 publicó el obispo de Berkeley:

El analista, o discurso dirigido a un matemático infiel, donde se examina si los objetos, principios e inferencias del análisis moderno están formulados de manera más clara, o deducidos de manera más evidente, que los misterios religiosos y los asuntos de la fe.

En esta fecha el cálculo infinitesimal tenía ya más de medio siglo de historia. La razón por la que sobrevivió inmune a estas críticas y a la vaguedad de sus fundamentos es que muchos de sus razonamientos infinitesimales terminaban en afirmaciones que no involucraban infinitésimos en absoluto, y que eran confirmados por la física y la geometría. Por ejemplo, consideremos la circunferencia formada por los puntos que satisfacen la ecuación

$$x^2 + y^2 = 25.$$

Aplicando la regla del producto que hemos “demostrado” antes al caso en que los dos factores son iguales obtenemos que $dx^2 = 2x dx$ e igualmente será $dy^2 = 2y dy$. Por otra parte, $d25 = 0$, pues al incrementar la variable x la constante 25 no se ve incrementada en absoluto. Si a esto añadimos que la diferencial de una suma es la suma de las diferenciales resulta la ecuación diferencial

$$2x dx + 2y dy = 0,$$

de donde a su vez

$$\frac{dy}{dx} = -\frac{x}{y}.$$

Esto significa que si tomamos, por ejemplo, el punto $(3, 4)$ de la circunferencia e incrementamos infinitesimalmente su coordenada x , la coordenada y disminuirá en $3/4 dx$. Notemos que esto es falso para cualquier incremento finito de la variable x , por pequeño que sea, pues si valiera para incrementos suficientemente pequeños resultaría que la circunferencia contendría un segmento de la recta

$$y - 4 = -\frac{3}{4}(x - 3),$$

lo cual no es el caso. Vemos que ésta se comporta igual que la circunferencia para variaciones infinitesimales de sus variables alrededor del punto $(3, 4)$, aunque difiere de ella para cualquier variación finita. La interpretación geométrica es que se trata de la recta tangente a la circunferencia por el punto $(3, 4)$.

El argumento será nebuloso y discutible, pero lo aplastante del caso es que nos proporciona un método sencillo para calcular la tangente a una circunferencia por uno cualquiera de sus puntos. De hecho el método se aplica a cualquier curva que pueda expresarse mediante una fórmula algebraica razonable, lo que supera con creces a las técnicas con las que contaba la geometría analítica antes del cálculo infinitesimal.

A lo largo del siglo XIX la matemática emprendió un proceso de fundamentación que terminó con una teoría formal donde todos los conceptos están perfectamente definidos a partir de unos conceptos básicos, los cuales a su vez están completamente gobernados por unos axiomas precisos. Las ambigüedades del cálculo infinitesimal fueron el motor principal de este proceso. En los años sesenta del siglo XX se descubrió que una delicada teoría lógica, conocida como *análisis no estándar*, permite definir rigurosamente cantidades infinitesimales con las que fundamentar el cálculo a la manera de Leibniz y L'Hôpital, pero no es ése el camino habitual ni el que nosotros vamos a seguir. Lo normal es erradicar los infinitésimos de la teoría, pero no así el formalismo infinitesimal. En ocasiones los símbolos dy , dx aparecen en ciertas definiciones “en bloque”, sin que se les pueda atribuir un significado independiente, como cuando se define la derivada de una función $y = y(x)$ mediante

$$\frac{dy}{dx} = \lim_{\Delta x \rightarrow 0} \frac{y(x + \Delta x) - y(x)}{\Delta x}.$$

De este modo, el cociente de diferenciales tiene el mismo significado que para Leibniz, en el sentido de que al calcularlo obtenemos el mismo número o la misma función que él obtenía, pero con la diferencia de que ya no se trata de un cociente de diferenciales, no es un cociente de nada. La definición anterior nos permite hablar de dy/dx , pero no de dy o de dx .

No obstante se puede ir más lejos y dar una definición adecuada de dx y dy de modo que se pueda probar la equivalencia

$$\frac{dy}{dx} = f(x) \iff dy = f(x) dx.$$

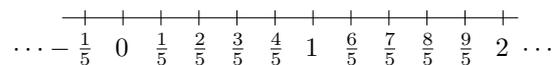
Es algo parecido al paso de una relación algebraica como $xy^2 = x + 4y^3$, donde x e y son, digamos, números reales indeterminados, a la misma expresión entendida como una igualdad de polinomios, donde ahora x e y son indeterminadas en un sentido matemático muy preciso. Por ejemplo, según una definición habitual del anillo de polinomios $\mathbb{R}[x, y]$, la indeterminada x es la aplicación de los pares de números naturales en \mathbb{R} dada por $x(1, 0) = 1$ y $x(i, j) = 0$ para cualquier otro par, es decir, algo que en nada recuerda a “un número real indeterminado”.

En este libro veremos cómo el análisis matemático moderno da sentido a las fórmulas y métodos clásicos evitando en todo momento las “cantidades infinitesimales”, a la vez que presentamos los resultados más importantes del análisis matemático real. Concretamente abordamos el cálculo diferencial e integral de una y varias variables y las ecuaciones diferenciales ordinarias.

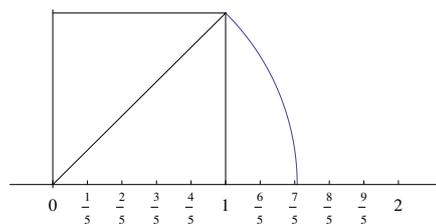
Capítulo I

Los números reales

El conjunto \mathbb{Q} de los números racionales tiene una representación geométrica: en una recta cualquiera, seleccionamos arbitrariamente dos de sus puntos, P_0 y P_1 y, a partir de ahí, a cada número racional $r = m/n$ le asignamos el punto P_r que resulta de dividir el segmento $\overline{P_0P_1}$ en n partes iguales y, desde P_0 , dar m pasos de dicha longitud hacia P_1 si $m > 0$ y en sentido opuesto si $m < 0$. La figura muestra la representación geométrica de los menores números racionales con denominador 5:



Pero la geometría nos enseña también que los números racionales no “llenan” la recta, sino que, por ejemplo, si abatimos sobre la recta la diagonal de un cuadrado de lado unitario, obtenemos un punto, que en la figura se ve que queda entre $7/5$ y $8/5$, que según el teorema de Pitágoras debería corresponder a un número x tal que $x^2 = 2$:



Pero sabemos que ningún número racional cumple $x^2 = 2$. Esto nos obliga a considerar una clase más amplia de números. Podemos pensar que la inexistencia de una raíz cuadrada de 2 en \mathbb{Q} delata la presencia de un “hueco” en \mathbb{Q} al que le podemos asignar una posición en la recta con toda precisión. La figura muestra que $7/5$ está “a la izquierda del hueco”, mientras que $8/5$ está “a la derecha del hueco”, lo cual deja poco margen a su situación, pero este margen se

puede reducir tanto como queramos considerando denominadores mayores. Por ejemplo, es fácil ver que

$$\left(\frac{141}{100}\right)^2 < 2 < \left(\frac{142}{100}\right)^2,$$

lo que nos indica que una hipotética raíz cuadrada de 2 debería estar situada entre las dos bases:

$$\frac{141}{100} < \sqrt{2} < \frac{142}{100},$$

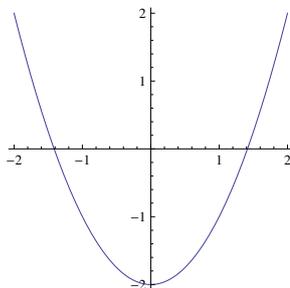
lo que “localiza el hueco” en \mathbb{Q} con una precisión de una centésima.

En la sección 2.2 de [G] hemos construido el cuerpo \mathbb{R} de los números reales, y hemos demostrado que son suficientes para medir cualquier segmento si admitimos que las rectas satisfacen la propiedad de Arquímedes, lo cual significa que no tienen puntos infinitamente alejados unos de otros. A su vez, esto se traduce en que la representación geométrica de \mathbb{Q} se puede extender de modo que cada punto de la recta tenga asociado un número real. Dicha asociación no es arbitraria, sino que se hace en concordancia con la estructura algebraica de \mathbb{R} , lo cual se plasma esencialmente en que la longitud del segmento que une los puntos P_α y P_β asociados a dos números reales α y β es precisamente $|\alpha - \beta|$ (teorema [G 2.18]).

No obstante, para desarrollar la geometría clásica no es preciso suponer que todos los números reales tienen un punto asociado en la recta, sino que es suficiente considerar rectas que, además de tener puntos correspondientes con los números racionales, tengan también algunos puntos “extra” correspondientes a unos pocos números irracionales más, que incluyan a $\sqrt{2}$ y otros similares cuya existencia es necesaria a causa del teorema de Pitágoras.

Por el contrario, las técnicas del análisis matemático que vamos a desarrollar aquí requieren trabajar con todos los números reales y, por consiguiente, para dar una interpretación geométrica a los resultados que vamos a demostrar, tenemos que pensar en rectas que cumplan el axioma de continuidad presentado en la sección 2.5 de [G], que afirma precisamente que todos los números reales tienen asignado un lugar en cualquier recta.

Para entender por qué el análisis necesita que “no haya huecos” en la recta vamos a presentar el mismo fenómeno desde otro punto de vista: Uno de nuestros objetos de estudio van a ser las funciones entre números, como $f(x) = x^2 - 2$. Consideremos $f : \mathbb{Q} \rightarrow \mathbb{Q}$ definida de este modo y observemos su gráfica:



Posiblemente el lector estará familiarizado con este tipo de representaciones gráficas de funciones y no es necesario añadir muchas explicaciones sobre su interpretación: la curva que se muestra está formada por los puntos que al proyectarse verticalmente sobre el eje horizontal llevan a un punto que se corresponde con un número racional x y al proyectarse horizontalmente sobre el eje vertical caen sobre el punto correspondiente a $f(x)$.

La curva de la gráfica corta al eje vertical cuando $x = 0$ en el punto -2 , y la figura sugiere que “debe cortar” al eje horizontal en dos puntos, uno situado entre $7/5$ y $8/5$ y el otro en posición simétrica respecto del 0 . Sin embargo, estamos de nuevo en el mismo caso: no existen números racionales x que cumplan $x^2 - 2 = 0$, luego nos vemos obligados a concluir que la curva mostrada no corta al eje horizontal... salvo que estemos dispuestos a considerar a f definida sobre un cuerpo mayor que \mathbb{Q} .

La diferencia está en que la geometría clásica sólo requiere que existan los irracionales que pueden obtenerse a partir de \mathbb{Q} mediante sumas, restas, productos, cocientes y extracción de raíces cuadradas, mientras que, al considerar funciones más sofisticadas que $x^2 - 1$, el análisis matemático se encuentra con la necesidad de “rellenar” muchos más de los “huecos” que los números racionales dejan en la recta. De hecho, necesita rellenarlos todos. En este capítulo vamos a estudiar los “huecos” de \mathbb{Q} con técnicas analíticas, mucho más minuciosas para este fin que las que proporciona la geometría.

Dado que las técnicas que vamos a presentar dan lugar a una construcción de \mathbb{R} alternativa a la presentada en [G], a lo largo de este capítulo no vamos a suponer conocido el cuerpo \mathbb{R} construido allí. En el apéndice A daremos otra construcción de \mathbb{R} basada en las técnicas y conceptos que vamos a presentar aquí. Esta construcción se debe a Cantor, mientras que la de [G] se debe a Dedekind.

1.1 Cuerpos métricos

Muchas de las ideas que vamos a presentar en este capítulo son aplicables en muchos otros contextos de interés, y por ello es conveniente trabajar en un marco general que realmente no supone ninguna dificultad añadida, sino simplemente observar que podemos aislar como axiomas las propiedades de \mathbb{Q} que vamos a utilizar. Por ejemplo, empezamos introduciendo el concepto de “intervalo”, que tiene sentido, no sólo en \mathbb{Q} , sino en cualquier conjunto totalmente ordenado X :

Definición 1.1 Sea X un conjunto totalmente ordenado. Llamaremos *intervalos* en X a los conjuntos siguientes, para todo $a, b \in X$:

$$\begin{aligned}]a, b[&= \{x \in X \mid a < x < b\}, & [a, b] &= \{x \in X \mid a \leq x \leq b\}, \\]a, b] &= \{x \in X \mid a < x \leq b\}, & [a, b[&= \{x \in X \mid a \leq x < b\}, \\]-\infty, b[&= \{x \in X \mid x < b\}, &]a, +\infty[&= \{x \in X \mid a < x\}, \\]-\infty, b] &= \{x \in X \mid x \leq b\}, & [a, +\infty[&= \{x \in X \mid a \leq x\}, \\ & &]-\infty, +\infty[&= X. \end{aligned}$$

El elemento a (en los intervalos en los que interviene) se llama *extremo inferior* del intervalo, mientras que b (cuando procede) es el *extremo superior*. Los intervalos de la forma $]a, b[$, incluso si a o b es infinito, se llaman *intervalos abiertos*, mientras que los de tipo $[a, b]$ se llaman *intervalos cerrados*.

Observemos que si X tiene máximo M , entonces

$$]a, +\infty[=]a, M], \quad [a, +\infty[= [a, M],$$

y si tiene mínimo m entonces

$$]-\infty, b[= [m, b[, \quad]-\infty, b] = [m, b],$$

y si tiene máximo y mínimo entonces $]-\infty, +\infty[= [m, M]$, por lo que los intervalos con extremos infinitos sólo son relevantes en ausencia de máximo o de mínimo. En tal caso son conjuntos no acotados, y por ello se llaman *intervalos no acotados*.

Por ejemplo, la gráfica de la página 2 “parece” dividir la recta en tres intervalos, pero no es así si identificamos la recta con \mathbb{Q} . En tal caso tenemos que $\mathbb{Q} = A \cup B \cup C$, donde

$$A = \{r \in \mathbb{Q} \mid r < 0, r^2 > 2\}, B = \{r \in \mathbb{Q} \mid r^2 < 2\}, C = \{r \in \mathbb{Q} \mid r > 0, r^2 > 2\}.$$

Se cumple que todo elemento de A es menor que todo elemento de B , y todo elemento de B es menor que todo elemento de C , pero A , B y C no son intervalos porque no tienen extremos. No hay ningún número racional en el que podamos decir que termina A y empieza B , ni otro donde termine B y empiece C . Los conjuntos que “parecen intervalos sin serlo” son una de muchas formas de “señalar huecos” en \mathbb{Q} . Vamos a ver varias más.

Ahora vamos a hacer algunas observaciones sobre \mathbb{Q} que en realidad son válidas sobre todo cuerpo ordenado R que cumpla una condición adicional que \mathbb{Q} satisface trivialmente:

Definición 1.2 Un anillo ordenado R es *arquimediano* si \mathbb{N} no está acotado¹ en R , es decir, si para todo $x \in R$ existe un $n \in \mathbb{N}$ tal que $x < n$.

Obviamente \mathbb{Q} es un cuerpo ordenado arquimediano, pues si $m/n \geq 0$, entonces $m/n < m+1$. Entre otros conceptos asociados a la propiedad arquimediana está el de ‘parte entera’:

Teorema 1.3 Si R es un anillo ordenado arquimediano, para cada $x \in R$ existe un único $m \in \mathbb{Z}$ tal que $m \leq x < m+1$.

DEMOSTRACIÓN: Si $x \geq 0$ tomamos el mínimo número natural k tal que $x < k$, que existe por la propiedad arquimediana. Entonces $k-1 \leq x < k$. Si

¹Recordemos de [Al, Sección 3.6] que todo anillo ordenado tiene característica 0 y, por consiguiente, contiene a \mathbb{Z} como subanillo, y todo cuerpo ordenado contiene a \mathbb{Q} como subcuerpo.

$x < 0$, entonces $-x > 0$ y tomamos el mínimo número natural k tal que $-x \leq k$, con lo que $k-1 < -x \leq k$, y a su vez $-k \leq x < -k+1$. En ambos casos hemos encontrado un entero m que cumple lo pedido ($m = k-1$ en el primer caso y $m = -k$ en el segundo).

Si también $m' \leq x < m'+1$, pero $m \neq m'$, podemos suponer que $m < m'$, pero entonces $m+1 \leq m'$ y tenemos una contradicción:

$$x < m+1 \leq m' \leq x.$$

Esto prueba la unicidad. ■

Definición 1.4 Si R es un anillo ordenado, llamaremos *parte entera* en R a la función $E : R \rightarrow \mathbb{Z}$ que a cada $x \in R$ le asigna el único número entero que cumple $E[x] \leq x < E[x] + 1$.

Así, por ejemplo, en \mathbb{Q} tenemos que $E[17/3] = 5$, $E[-17/3] = -6$.

El elemento $F[x] = x - E[x]$ se llama *parte fraccionaria* de x , de modo que todo $x \in R$ se descompone de forma única como

$$x = E[x] + F[x], \quad \text{con } E[x] \in \mathbb{Z}, \quad 0 \leq F[x] < 1.$$

En el caso de un cuerpo podemos decir más:

Teorema 1.5 Si R es un cuerpo ordenado arquimediano, entonces, para todo par de elementos $a < b$ en R existe un $q \in \mathbb{Q}$ tal que $a < q < b$.

DEMOSTRACIÓN: Sea n un número natural tal que

$$\frac{1}{b-a} < n,$$

con lo que $1/n < b-a$. Sea $m = E[na] + 1$, de modo que $m-1 \leq na < m$, luego

$$a < \frac{m}{n} = \frac{m-1}{n} + \frac{1}{n} < a + b - a = b.$$

Así, $q = m/n$ cumple $a < q < b$. ■

Nota En lo que sigue vamos a presentar muchos resultados válidos para “un cierto cuerpo ordenado arquimediano R ”. La finalidad de este marco de trabajo es incluir simultáneamente los dos únicos casos que nos van a interesar: $R = \mathbb{Q}$ y $R = \mathbb{R}$, pero sin suponer conocido \mathbb{R} , para que los resultados que presentamos puedan ser usados más tarde en una construcción de \mathbb{R} . ■

Sea R cualquier cuerpo ordenado arquimediano y sea $x_0 \in R$ tal que $x_0 > 0$. Vamos a razonar en general, pero, paralelamente, vamos a ilustrar el razonamiento con el supuesto concreto en que $x_0^2 = 2$.

Podemos descomponer $x_0 = a_0 + x_1$, donde $a_0 = E[x_0]$ y $x_1 = F[x_0]$. Así, $x_0 \in [a_0, a_0 + 1]$. Además, el número natural a_0 admite un desarrollo decimal:²

$$a_0 = \sum_{n=0}^m c_n 10^n, \quad 0 \leq c_n < 10.$$

Por otra parte, $x_1 = F[x_0]$ cumple $0 \leq x_1 < 1$.

En el ejemplo concreto, como $1^2 < x_0^2 = 2 < 2^2$, es $1 < x_0 < 2$, luego $a_0 = 1$ y tenemos que $x_0 \in [1, 2]$.

Ahora dividimos el intervalo $[0, 1]$ en 10 partes iguales:

$$0 < \frac{1}{10} < \frac{2}{10} < \dots < \frac{10}{10} = 1$$

y determinamos en cuál de los 10 subintervalos está x_1 , es decir, para qué c_{-1} se cumple

$$\frac{c_{-1}}{10} \leq x_1 < \frac{c_{-1} + 1}{10}.$$

Equivalentemente, $c_{-1} = E[10x_1]$, de modo que $0 \leq c_{-1} < 10$ y además se cumple que $10x_1 = c_{-1} + F[10x_1]$, lo que a su vez se traduce en que

$$x_1 = \frac{c_{-1}}{10} + x_2, \quad \text{con } x_2 = \frac{F[10x_1]}{10}, \quad 0 \leq x_2 < \frac{1}{10}.$$

En total:

$$x_0 = \sum_{n=0}^m c_n 10^n + \frac{c_{-1}}{10} + x_2,$$

luego

$$a_1 = \sum_{n=0}^m c_n 10^n + \frac{c_{-1}}{10} \leq x_0 < \sum_{n=0}^m c_n 10^n + \frac{c_{-1} + 1}{10}.$$

Así hemos situado a x_0 en un intervalo $[a_1, a_1 + 1/10]$. Más precisamente:

$$x_0 \in [a_1, a_1 + \frac{1}{10}] \subset [a_0, a_0 + 1]$$

$$a_0 \leq a_1 \leq x_0 < a_1 + \frac{1}{10} \leq a_0 + 1.$$

En nuestro ejemplo, como

$$\left(1 + \frac{4}{10}\right)^2 = \left(\frac{14}{10}\right)^2 < 2 < \left(\frac{15}{10}\right)^2 = \left(1 + \frac{5}{10}\right)^2,$$

tenemos que

$$1 + \frac{4}{10} < x_0 < 1 + \frac{5}{10},$$

luego $c_0 = 1$, $c_{-1} = 4$.

²Todo lo que sigue es válido si consideramos desarrollos en cualquier base $k \geq 2$, de modo que 10 puede interpretarse como $10_k = k$. En el ejemplo concreto tomamos la base 10 usual.

Seguidamente dividimos el intervalo $[0, 1/10]$ en 10 partes iguales:

$$0 < \frac{1}{100} < \frac{2}{100} < \dots < \frac{10}{100} = \frac{1}{10}$$

y determinamos en cuál de ellas se encuentra x_2 , lo que nos da un c_{-2} tal que

$$\frac{c_{-2}}{10^2} \leq x_2 < \frac{c_{-2} + 1}{10^2}.$$

Equivalentemente, $c_{-2} = E[10^2 x_2]$, con lo que $0 \leq c_{-2} < 10$ y

$$x_2 = \frac{c_{-2}}{10^2} + x_3, \quad \text{con } x_3 = \frac{F[10^2 x_2]}{10^2}, \quad 0 \leq x_3 < \frac{1}{10^2}.$$

En total:

$$x_0 = \sum_{n=0}^m c_n 10^n + \frac{c_{-1}}{10} + \frac{c_{-2}}{10^2} + x_3,$$

luego

$$a_2 = \sum_{n=0}^m c_n 10^n + \frac{c_{-1}}{10} + \frac{c_{-2}}{10^2} \leq x_0 < \sum_{n=0}^m a_n 10^n + \frac{c_{-1}}{10} + \frac{c_{-2} + 1}{10^2}$$

y así hemos situado a x_0 en un intervalo

$$x_0 \in [a_2, a_2 + \frac{1}{10^2}] \subset [a_1, a_1 + \frac{1}{10}] \subset [a_0, a_0 + 1].$$

En nuestro ejemplo sucede que

$$\left(1 + \frac{4}{10} + \frac{1}{10^2}\right)^2 = \left(\frac{141}{100}\right)^2 < 2 < \left(\frac{142}{100}\right)^2 = \left(1 + \frac{4}{10} + \frac{2}{100}\right)^2,$$

por lo que

$$1 + \frac{4}{10} + \frac{1}{10^2} = \frac{141}{100} < x_0 < \frac{142}{100} = 1 + \frac{4}{10} + \frac{2}{100},$$

es decir, $c_0 = 1$, $c_{-1} = 4$, $c_{-2} = 1$.

Este proceso se puede repetir indefinidamente. Concretamente, una simple inducción basada en el razonamiento anterior prueba el teorema siguiente:

Teorema 1.6 *Si $x = x_0 > 0$ es un elemento de un cuerpo ordenado arquimediano R , consideramos el desarrollo en base 10 de $E[x] = \sum_{n=0}^m c_n 10^n$ y, para cada $n \geq 1$, definimos $x_n = F[10^{n-1} x_{n-1}]/10^{n-1}$, $c_{-n} = E[10^n x_n]$,*

$$a_N = \sum_{n=0}^m c_n 10^n + \sum_{n=1}^N c_{-n} 10^{-n} \in \mathbb{Q} \quad \text{y} \quad b_N = a_N + 10^{-N} \in \mathbb{Q}.$$

Entonces $0 \leq c_n < 10$ y la sucesión $\{[a_N, b_N]\}_{N \in \mathbb{N}}$ de intervalos cumple

$$x \in [a_{N+1}, b_{N+1}] \subset [a_N, b_N]$$

para todo N .

Los intervalos hay que entenderlos en R . Si los consideramos en \mathbb{Q} y $x \notin \mathbb{Q}$, entonces los intervalos “rodean un hueco” de \mathbb{Q} al que delimitan con precisión creciente, pues $b_N - a_N = 1/10^N$.

Dados números naturales $0 \leq a_i < 10$, usaremos la *notación decimal*.³

$$a_m \cdots a_0 . a_{-1} \cdots a_{-N} = \sum_{n=0}^m a_n 10^n + \sum_{n=1}^N a_{-n} 10^{-n} \in \mathbb{Q}.$$

En estos términos hemos probado que

$$x_0 \in [1.41, 1.42] \subset [1.4, 1.5] \subset [1, 2],$$

y la sucesión de intervalos puede prolongarse indefinidamente, por ejemplo, hasta

$$1.41421356237310 < x_0 < 1.41421356237311$$

Pero lo más destacado es que, aunque hemos supuesto que x_0 era una hipotética raíz cuadrada de 2 en un hipotético cuerpo ordenado arquimediano R , todas las cuentas que determinan la sucesión de decimales dependen únicamente de la aritmética de \mathbb{Q} . Vemos así que, aunque en \mathbb{Q} no haya raíces cuadradas de 2, lo cierto es que \mathbb{Q} “determina” una raíz cuadrada de 2. Nos gustaría escribir:

$$\sqrt{2} = 1.41421356237310 \dots,$$

pero dar sentido a esos puntos suspensivos requiere una serie de conceptos no triviales que vamos a presentar seguidamente.

Sucesiones Recordemos que una sucesión en un conjunto M es una aplicación $x : \mathbb{N} \rightarrow M$, que representaremos más habitualmente como $\{x_n\}_{n \in \mathbb{N}}$ o, para adaptarnos a lo habitual en análisis, $\{x_n\}_{n=0}^{\infty}$. Con cualquiera de estas notaciones, lo que tenemos es una sucesión infinita

$$x_0, \quad x_1, \quad x_2, \quad x_3, \quad x_4, \quad \dots$$

de elementos de M . En realidad, que los términos de la sucesión estén numerados a partir de 0 o de cualquier otro número natural va a ser irrelevante, pues sólo nos van a interesar las propiedades de las sucesiones que dependen de sus términos “avanzados”. Si queremos referirnos a la sucesión de números racionales

$$1, \quad \frac{1}{2}, \quad \frac{1}{3}, \quad \frac{1}{4}, \quad \frac{1}{5}, \quad \dots,$$

tenemos dos opciones, o consideramos $\{1/n\}_{n=1}^{\infty}$, que no está definida en $n = 0$, o bien consideramos $\{1/(n+1)\}_{n=0}^{\infty}$, que es un objeto diferente desde un punto de vista conjuntista (es una aplicación con dominio \mathbb{N} en lugar de $\mathbb{N} \setminus \{0\}$), pero todo cuanto digamos va a ser independiente de si la opción que elijamos.

³Esto es válido si entendemos que 10 representa a cualquier número natural k en base k . Cuando 10 no tiene el significado habitual se indica con un subíndice. Por ejemplo,

$$11.011_2 = 1 \cdot 2 + 1 \cdot 1 + 0 \cdot 2^{-1} + 1 \cdot 2^{-2} + 1 \cdot 2^{-3} = \frac{27}{8} = 3.375.$$

Una sucesión $\{x_n\}_{n=0}^{\infty}$ en un conjunto ordenado es *monótona creciente* si cumple

$$x_0 \leq x_1 \leq x_2 \leq x_3 \leq \dots$$

es decir, si cuando $m \leq n$ se cumple $x_m \leq x_n$. Si cambiamos las desigualdades por $<$ se dice que es *estrictamente creciente*. La sucesión es *monótona decreciente* si cumple

$$x_0 \geq x_1 \geq x_2 \geq x_3 \geq \dots$$

o, lo que es lo mismo, si cuando $m \leq n$ entonces $x_n \leq x_m$. De nuevo, la sucesión es estrictamente decreciente si las desigualdades son estrictas. Diremos que una sucesión es *monótona* cuando es monótona creciente o monótona decreciente.

Por ejemplo, la sucesión $\{1/n\}_{n=1}^{\infty}$ es estrictamente decreciente.

Una *subsucesión* de una sucesión dada $\{x_n\}_{n=0}^{\infty}$ es una sucesión de la forma $\{x_{n_k}\}_{k=0}^{\infty}$, donde $\{n_k\}_{k=0}^{\infty}$ es una sucesión estrictamente creciente de números naturales.

Por ejemplo,

$$\frac{1}{2}, \frac{1}{4}, \frac{1}{6}, \frac{1}{8}, \dots$$

es una subsucesión de $\{1/n\}_{n=1}^{\infty}$, concretamente, la determinada por la sucesión de índices $n_k = 2k$.

Notemos que, en general, si $\{n_k\}_{k=0}^{\infty}$ es una sucesión estrictamente creciente de números naturales, siempre se cumple que $k \leq n_k$.

En efecto, si para algún k se cumpliera $n_k < k$, podemos tomar el mínimo posible, y entonces $n_{n_k} < n_k$, con $n_k < k$, lo que contradice la minimalidad de k .

En algunas ocasiones el teorema siguiente puede simplificar las cosas:

Teorema 1.7 *Toda sucesión en un conjunto totalmente ordenado contiene una subsucesión monótona.*

DEMOSTRACIÓN: Sea $\{a_n\}_{n=0}^{\infty}$ una sucesión en un conjunto totalmente ordenado. Sea A el conjunto de las imágenes de la sucesión. Si A es finito es obvio que A tiene una subsucesión constante, luego monótona. Supongamos que A es infinito.

Si todo subconjunto no vacío de A tiene mínimo podemos tomar x_0 igual al mínimo de A , luego x_1 igual al mínimo de $A \setminus \{x_0\}$, luego x_2 igual al mínimo de $A \setminus \{x_0, x_1\}$, y así obtenemos puntos $x_0 < x_1 < x_2 < \dots$, es decir, obtenemos un subconjunto de A sin máximo.

Así pues, o bien existe un subconjunto de A sin mínimo o bien existe un subconjunto de A sin máximo. Los dos casos se tratan igual. Supongamos que hay un subconjunto de A sin mínimo. Llamémoslo B .

Sea a_{n_0} un elemento cualquiera de B . Como B no tiene mínimo contiene infinitos términos de la sucesión bajo a_{n_0} , pero sólo un número finito de ellos tienen índice anterior a n_0 , luego existe un cierto a_{n_1} en B de manera que

$a_{n_1} < a_{n_0}$ y $n_0 < n_1$. Podemos repetir recurrentemente este proceso y obtener una subsucesión

$$a_{n_0} > a_{n_1} > a_{n_2} > a_{n_3} > a_{n_4} > a_{n_5} > \dots$$

monótona decreciente. ■

Si $\{a_n\}_{n=0}^{\infty}$ es una sucesión en un anillo A , podemos considerar la sucesión

$$\sum_{n=0}^{\infty} a_n = \left\{ \sum_{n=0}^N a_n \right\}_{N=0}^{\infty},$$

es decir, la sucesión

$$a_0, \quad a_0 + a_1, \quad a_0 + a_1 + a_2, \quad a_0 + a_1 + a_2 + a_3, \quad \dots$$

Las sucesiones de esta forma se llaman *series infinitas* y los términos que la componen, $S_N = \sum_{n=0}^N a_n$, se llaman *sumas parciales* de la serie.

Por ejemplo, la serie $\sum_{n=0}^{\infty} a^n$ se llama *serie geométrica* de razón a . Es fácil ver que

$$a^n - 1 = (a - 1)(a^{n-1} + \dots + a + 1),$$

de donde las sumas parciales de una serie geométrica de razón $a \neq 1$ son de la forma

$$\sum_{n=0}^N a^n = \frac{1 - a^{N+1}}{1 - a},$$

es decir: $\sum_{n=0}^{\infty} a^n = \left\{ \frac{1 - a^{N+1}}{1 - a} \right\}_{N=0}^{\infty}$.

Ahora podemos decir que el teorema 1.6 asocia a cada elemento $x > 0$ en un cuerpo ordenado arquimediano R una serie finita $\sum_{n=0}^m c_n 10^n$ y otra infinita $\sum_{n=0}^{\infty} c_{-n} 10^{-n}$ de números racionales, donde $0 \leq c_n < 10$, que “en cierto sentido” determinan a x . Pero todavía tenemos que precisar en qué sentido exacto podemos decir que estos desarrollos decimales determinan el x de partida.

La idea es que las sumas a_N están “cada vez más cerca” de x . Más precisamente, como $x \in [a_0, a_0 + 1]$, tenemos que a_0 está a menos de una unidad de x , como $x \in [a_1, a_1 + 1/10]$, tenemos que a_1 está a menos de una décima de x , e igualmente a_2 está a menos de una centésima, etc.

Valores absolutos Al hablar de proximidad entre números estamos usando implícitamente que la distancia entre dos números racionales x e y o, más en general, entre dos elementos de un cuerpo ordenado arquimediano, es $|x - y|$, lo cual tiene una justificación geométrica que ya hemos señalado. Vamos a explicitar este tipo de consideraciones, pero para ello conviene axiomatizar las propiedades del valor absoluto de \mathbb{Q} :

Definición 1.8 Un *valor absoluto* en un cuerpo K es cualquier aplicación $|\cdot| : K \rightarrow R$, donde R es un cuerpo ordenado arquimediano, que cumpla las propiedades siguientes:

- a) $|x| \geq 0$ y el único elemento de K que cumple $|x| = 0$ es $x = 0$.
- b) $|x + y| \leq |x| + |y|$.
- c) $|xy| = |x||y|$.

Un *cuerpo métrico* es un par $(K, |\cdot|)$, donde K es un cuerpo y $|\cdot|$ es un valor absoluto en K .

Sabemos que el valor absoluto en cualquier cuerpo ordenado arquimediano K cumple estas propiedades. (En este caso el valor absoluto toma imágenes en $R = K$, y la propiedad arquimediana de K sólo hace falta para cumplir el requisito impuesto de que R sea arquimediano.)

Ahora podemos probar que algunas propiedades conocidas del valor absoluto en \mathbb{Q} (o en un cuerpo ordenado) son válidas en cualquier cuerpo métrico:

- $|1| = |-1| = 1$,
pues $|1| = |1 \cdot 1| = |1| \cdot |1|$ y, como por a) $|1| \neq 0$, tiene que ser $|1| = 1$. Por otra parte, $|-1|^2 = |(-1)^2| = |1| = 1$, luego $|-1| = \pm 1$, pero los valores absolutos son positivos.
- $|x| = |-x|$,
pues $|-x| = |-1||x| = |x|$.
- Si $x \neq 0$ tiene inverso para el producto, entonces $|x^{-1}| = |x|^{-1}$,
pues $|x||x^{-1}| = |xx^{-1}| = |1| = 1$.
- $||x| - |y|| \leq |x - y|$.
En efecto, $|x| = |x - y + y| \leq |x - y| + |y|$, luego $|x| - |y| \leq |x - y|$. Igualmente probamos que $|y| - |x| \leq |y - x| = |x - y|$, luego

$$-|x - y| \leq |x| - |y| \leq |x - y|,$$

y, por las propiedades del valor absoluto en el cuerpo ordenado R , esto implica que $||x| - |y|| \leq |x - y|$.

Espacios métricos En cada cuerpo métrico K podemos definir la distancia entre dos de sus elementos mediante la fórmula

$$d(x, y) = |x - y|,$$

de modo que en el caso de \mathbb{Q} tenemos la distancia natural desde el punto de vista geométrico. Ahora bien, la mayor parte de los resultados que vamos a obtener sirven igualmente para otras distancias definidas en contextos muy distintos con tal de que cumplan unos pocos axiomas básicos:

Definición 1.9 Una *distancia* en un conjunto M es una aplicación

$$d : M \times M \longrightarrow R,$$

donde R es un cuerpo ordenado arquimediano,⁴ que cumpla las propiedades siguientes:

- a) $d(x, y) \geq 0$ y $d(x, y) = 0$ únicamente cuando $x = y$,
- b) $d(x, y) = d(y, x)$,
- c) $d(x, z) \leq d(x, y) + d(y, z)$.

Es obvio que las tres propiedades son exigencias razonables para cualquier aplicación que realmente pueda ser considerada como medida de una distancia. La tercera recibe el nombre de *desigualdad triangular*. De ellas se deduce una cuarta:

$$|d(x, y) - d(x, z)| \leq d(y, z).$$

En efecto, $d(x, y) \leq d(x, z) + d(z, y)$, luego $d(x, y) - d(x, z) \leq d(y, z)$. Invirtiendo los papeles de y, z obtenemos igualmente que $d(x, z) - d(x, y) \leq d(y, z)$, luego

$$-d(y, z) \leq d(x, y) - d(x, z) \leq d(y, z),$$

y esto es equivalente a la desigualdad que queremos probar.

Un *espacio métrico* es un par (M, d) , donde M es un conjunto y d es una distancia en M .

La conexión con lo anterior es que si K es un cuerpo métrico, entonces la aplicación $d : K \times K \longrightarrow R$ dada por $d(x, y) = |x - y|$ es una distancia en K . En efecto, los axiomas de distancia se demuestran trivialmente a partir de los axiomas de valor absoluto.

En lo sucesivo consideraremos siempre a cada cuerpo métrico como espacio métrico con esta distancia.

Así pues, a partir de aquí podemos razonar a cuatro niveles de generalidad: podemos definir y demostrar cosas sobre \mathbb{Q} en particular, o bien sobre un cuerpo ordenado arquimediano arbitrario, sin necesidad de que sea el propio \mathbb{Q} , o bien sobre un cuerpo métrico arbitrario, sin necesidad de que su valor absoluto provenga de una relación de orden, o bien sobre un espacio métrico arbitrario, sin necesidad de que su distancia provenga de un valor absoluto. Los argumentos serían exactamente los mismos aunque los presentáramos con menos generalidad de la posible.

⁴En la definición usual se toma $R = \mathbb{R}$, cosa que no podemos hacer aquí porque todavía no hemos definido \mathbb{R} . Luego veremos que, en contra de lo que podría parecer, esta definición no es más general.

1.2 Convergencia de sucesiones

Ya podemos introducir uno de los conceptos fundamentales que vamos a manejar a lo largo de este libro:

Definición 1.10 Una sucesión $\{x_n\}_{n=0}^{\infty}$ en un espacio métrico M converge a un límite $l \in M$ si⁵ para todo $\epsilon > 0$ existe un m tal que para todo $n \geq m$ se cumple $d(x_n, l) < \epsilon$.

Es fundamental que el lector asimile plenamente esta definición, sea meditando sobre ella, sea contrastándola con los resultados que exponemos a continuación. Informalmente, que una sucesión converja a un límite l significa que sus términos se aproximan cada vez más a l hasta hacerse prácticamente indistinguibles de l . Por ejemplo, si una sucesión en \mathbb{Q} converge a 1 tenemos que, para $\epsilon = 1/1\,000\,000$, existe un $m \in \mathbb{N}$ tal que todos los términos posteriores a x_m se diferencian de 1 en a lo sumo una millonésima, y si esta aproximación no nos parece suficiente, siempre podemos tomar un ϵ menor, pues aumentando m si es necesario tendremos garantizada la aproximación que deseemos.

Una sucesión no altera su convergencia o no convergencia porque se alteren sus primeros términos (cualquier número finito de ellos). Notemos también que la definición de convergencia se aplica igualmente a sucesiones definidas a partir de $n = 0$ o de $n = 17$.

Observemos que una sucesión no puede converger a más de un límite, pues si $\{x_n\}_{n=0}^{\infty}$ converge a dos puntos l y l' de un espacio métrico, entonces, para $\epsilon = d(l, l')$, existe un $m \in \mathbb{N}$ tal que si $n \geq m$ se cumple⁶ $d(x_n, l) < \epsilon/2$ y $d(x_n, l') < \epsilon/2$, y entonces tenemos una contradicción:

$$d(l, l') \leq d(l, x_n) + d(x_n, l') < \epsilon/2 + \epsilon/2 = d(l, l').$$

Si una sucesión $\{x_n\}_{n=0}^{\infty}$ en un espacio métrico M es convergente, representaremos su único límite como

$$\lim_n x_n \in M.$$

Nota En el caso de series infinitas, es decir, de sucesiones de la forma $\sum_{n=0}^{\infty} a_n$ en un anillo (dotado de una métrica), en caso de que tengan límite éste se representa igualmente por $\sum_{n=0}^{\infty} a_n$ y, en lugar de límite, se le llama *suma* de la serie. El contexto dejará siempre claro si nos estamos refiriendo al límite o a la sucesión de sumas parciales de la serie, y ello no debería causar confusión a un lector que sea consciente de este doble significado posible. ■

⁵En lo sucesivo sobrentenderemos que la letra ϵ representa siempre un elemento del cuerpo R donde toma imágenes la distancia, mientras que m, n representarán siempre números naturales.

⁶En principio tendríamos un m_1 tal que para $n \geq m_1$ se cumpliría la definición de convergencia a l y otro m_2 a partir del cual se cumpliría la definición de convergencia a l' , pero si tomamos $m = \max\{m_1, m_2\}$, a partir de m se cumplen las dos. En lo sucesivo aplicaremos tácitamente este argumento con frecuencia.

Ejemplos Veamos algunos casos sencillos de convergencia:

- Una sucesión constante $\{c\}_{n=0}^{\infty}$ en un espacio métrico converge a c .
En efecto, cumple trivialmente la definición porque en este caso tenemos que $d(x_n, l) = d(c, c) = 0 < \epsilon$.
- Más en general, lo mismo vale para las sucesiones finalmente constantes.
- Toda subsucesión de una sucesión convergente converge al mismo límite.
En efecto, si se cumple $d(x_n, l) < \epsilon$ para todo $n \geq m$, también vale $d(x_{n_k}, l) < \epsilon$ para $k \geq m$, ya que entonces $n_k \geq k \geq m$.
- La sucesión $\{(-1)^n\}_{n=0}^{\infty}$ en \mathbb{Q} no es convergente.

Se trata de la sucesión $1, -1, 1, -1, \dots$ y no converge a ningún límite l , porque si así fuera, tomando $\epsilon = 1$, tendríamos un natural m tal que si $n \geq m$ se cumpliría $|x_n - l| < 1$. En particular, tomando un n par y otro impar, $|1 - l| < 1$ y $|-1 - l| < 1$, luego

$$2 = |1 - (-1)| = |1 - l + l - (-1)| \leq |1 - l| + |l - (-1)| < 1 + 1 = 2,$$

contradicción.

- Si K es un cuerpo ordenado arquimediano,⁷ entonces $\lim_n \frac{1}{n} = 0$.
En efecto, dado $\epsilon > 0$, la propiedad arquimediana nos da un $m \in K$ tal que $1/\epsilon < m$, luego si $n \geq m$, se cumple también que $1/\epsilon < n$, luego $|1/n| = 1/n < \epsilon$, como requiere la definición de límite.
- Si K es un cuerpo métrico y $x \in K$ cumple $|x| < 1$, entonces $\lim_n x^n = 0$.
En efecto, si $x = 0$ la sucesión es constante igual a 0 y la conclusión es trivial. En otro caso sea $\alpha = 1/|x| > 1$. Entonces

$$\alpha^n = (1 + \alpha - 1)^n = \sum_{i=0}^n \binom{n}{i} (\alpha - 1)^i \geq 1 + n(\alpha - 1),$$

porque todos los términos de la suma son positivos. Dado $\epsilon > 0$, existe un natural m tal que

$$\frac{1/\epsilon - 1}{\alpha - 1} < m,$$

y así si $n \geq m$ se cumple $\alpha^n \geq 1 + n(\alpha - 1) > 1/\epsilon$, luego $|x^n - 0| = |x|^n < \epsilon$, como exige la definición de límite. ■

Por ejemplo, ahora ya podemos probar que todo elemento de un cuerpo ordenado está determinado por su desarrollo decimal:

⁷La propiedad arquimediana es esencial para este ejemplo. La sucesión $1/n$ no tiene por qué converger en todo cuerpo métrico. De hecho, podemos definir un cuerpo métrico *arquimediano* como un cuerpo métrico en el que la sucesión $1/n$ tiende a 0. Entonces este ejemplo prueba que los cuerpos métricos asociados a cuerpos ordenados arquimedianos son arquimedianos en este sentido.

Teorema 1.11 Si $x = x_0 > 0$ es un elemento de un cuerpo ordenado arquimediano R , consideramos el desarrollo en base 10 de $E[x] = \sum_{n=0}^m c_n 10^n$ y, para cada $n \geq 1$, definimos $x_n = F[10^{n-1}x_{n-1}]/10^{n-1}$, $c_{-n} = E[10^n x_n]$, Entonces $0 \leq c_n < 10$ y

$$x = \sum_{n=0}^m c_n 10^n + \sum_{n=1}^{\infty} c_{-n} 10^{-n}.$$

DEMOSTRACIÓN: Sea $y = x - \sum_{n=0}^m c_n 10^n$ y sea $a'_N = \sum_{n=1}^N c_{-n} 10^{-n}$. El teorema 1.6 afirma que $a_N \leq x \leq a_N + 10^{-N}$, de donde $a'_N \leq y \leq a'_N + (1/10)^N$. Por lo tanto $|y - a'_N| \leq (1/10)^N$.

Como $|1/10| < 1$, sabemos que $\lim_N (1/10)^N = 0$ luego, dado $\epsilon > 0$, existe un $N_0 \in \mathbb{N}$ tal que si $N \geq N_0$ entonces

$$|y - a'_N| \leq (1/10)^N < \epsilon,$$

y esto significa que $y = \lim_N a'_N = \sum_{n=1}^{\infty} c_{-n} 10^{-n}$, luego x se expresa como indica el enunciado. ■

Así pues, todo elemento positivo de un cuerpo ordenado arquimediano (piénsese “todo número real positivo”) puede expresarse como suma infinita de números racionales, concretamente a través de su desarrollo decimal. Ahora ya podemos afirmar con precisión que, si en un cuerpo ordenado arquimediano existe una raíz cuadrada de 2 positiva, tiene que ser

$$\sqrt{2} = 1.41421356237310\dots$$

donde el miembro de la derecha representa el límite de una serie infinita en las condiciones del teorema anterior y los puntos suspensivos abrevian una cierta definición recurrente de la sucesión de cifras decimales, que puede calcularse estrictamente en términos de la aritmética de \mathbb{Q} .

Para extraer más consecuencias necesitamos algunos resultados más sobre cálculo de límites. Vamos a necesitar el concepto siguiente:

Definición 1.12 Una sucesión $\{x_n\}_{n=0}^{\infty}$ en un espacio métrico M está *acotada* si existe un $x \in M$ y un $C \in \mathcal{R}$ de modo que para todo $n \in \mathbb{N}$ se cumple $d(x_n, x) \leq C$.

La idea es muy simple: una sucesión está acotada si nunca se aleja más de una distancia C de un punto fijo x . Conviene observar que el punto x podemos elegirlo, en el sentido de que si una sucesión cumple la definición de sucesión acotada con un cierto $x \in M$, entonces la cumple también con cualquier otro $y \in M$, pues

$$d(x_n, y) \leq d(x_n, x) + d(x, y) \leq C + d(x, y),$$

y basta tomar $C' = C + d(x, y)$.

En particular, una sucesión $\{x_n\}_{n=0}^{\infty}$ está acotada en un cuerpo métrico si y sólo si existe un $C \in R$ tal que para todo $n \in \mathbb{N}$ se cumple $|x_n| \leq C$ (donde hemos tomado $x = 0$ como punto de referencia para fijar la cota).

Teorema 1.13 *Toda sucesión convergente en un espacio métrico está acotada.*

DEMOSTRACIÓN: Sea $\{x_n\}_{n=0}^{\infty}$ una sucesión en un espacio métrico M convergente a un límite l . Entonces, tomando $\epsilon = 1$, tenemos que existe un m tal que si $n \geq m$, entonces $d(x_n, l) < 1$. Sea

$$C = \max\{1, d(x_0, l), \dots, d(x_{m-1}, l)\}.$$

Así es claro que $d(x_n, l) \leq C$ para todo natural n , luego la sucesión está acotada. ■

Veamos una aplicación:

Teorema 1.14 *Sea K un cuerpo métrico y $\{a_n\}_{n=0}^{\infty}$, $\{b_n\}_{n=0}^{\infty}$, dos sucesiones convergentes en K . Entonces las sucesiones $\{a_n + b_n\}_{n=0}^{\infty}$ y $\{a_n b_n\}_{n=0}^{\infty}$ también son convergentes y*

$$\lim_n (a_n + b_n) = \lim_n a_n + \lim_n b_n, \quad \lim_n (a_n b_n) = \lim_n a_n \lim_n b_n.$$

DEMOSTRACIÓN: Llamamos $l = \lim_n a_n$, $l' = \lim_n b_n$. Para la suma observamos que

$$|a_n + b_n - l - l'| \leq |a_n - l| + |b_n - l'|,$$

luego, dado $\epsilon > 0$, tomando un $m \in \mathbb{N}$ tal que $|a_n - l| < \epsilon/2$ y $|b_n - l'| < \epsilon/2$ siempre que $n \geq m$, tenemos también que $|a_n + b_n - l - l'| < \epsilon$, luego la suma converge a $l + l'$.

Para el producto usamos que

$$|a_n b_n - ll'| = |a_n b_n - a_n l' + a_n l' - ll'| \leq |a_n| |b_n - l'| + |a_n - l| |l'|,$$

así como que ambas sucesiones están acotadas, de modo que existe un $C \in R$ tal que $|a_n| \leq C$, $|b_n| \leq C$ para todo n . Podemos suponer además que $|l'| \leq C$. Así, dado $\epsilon > 0$, tomamos m tal que si $n \geq m$ se cumpla $|a_n - l| < \epsilon/2C$, $|b_n - l'| < \epsilon/2C$, y entonces

$$|a_n b_n - ll'| < C \frac{\epsilon}{2C} + \frac{\epsilon}{2C} C = \epsilon,$$

luego el producto converge a ll' . ■

Un caso particular es que si $\{a_n\}_{n=0}^{\infty}$ converge a l y $c \in K$, entonces

$$\lim_n ca_n = c \lim_n a_n.$$

Basta aplicar el teorema anterior a las sucesiones $\{c\}_{n=0}^{\infty}$ y $\{a_n\}_{n=0}^{\infty}$.

Otro caso particular es que podemos partir las series infinitas:

$$\sum_{n=0}^{\infty} a_n = \sum_{n=0}^k a_n + \sum_{n=k+1}^{\infty} a_n,$$

donde hay que entender que la serie de la izquierda converge si y sólo si lo hace la de la derecha, y en tal caso se da la igualdad indicada entre las sumas.

En efecto, si llamamos $S_N = \sum_{n=0}^N a_n$, teniendo en cuenta que una sucesión converge si y sólo si lo hace la sucesión que resulta de suprimir sus primeros términos (y en tal caso converge al mismo límite), basta probar que $\{S_N\}_{N=k+1}^{\infty}$ converge si y sólo si lo hace $\{S_N - S_k\}_{N=k+1}^{\infty}$, y que en tal caso

$$\lim_N S_N = \lim_N S_k + \lim_N (S_N - S_k),$$

pero una implicación es justo lo que afirma el teorema anterior (notemos que la sucesión central es constante), y la otra también si reformulamos la igualdad como

$$\lim_N (S_N - S_k) = \lim_N (-S_k) + \lim_N S_N.$$

Si una sucesión converge a 0, para que converja un producto basta con que la otra sucesión esté acotada:

Teorema 1.15 *Si $\{x_n\}_{n=0}^{\infty}$ es una sucesión acotada en un cuerpo métrico y $\{y_n\}_{n=0}^{\infty}$ converge a 0, entonces $\{x_n y_n\}_{n=0}^{\infty}$ converge a 0.*

DEMOSTRACIÓN: Sea $C \in R$ tal que $|x_n| \leq C$ para todo n . Dado $\epsilon > 0$, existe un $m \in \mathbb{N}$ tal que si $n \geq m$, entonces $|y_n| < \epsilon/C$, luego

$$|x_n y_n| = |x_n| |y_n| < C\epsilon/C = \epsilon.$$

Esto significa que la sucesión producto tiende a 0. ■

Teorema 1.16 *Sea R un cuerpo ordenado arquimediano y sea $\{x_n\}_{n=0}^{\infty}$ una sucesión convergente en R tal que existe $a \in R$ y $m \in \mathbb{N}$ de modo que $a \leq x_n$ (resp. $x_n \leq a$) para todo $n \geq m$. Entonces $a \leq \lim_n x_n$ (resp. $\lim_n x_n \leq a$).*

DEMOSTRACIÓN: Probaremos el caso en que $a \leq x_n$. El otro es análogo. Supongamos que el límite cumple $l < a$. Entonces tomamos $\epsilon = a - l$ y, por la definición de convergencia, existe un $n \geq m$ tal que $|x_n - l| < \epsilon$ y, como $l < a \leq x_n$, esto equivale a $x_n - l < a - l$, con lo que $x_n < a$, contradicción. ■

Como consecuencia:

Teorema 1.17 *Sea R un cuerpo ordenado arquimediano y $\{a_n\}_{n=0}^{\infty}$, $\{b_n\}_{n=0}^{\infty}$ dos sucesiones convergentes en R tales que $a_n \leq b_n$ para todo n . Entonces se cumple que $\lim_n a_n \leq \lim_n b_n$.*

DEMOSTRACIÓN: Si llamamos $c_n = b_n - a_n \geq 0$, el teorema anterior nos da que

$$\lim_n b_n - \lim_n a_n = \lim_n c_n \geq 0. \quad \blacksquare$$

Teorema 1.18 (Criterio del emparedado) *Sea R un cuerpo ordenado arquimediano y $\{a_n\}_{n=0}^\infty$, $\{b_n\}_{n=0}^\infty$, $\{c_n\}_{n=0}^\infty$ tres sucesiones en R que cumplan $a_n \leq b_n \leq c_n$ para todo n . Si la primera y la tercera convergen al mismo límite l , entonces también $\lim_n b_n = l$.*

DEMOSTRACIÓN: Dado $\epsilon > 0$, existe un m tal que si $n \geq m$ entonces $|a_n - l| < \epsilon$ y $|c_n - l| < \epsilon$. Entonces

$$-\epsilon < a_n - l \leq b_n - l \leq c_n - l < \epsilon,$$

luego $|b_n - l| < \epsilon$, como requiere la definición de límite. \blacksquare

Con esto ya estamos en condiciones de estudiar los desarrollos decimales. Dedicamos a ello la sección siguiente, pero terminamos ésta introduciendo un concepto que permite explicitar una idea que subyace en todas las consideraciones que estamos exponiendo:

Definición 1.19 Si M es un espacio métrico, un conjunto $D \subset M$ es *denso* en M si para todo $x \in M$ y todo $\epsilon > 0$ existe un $d \in D$ tal que $d(x, d) < \epsilon$.

Así, que D sea denso en M quiere decir que todo punto de M , aunque pueda no estar en D , tiene puntos de D tan cercanos como se desee. Podemos caracterizar la densidad en términos de sucesiones:

Teorema 1.20 *Un subconjunto D de un espacio métrico M es denso si y sólo si todo punto de M es el límite de una sucesión en D .*

DEMOSTRACIÓN: Si D es denso y $l \in M$, para cada $n \in \mathbb{N}$ tomamos $d_n \in D$ tal que $d(d_n, l) < \frac{1}{n+1}$. Tenemos así una sucesión $\{d_n\}_{n=0}^\infty$ en D que ciertamente converge a l , pues, dado $\epsilon > 0$, como R es arquimediano existe un $m \in \mathbb{N}$ tal que $1/\epsilon < m$, luego $1/m < \epsilon$, y si $n \geq m$ entonces

$$d(d_n, l) < \frac{1}{n+1} < \frac{1}{m} < \epsilon.$$

Recíprocamente, si D tiene la propiedad indicada, dado $x \in M$, tomamos una sucesión $\{d_n\}_{n=0}^\infty$ en D que converja a x . Dado $\epsilon > 0$ existe un $m \in \mathbb{N}$ tal que $d(d_m, x) < \epsilon$, y esto prueba que D es denso. \blacksquare

En los cuerpos ordenados los conjuntos densos tienen una caracterización más simple:

Teorema 1.21 *Un subconjunto D de un cuerpo ordenado arquimediano R es denso si y sólo si para todo par $a < b$ de elementos de R existe un $d \in D$ tal que $a < d < b$.*

DEMOSTRACIÓN: Si D es denso y $a < b$, consideramos $c = (a + b)/2$ y $\epsilon = (b - a)/2 > 0$. Existe un $d \in D$ tal que $|d - c| < \epsilon$, pero esto quiere decir que $-\epsilon < d - c < \epsilon$, o también $a - c < d - c < b - c$, luego $a < d < b$.

Recíprocamente, si D cumple la condición del enunciado, dado $x \in R$ y $\epsilon > 0$, existe un $d \in D$ tal que $x < d < x + \epsilon$, y entonces $|d - x| = d - x < \epsilon$. ■

Así pues, en estos términos el teorema 1.5 afirma que \mathbb{Q} es denso en todo cuerpo ordenado arquimediano. Ya habíamos comprobado esto indirectamente, pues sabemos que todo elemento de un cuerpo ordenado arquimediano es el límite de una sucesión en \mathbb{Q} (su desarrollo decimal).

1.3 Desarrollos decimales

Resumamos lo que hemos obtenido hasta el momento sobre desarrollos decimales:

Definición 1.22 Dadas una sucesión finita $\{c_n\}_{n=0}^m$ y una sucesión infinita $\{c_{-n}\}_{n=1}^{\infty}$ de números naturales $0 \leq c_n < 10$, llamamos *desarrollo decimal* asociado a dichas sucesiones a la sucesión siguiente⁸ de números racionales:

$$\sum_{n=0}^m c_n 10^n + \sum_{n=1}^{\infty} c_{-n} 10^{-n},$$

es decir, la sucesión cuyo término N -simo es

$$\sum_{n=0}^m c_n 10^n + \sum_{n=1}^N c_{-n} 10^{-n}.$$

Una sucesión de este tipo puede converger o no en \mathbb{Q} y, aunque no converja en \mathbb{Q} , puede converger en otro cuerpo ordenado arquimediano R (que necesariamente contiene a \mathbb{Q} como subcuerpo denso).

Nuestro propósito a medio plazo es demostrar que todos los desarrollos decimales convergen en \mathbb{R} . De momento, lo que tenemos probado (teorema 1.11) es que si R es cualquier cuerpo ordenado arquimediano y $x = x_0 \in R$, entonces existe un desarrollo decimal convergente a x , concretamente, el determinado por

$$E[x] = \sum_{n=0}^m c_n 10^n, \quad x_n = F[10^{n-1}x_{n-1}]/10^{n-1}, \quad c_{-n} = E[10^n x_n].$$

Cuando un desarrollo decimal converge, usamos la misma notación para representar su límite, es decir, escribimos

$$x = \sum_{n=0}^m c_n 10^n + \sum_{n=1}^{\infty} c_{-n} 10^{-n},$$

donde ahora el miembro derecho no debe entenderse como la sucesión que hemos descrito antes, sino como su límite.

⁸Insistimos en que 10 puede interpretarse como cualquier número natural $k \geq 2$ escrito en base k , aunque en todos los ejemplos concretos entenderemos que es el 10 usual.

Un caso particular de desarrollos decimales que convergen trivialmente son aquellos cuya sucesión de cifras es finalmente igual a 0, pues entonces el desarrollo decimal (como sucesión) es finalmente constante. A dichos desarrollos decimales los llamamos *exactos*, y hemos introducido la notación

$$c_m \cdots c_0 . c_{-1} \cdots c_{-N} = \sum_{n=0}^m c_n 10^n + \sum_{n=1}^N c_{-n} 10^{-n} \in \mathbb{Q}.$$

Por ejemplo, $395.178 = 3 \cdot 100 + 9 \cdot 10 + 5 + 1/10 + 7/100 + 8/1000$.

Es costumbre usar una notación análoga para desarrollos infinitos con puntos suspensivos. Así, por definición,

$$0.3333 \dots = \sum_{n=1}^{\infty} 3 \cdot 10^{-n},$$

mientras que $5.101001000100001 \dots$ representa el desarrollo decimal (que en principio puede ser convergente o no, según el cuerpo ordenado considerado) determinado por la sucesión de cifras dada por $c_0 = 5$ y

$$c_{-n} = \begin{cases} 1 & \text{si } n = \frac{k(k+1)}{2} \text{ para cierto } k \geq 1, \\ 0 & \text{en otro caso.} \end{cases}$$

En general, el uso de los puntos suspensivos sólo estará justificado en la medida en que quede claro por el contexto cuál es el criterio que determina las cifras decimales siguientes.

El primer paso para estudiar los desarrollos decimales es el teorema siguiente sobre suma de series geométricas:

Teorema 1.23 *Sea K un cuerpo métrico y $x \in K$ tal que $|x| < 1$. Entonces*

$$\sum_{n=0}^{\infty} x^n = \frac{1}{1-x}.$$

DEMOSTRACIÓN: Hemos visto que $\sum_{n=0}^{\infty} x^n = \left\{ \frac{1-x^{N+1}}{1-x} \right\}_{N=0}^{\infty}$, hemos probado que $\lim_N x^N = 0$, luego $\lim_N -x^N = 0$, luego $\lim_N 1 - x^N = 1$, luego

$$\lim_N \frac{1}{1-x} (1 - x^{N+1}) = \frac{1}{1-x}. \quad \blacksquare$$

Un poco más en general, si $|x| < 1$:

$$\sum_{n=k}^{\infty} x^n = x^k \sum_{n=0}^{\infty} x^n = \frac{x^k}{1-x}.$$

La primera igualdad hay que entenderla como una igualdad de sucesiones, y la segunda nos da el límite.

Por ejemplo, tomando $x = 1/10$ obtenemos que en \mathbb{Q} se cumple

$$\frac{1}{9} = \sum_{n=1}^{\infty} 10^{-n} = 0.1111\dots$$

Multiplicando por 3, vemos que

$$\frac{1}{3} = \sum_{n=1}^{\infty} 3 \cdot 10^{-n} = 0.3333\dots$$

Más sorprendente es lo que sucede al multiplicar por 9:

$$0.9999\dots = \sum_{n=1}^{\infty} 9 \cdot 10^{-n} = 1 = 1.0000\dots$$

Vemos así que un mismo número racional puede tener dos desarrollos decimales distintos. El fenómeno es un poco más general. Por ejemplo,

$$\begin{aligned} 4.739999\dots &= 4 + 7 \cdot 10^{-1} + 3 \cdot 10^{-2} + \sum_{n=3}^{\infty} 9 \cdot 10^{-n} = \\ &= 4.73 + 9 \frac{10^{-3}}{1 - 10^{-1}} = 4.73 + 10^{-2} = 4.74. \end{aligned}$$

Así, si una sucesión $\{a_n\}_{n=1}^{\infty}$ es finalmente igual a 9 (digamos a partir de a_n), entonces

$$0.a_1 a_2 \dots a_{n-1} 999\dots = 0.a_1 a_2 \dots a_{n-2} (a_{n-1} + 1),$$

pues

$$\begin{aligned} 0.a_1 a_2 \dots a_{n-1} 999\dots &= 0.a_1 a_2 \dots a_{n-1} + \sum_{r=n}^{\infty} 9 \cdot 10^{-r} \\ &= 0.a_1 a_2 \dots a_{n-1} + \frac{1}{10^{n-1}} = 0.a_1 a_2 \dots a_{n-2} (a_{n-1} + 1). \end{aligned}$$

O sea, toda expresión decimal finalmente igual a 9 admite una expresión decimal exacta sumando 1 a la cifra anterior al primer 9 de la cola constante.

Vamos a probar que éste es el único caso en el que dos desarrollos decimales pueden coincidir.

Teorema 1.24 *Si R es un cuerpo ordenado arquimediano, todo $x \in R$ positivo admite un único desarrollo decimal*

$$x = \sum_{n=0}^m c_n 10^n + \sum_{n=1}^{\infty} c_{-n} 10^{-n}$$

de modo que $c_m \neq 0$ (o bien $m = 0$) y la sucesión $\{c_{-n}\}_{n=1}^{\infty}$ no sea finalmente constante igual⁹ a 9.

⁹El resultado vale para desarrollos decimales en cualquier base si entendemos que 9 representa el número anterior a la base.

DEMOSTRACIÓN: Sabemos que x admite un desarrollo de este tipo (si la sucesión $\{c_{-n}\}_{n=1}^{\infty}$ fuera finalmente constante igual a 9 hemos visto cómo transformar el desarrollo decimal en otro que no cumple eso). Observemos ahora que en un desarrollo de este tipo (en particular convergente) se cumple que

$$0 \leq \sum_{n=1}^{\infty} c_{-n} 10^{-n} < 1.$$

En efecto, la primera desigualdad se sigue de 1.17 (comparando con la sucesión constante igual a 0), y para la segunda consideramos un n_0 tal que $c_{-n_0} < 9$. Entonces

$$\begin{aligned} \sum_{n=1}^{\infty} c_{-n} 10^{-n} &= \sum_{n=1}^{n_0} c_{-n} 10^{-n} + \sum_{n=n_0+1}^{\infty} c_{-n} 10^{-n} < \\ \sum_{n=1}^{n_0} 9 \cdot 10^{-n} + \sum_{n=n_0+1}^{\infty} 9 \cdot 10^{-n} &= \sum_{n=1}^{\infty} 9 \cdot 10^{-n} = 1. \end{aligned}$$

Hemos tenido que partir la serie para obtener una desigualdad estricta a partir de la primera parte (que es una suma finita), ya que el teorema 1.17 no es cierto para desigualdades estrictas.

Por lo tanto, $\sum_{n=0}^m c_{-n} 10^{-n} = F[x]$. Supongamos que

$$x = \sum_{n=0}^m c_n 10^n + \sum_{n=1}^{\infty} c_{-n} 10^{-n} = \sum_{n=0}^m c'_n 10^n + \sum_{n=1}^{\infty} c'_{-n} 10^{-n},$$

en las condiciones del enunciado. Entonces, por la unicidad de la parte entera y la parte fraccionaria tenemos que

$$\sum_{n=0}^m c_n 10^n = \sum_{n=0}^m c'_n 10^n, \quad \sum_{n=1}^{\infty} c_{-n} 10^{-n} = \sum_{n=1}^{\infty} c'_{-n} 10^{-n},$$

y de la primera igualdad se sigue que $c_n = c'_n$ para $n \geq 0$, por la unicidad de los desarrollos decimales de los números naturales. Supongamos que existe un n_0 tal que $c_{-n_0} \neq c'_{-n_0}$. Podemos tomar el mínimo posible. Entonces

$$\begin{aligned} \sum_{n=1}^{n_0-1} c_{-n} 10^{-n} + c_{-n_0} 10^{-n_0} + \sum_{n=n_0+1}^{\infty} c_{-n} 10^{-n} &= \\ \sum_{n=1}^{n_0-1} c_{-n} 10^{-n} + c'_{-n_0} 10^{-n_0} + \sum_{n=n_0+1}^{\infty} c'_{-n} 10^{-n}, \end{aligned}$$

luego

$$c_{-n_0} 10^{-n_0} \leq c_{-n_0} 10^{-n_0} + \sum_{n=n_0+1}^{\infty} c_{-n} 10^{-n} = c'_{-n_0} 10^{-n_0} + \sum_{n=n_0+1}^{\infty} c'_{-n} 10^{-n} <$$

$$c'_{-n_0} 10^{-n_0} + \sum_{n=n_0+1}^{\infty} 9 \cdot 10^{-n} = c'_{-n_0} 10^{-n_0} + 10^{-n_0} = (c'_{-n_0} + 1) 10^{-n_0}.$$

Por lo tanto $c_{-n_0} < c'_{-n_0} + 1$, luego $c_{-n_0} \leq c'_{-n_0}$, e igualmente probamos la desigualdad opuesta, con lo que $c_{-n_0} = c'_{-n_0}$, contradicción.

Notemos que la desigualdad estricta en la última cadena de expresiones se justifica partiendo la serie como en la primera parte de la prueba, usando que algún $c'_{-n} \neq 9$. ■

En particular, esto vale para \mathbb{Q} , es decir, todo número racional admite un desarrollo decimal, que en la práctica se puede calcular mediante una simple extensión del algoritmo de Euclides para dividir números naturales:

$$\begin{array}{r}
 5.00000000 \mid 14 \\
 \underline{42} \\
 80 \\
 \underline{70} \\
 100 \\
 \underline{98} \\
 20 \\
 \underline{14} \\
 60 \\
 \underline{56} \\
 40 \\
 \underline{28} \\
 120 \\
 \underline{112} \\
 80 \\
 \underline{70} \\
 10
 \end{array}$$

Concluimos que $5/14 = 0.3571428571428571428\dots$. Notemos que, como los restos posibles son un número finito, tras un número finito de pasos hemos de obtener un resto ya obtenido previamente, y como cada cifra del cociente depende exclusivamente del último resto, resulta que las cifras se repiten cíclicamente. En el caso de $5/14$ el grupo de cifras que se repite es 571428 .

Para indicar esto se suele usar la notación $5/14 = 0.3\overline{571428}$.

El bloque 428 se suele llamar *período* del número. El bloque de cifras decimales previas al período (que puede no existir) se llama *anteperíodo* (3 en este caso), luego en la expresión decimal de un número racional podemos distinguir la parte entera, el anteperíodo y el período.

Recíprocamente, todo número cuya expresión decimal sea de esta forma es un número racional. Veámoslo con un ejemplo, aunque el método vale claramente en general.¹⁰ $r = 37.195\overline{513}$. Vamos a encontrar una fracción igual a r . Para ello multiplicamos por un 1 seguido de tantos ceros como cifras hay en el período más el anteperíodo:

$$(10000)r = 3719513.\overline{513}.$$

¹⁰Notemos que el razonamiento presupone que todo desarrollo decimal finalmente periódico es convergente. No merece la pena demostrarlo ahora, porque más adelante veremos que todos los desarrollos decimales convergen en \mathbb{R} , y el argumento que presentamos aquí prueba que, de hecho, los finalmente periódicos convergen a números racionales.

Le restamos r multiplicado por un 1 seguido de tantos ceros como cifras hay en el anteperiodo:

$$(1\,000\,000)r - (1\,000)r = 3\,719\,513 + 0.\overline{513} - 3\,719 - 0.\overline{513} = 3\,715\,794.$$

Por consiguiente

$$r = \frac{3\,715\,794}{999\,000}.$$

■

Así, pues, los desarrollos decimales que convergen en \mathbb{Q} son precisamente los *finalmente periódicos*. Cualquier otro desarrollo decimal no converge en \mathbb{Q} y está señalando un “agujero” en \mathbb{Q} que “está pidiendo” ser rellenado. Para precisar el modo en que los desarrollos decimales señalan huecos conviene introducir un nuevo concepto:

Definición 1.25 Sea R un cuerpo ordenado arquimediano. Una *sucesión de intervalos encajados* en R es una sucesión $\{[a_n, b_n]\}_{n=0}^{\infty}$ de intervalos cerrados en R tales que $a_n \leq a_{n+1} < b_{n+1} \leq b_n$ para todo n y $\lim_n (b_n - a_n) = 0$.

Teorema 1.26 Sea R un cuerpo ordenado arquimediano y $\{[a_n, b_n]\}_{n=0}^{\infty}$ una sucesión de intervalos encajados en K . Si existe un $x \in \bigcap_{n=0}^{\infty} [a_n, b_n]$, entonces sólo hay un x que cumpla esto, y además $x = \lim_n a_n = \lim_n b_n$.

DEMOSTRACIÓN: Estamos suponiendo que $a_n \leq x \leq b_n$ para todo n . Entonces

$$0 \leq x - a_n \leq b_n - a_n, \quad 0 \leq b_n - x \leq b_n - a_n$$

y basta aplicar dos veces el criterio del emparedado para concluir que

$$\lim_n (x - a_n) = 0, \quad \lim_n (b_n - x) = 0,$$

lo cual equivale a que $x = \lim_n a_n = \lim_n b_n$, y esto prueba que x es único, pues una misma sucesión no puede converger a dos límites distintos. ■

Ahora, si R es un cuerpo ordenado arquimediano y $x \in R$ es positivo, el teorema 1.6 nos da que si

$$x = \sum_{n=0}^m c_n 10^n + \sum_{n=1}^{\infty} c_{-n} 10^{-n}$$

y llamamos

$$a_N = \sum_{n=0}^m c_n 10^n + \sum_{n=1}^N c_{-n} 10^{-n}, \quad b_N = a_N + 10^{-N},$$

se cumple que $a_N \leq x \leq b_N$, con lo que $\{[a_N, b_N]\}_{N=0}^{\infty}$ es una sucesión de intervalos encajados en K (aquí usamos que $\lim_N 10^{-N} = 0$) tal que

$$\bigcap_{N=0}^{\infty} [a_N, b_N] = \{x\},$$

considerando los intervalos en \mathbb{R} , pero si consideramos los mismos intervalos en \mathbb{Q} y $x \notin \mathbb{Q}$, entonces

$$\bigcap_{N=0}^{\infty} [a_N, b_N] = \emptyset.$$

Vemos así que un desarrollo decimal, no sólo “se acerca” a un hueco, sino que determina intervalos encajados que “encierran” el “hueco” en un margen cada vez más estrecho, de modo que determinan totalmente la posición que debería tener un número que “rellene” ese hueco.

1.4 Sucesiones de Cauchy

Finalmente estamos en condiciones de presentar la propiedad que debe tener un cuerpo ordenado arquimediano R para que podamos decir que es \mathbb{R} . Entendamos esto: en la sección 2.2 de [G] hemos construido un cuerpo ordenado arquimediano al que hemos llamado \mathbb{R} , pero es posible dar otras construcciones, formalmente muy distintas, de otros cuerpos, muy distintos en términos conjuntistas (en el apéndice A tenemos una) que igualmente pueden tomarse como definición de \mathbb{R} , porque tienen la misma propiedad de completitud que vamos a definir aquí y, según demostraremos, dos cuerpos ordenados arquimedianos que tengan dicha propiedad son necesariamente isomorfos, luego son a efectos prácticos “el mismo cuerpo”, aunque tengan elementos de naturaleza muy distinta (por ejemplo, los números reales construidos en [G] son subconjuntos de \mathbb{Q} , mientras que los construidos en el Apéndice A son conjuntos de sucesiones de números racionales).

Consideremos la sucesión de números racionales

$$1, \quad 1.4, \quad 1.41, \quad 1.414, \quad 1.4142, \quad 1.41421, \quad \dots$$

que “debería” converger a $\sqrt{2}$, pero no converge a nada en \mathbb{Q} . Esta sucesión de números racionales debe ser convergente en cualquier cuerpo que aspire a ser tenido como “un ejemplo” de \mathbb{R} . Por el contrario, sucesiones como $1, 2, 3, 4, \dots$ o como $1, 2, 3, 1, 2, 3, \dots$, no se están acercando a nada en \mathbb{Q} , donde “a nada” incluye que ni siquiera se están acercando a ningún “hueco”, sino que una se hace cada vez más grande y la otra “divaga” sin dirigirse a ninguna parte en concreto.

Tenemos así tres ejemplos de sucesiones de números racionales que no convergen, pero la primera “debe” converger en \mathbb{R} y las otras dos no deben converger ni en \mathbb{R} ni en ninguna parte. La propiedad que caracteriza a \mathbb{R} afirma que todas las sucesiones que deben converger convergen, y lo que vamos a hacer ahora es precisar cómo hay que entender eso de que una sucesión que no converge en \mathbb{Q} deba, pese a ello, converger en \mathbb{R} . Esto nos lleva al concepto de sucesión de Cauchy:

Definición 1.27 Una sucesión $\{x_n\}_{n=0}^{\infty}$ en un espacio métrico M es de Cauchy si para todo $\epsilon > 0$ existe un $m \in \mathbb{N}$ tal que si $n, n' \geq m$ entonces $d(x_n, x_{n'}) < \epsilon$.

Una sucesión es de Cauchy cuando sus términos se aproximan entre sí, hasta hacerse indistinguibles unos de otros, de modo que, superado nuestro umbral de discernimiento, todos ellos se ven “en el mismo punto”. Si la sucesión es convergente, ese punto donde “se aglomeran” los puntos de la sucesión es el límite, pero la definición de sucesión de Cauchy no exige que exista tal límite. Si una sucesión de Cauchy no tiene límite, entonces sus términos se están “aglomerando” alrededor de “nada”, alrededor de “un hueco”, de “un agujero microscópico” en el espacio métrico considerado. (Pero aglomerarse alrededor de nada es más que “divagar”, que es lo que hace la sucesión $1, 2, 3, 1, 2, 3, \dots$, y ésta es la diferencia que marca el concepto de sucesión de Cauchy.)

Teorema 1.28 *Toda sucesión convergente en un espacio métrico es de Cauchy, y toda sucesión de Cauchy está acotada.*

DEMOSTRACIÓN: Si $\{x_n\}_{n=0}^{\infty}$ converge a l , dado $\epsilon > 0$ existe un $m \in \mathbb{N}$ tal que si $n \geq m$ entonces $d(x_n, l) < \epsilon/2$, luego si $n, n' \geq m$ tenemos que

$$d(x_n, x_{n'}) \leq d(x_n, l) + d(l, x_{n'}) < \epsilon/2 + \epsilon/2 = \epsilon,$$

luego la sucesión es de Cauchy.

Por otra parte, si la sucesión es de Cauchy, existe un $m \in \mathbb{N}$ tal que si $n \geq m$ se cumple $d(x_n, x_m) \leq 1$. Sea $C = \max(\{d(x_k, x_m) \mid k < m\} \cup \{1\})$. Es claro entonces que $d(x_n, x_m) \leq C$ para todo n , luego la sucesión está acotada. ■

Por ejemplo, como cabía esperar, todos los desarrollos decimales son sucesiones de Cauchy:

Teorema 1.29 *Para toda sucesión de números naturales tal que $0 \leq c_{-n} < 10$, la sucesión $\sum_{n=1}^{\infty} c_{-n}10^{-n}$ es de Cauchy en \mathbb{Q} .*

DEMOSTRACIÓN: Sea $S_N = \sum_{n=1}^N c_{-n}10^{-n}$. Tenemos que probar que la sucesión $\{S_N\}_{N=1}^{\infty}$ es de Cauchy. Sea $S'_N = \sum_{n=1}^N 9 \cdot 10^{-n}$. Entonces la sucesión $\{S'_N\}_{N=1}^{\infty}$ es de Cauchy, pues converge a 1. Además, si $N < M$,

$$|S_M - S_N| = \sum_{n=N}^M c_{-n}10^{-n} \leq \sum_{n=N}^M 9 \cdot 10^{-n} = |S'_M - S'_N|.$$

Así pues, dado $\epsilon > 0$, existe un N_0 tal que si $M, N \geq N_0$, se cumple que

$$|S_M - S_N| \leq |S'_M - S'_N| < \epsilon,$$

como exige la definición de sucesión de Cauchy. ■

Veamos otra relación entre las sucesiones de Cauchy y las convergentes:

Teorema 1.30 *Si una sucesión de Cauchy tiene una subsucesión convergente entonces converge al mismo límite.*

DEMOSTRACIÓN: Sea $\{x_n\}_{n=0}^{\infty}$ una sucesión de Cauchy en un espacio métrico M y sea $\{x_{n_k}\}_{k=0}^{\infty}$ una subsucesión convergente a l . Entonces, dado $\epsilon > 0$, existe un $m \in \mathbb{N}$ tal que si $n, n' \geq m$, entonces $d(x_n, x_{n'}) < \epsilon/2$ y si $k \geq m$ entonces $d(x_{n_k}, l) < \epsilon/2$. Así,

$$d(x_n, l) \leq d(x_n, x_{n_k}) + d(x_{n_k}, l) < \epsilon/2 + \epsilon/2 = \epsilon. \quad \blacksquare$$

Esto prueba, por ejemplo, que la sucesión $1, 2, 3, 1, 2, 3, \dots$, no es de Cauchy, pues tiene subsucesiones constantes y, en cambio, no es convergente.

Al igual que la convergencia, la propiedad de Cauchy se conserva por sumas y productos:

Teorema 1.31 *Sea K un cuerpo métrico y $\{a_n\}_{n=0}^{\infty}$, $\{b_n\}_{n=0}^{\infty}$, dos sucesiones de Cauchy en K . Entonces las sucesiones $\{a_n + b_n\}_{n=0}^{\infty}$ y $\{a_n b_n\}_{n=0}^{\infty}$ también son de Cauchy.*

DEMOSTRACIÓN: Para la suma observamos que

$$|a_{n'} + b_{n'} - a_n - b_n| \leq |a_{n'} - a_n| + |b_{n'} - b_n|,$$

luego, dado $\epsilon > 0$, tomando un $m \in \mathbb{N}$ tal que $|a_{n'} - a_n| < \epsilon/2$ y $|b_{n'} - b_n| < \epsilon/2$ siempre que $n', n \geq m$, tenemos también que $|a_{n'} + b_{n'} - a_n - b_n| < \epsilon$, luego la suma es de Cauchy.

Para el producto usamos que

$$|a_{n'} b_{n'} - a_n b_n| = |a_{n'} b_{n'} - a_{n'} b_n + a_{n'} b_n - a_n b_n| \leq |a_{n'}| |b_{n'} - b_n| + |a_{n'} - a_n| |b_n|,$$

así como que ambas sucesiones están acotadas, de modo que existe un $C \in \mathbb{R}$ tal que $|a_n| \leq C$, $|b_n| \leq C$ para todo n . Por lo tanto, dado $\epsilon > 0$, basta tomar un $m \in \mathbb{N}$ tal que $|a_{n'} - a_n| \leq \epsilon/2C$, $|b_{n'} - b_n| \leq \epsilon/2C$ siempre que $n, n' \geq m$. Obtenemos entonces que

$$|a_{n'} b_{n'} - a_n b_n| < C\epsilon/2C + C\epsilon/2C = \epsilon. \quad \blacksquare$$

Ahora ya podemos definir la noción de completitud que caracterizará a \mathbb{R} , que es precisamente la que habíamos indicado, la que afirma que todas las sucesiones “que deben converger” convergen. Esta propiedad de completitud tiene sentido en un espacio métrico arbitrario:

Definición 1.32 Un espacio métrico es *completo* si en él toda sucesión de Cauchy es convergente. En particular podemos hablar de cuerpos métricos y cuerpos ordenados arquimedianos completos.

En el caso de un cuerpo ordenado arquimediano, esto tiene varias caracterizaciones de interés:

Teorema 1.33 *Sea R un cuerpo ordenado arquimediano. Las afirmaciones siguientes son equivalentes:*

- a) *Toda sucesión de Cauchy en R es convergente.*
- b) *Si $\{[a_n, b_n]\}_{n=0}^{\infty}$ es una sucesión de intervalos encajados en R , existe un único $l \in K$ tal que $a_n \leq l \leq b_n$ para todo n .*
- c) *Todo subconjunto no vacío de R acotado superiormente tiene supremo.*

DEMOSTRACIÓN: a) \Rightarrow b) Supongamos que R es completo y consideremos una sucesión de intervalos encajados. Observamos que $\{a_n\}_{n=0}^{\infty}$ es de Cauchy pues, dado $\epsilon > 0$, tomamos m tal que $|b_m - a_m - 0| < \epsilon$, es decir, $b_m - a_m < \epsilon$. Entonces, si $n' \geq n \geq m$, tenemos que $a_m \leq a_n \leq a_{n'} \leq b_{n'} \leq b_m$, luego $a_{n'} - a_n \leq b_m - a_m < \epsilon$, luego $|a_{n'} - a_n| < \epsilon$ y, con el valor absoluto, la desigualdad vale igualmente si $n' \leq n$.

Por hipótesis existe $l = \lim_n a_n$. Como, para cada $m \in \mathbb{N}$ se cumple que si $n \geq m$ entonces $a_m \leq a_n \leq b_m$, el teorema 1.16 nos da que $a_m \leq l \leq b_m$. Ahora basta aplicar 1.26.

b) \Rightarrow c) Sea $A \subset R$ un subconjunto no vacío acotado superiormente. Sea $a \in A$ y sea b una cota superior. Si a es también una cota superior entonces es el máximo de A , luego también su supremo. Supongamos, pues, que a no es una cota superior de A . En particular $a < b$.

Vamos a construir recurrentemente una sucesión $\{(a_n, b_n)\}_{n=0}^{\infty}$ en $R \times R$ de modo que $a_n < b_n$ y cada b_n sea una cota superior de A , pero ningún a_n lo sea. Tomamos $(a_0, b_0) = (a, b)$, supuesto definido (a_n, b_n) , consideramos $c = (a_n + b_n)/2$, de modo que $a_n < c < b_n$, y distinguimos dos casos:

Si c no es cota superior de A , definimos $a_{n+1} = c$, $b_{n+1} = b_n$, mientras que si c es cota superior de A tomamos $a_{n+1} = a_n$, $b_{n+1} = c$.

Teniendo en cuenta que si $c = (a + b)/2$ entonces $b - c = c - a = (b - a)/2$, una simple inducción prueba que

$$b_n - a_n = \frac{b - a}{2^n}.$$

Esto implica que $\lim_n (b_n - a_n) = 0$, luego $\{[a_n, b_n]\}_{n=0}^{\infty}$ es una sucesión de intervalos encajados.

Por b) existe un $l \in R$ tal que $a_n \leq l \leq b_n$ para todo n . Vamos a probar que l es el supremo de A . Tiene que ser una cota superior, porque si no lo fuera existiría un $x \in A$ tal que $l < x$. Pero entonces sería $a_n \leq l < x \leq b_n$ y $\epsilon = x - l$ contradiría la convergencia de $b_n - a_n$ a 0.

Por otro lado, toda cota superior c de A cumple $l \leq c$, pues si fuera $c < l$ tendríamos que $a_n \leq c < l \leq b_n$ y de nuevo tenemos una contradicción con $\epsilon = l - c$.

c) \Rightarrow a) Dada una sucesión de Cauchy en R , para probar que es convergente, basta probar que lo es una cualquiera de sus subsucesiones. Por el teorema 1.7

podemos tomar una subsucesión monótona. Como toda sucesión de Cauchy es acotada y obviamente toda subsucesión de una sucesión acotada está acotada, basta probar que toda sucesión $\{x_n\}_{n=0}^{\infty}$ monótona y acotada es convergente.

No perdemos generalidad si suponemos que es monótona creciente, ya que si es decreciente la sucesión $\{-x_n\}_{n=0}^{\infty}$ es creciente, y si ésta converge a l , la sucesión inicial converge a $-l$ (porque la condición $|-x_n - l| < \epsilon$ es equivalente a $|x_n - (-l)| < \epsilon$).

Es claro que el conjunto $A = \{x_n \mid n \in \mathbb{N}\}$ es no vacío y está acotado superiormente, luego tiene supremo l . Vamos a probar que dicho supremo es el límite de la sucesión. Dado $\epsilon > 0$, por definición de supremo se cumple que $l - \epsilon$ no es cota superior de A , luego existe un $m \in \mathbb{N}$ tal que $l - \epsilon < x_m \leq l$, luego si $n \geq m$ tenemos que $l - \epsilon < x_m \leq x_n \leq l$, luego $|x_n - l| = l - x_n < l - (l - \epsilon) = \epsilon$. ■

Hay que destacar que en la tercera de las afirmaciones del teorema anterior sólo interviene la relación de orden, por lo que la propiedad correspondiente tiene sentido en conjuntos totalmente ordenados arbitrarios, aunque no sean cuerpos. Nos ocupamos de ella en la sección siguiente.

1.5 Cuerpos ordenados completos

Abandonamos momentáneamente los cuerpos ordenados para considerar un conjunto totalmente ordenado arbitrario:

Teorema 1.34 *Si X es un conjunto totalmente ordenado, las afirmaciones siguientes son equivalentes:*

- a) *Todo subconjunto de X no vacío y acotado superiormente tiene supremo.*
- b) *Todo subconjunto de X no vacío y acotado inferiormente tiene ínfimo.*
- c) *Un conjunto $I \subset X$ es un intervalo si y sólo si cuando $a < c < b$ con $a, b \in I$ y $c \in X$, entonces $c \in I$.*
- d) *Si $X = A \cup B$ de modo que $A \neq \emptyset \neq B$ y todo elemento de A es menor que todo elemento de B , entonces, o bien A tiene máximo, o bien B tiene mínimo.*

DEMOSTRACIÓN: a) \Rightarrow b) Sea B un subconjunto de X no vacío y acotado inferiormente. Sea A el conjunto de las cotas inferiores de B . Como B está acotado, A no es vacío, y como B no es vacío, cualquiera de sus elementos es una cota superior de A , luego A tiene supremo i , que es el ínfimo de B , pues todo elemento $b \in B$ es una cota superior de A , por lo que $i \leq b$, y si c es una cota inferior de B , entonces $c \in A$, luego $c \leq i$.

Análogamente se prueba que b) \Rightarrow a). Veamos que a) y b) implican c). Es inmediato que todo intervalo tiene la propiedad indicada. Se trata de probar el

recíproco. Supongamos, pues que I cumple la condición. Si $I = \emptyset$, entonces $I =]-\infty, +\infty[$ si $X = \emptyset$, o bien $I =]a, a[$ si existe un $a \in X$, luego es un intervalo.

Supongamos, pues, que $I \neq \emptyset$. Si I no está acotado ni superior ni inferiormente, entonces $I =]-\infty, +\infty[$, pues, para todo $x \in X$, como no es ni una cota superior ni una cota inferior de I , existen $a, b \in I$ tales que $a < x < b$, luego $x \in I$.

Si I tiene cota superior pero no inferior, tomamos $b = \sup I$ y observamos que $]-\infty, b[\subset I \subset]-\infty, b]$, con lo que I será uno de los dos intervalos según si $b \notin I$ o bien $b \in I$.

En efecto, si $x \in]-\infty, b[$, como x no es cota inferior de I existe un $a \in I$ tal que $a < x$ y, como b es el supremo de I , no puede ser que x sea una cota superior, luego existe un $c \in I$ tal que $a < x < c$, luego $x \in I$. La otra inclusión es trivial, puesto que b es una cota superior.

Los casos restantes (combinaciones de que I tenga o no tenga cota superior e inferior) se tratan análogamente.

c) \Rightarrow d) Si tenemos $u < x < v$ con $u, v \in A$ y $x \in X$, no puede ser $x \in B$, puesto que tendría que ser mayor que v , luego tiene que ser $x \in A$, luego A es un intervalo, y análogamente se razona que B lo es. Es claro que la única opción para A es ser un intervalo de la forma $A =]-\infty, c[$ o bien $A =]-\infty, c]$, en cuyo caso, B tiene que ser respectivamente de la forma $B = [c, +\infty[$ o bien $B =]c, +\infty[$. En el primer caso B tiene mínimo, y en el segundo A tiene máximo.

d) \Rightarrow a) Sea $C \subset X$ un conjunto no vacío y acotado superiormente. Llamemos A al conjunto de elementos de X que no son cotas superiores de C y B al conjunto de los elementos que sí que lo son. Obviamente $X = A \cup B$ y si $a \in A$ y $b \in B$, tenemos que a no es una cota superior de C , luego existe un $c \in C$ tal que $a < c$ y, como b es cota superior, $a < c \leq b$. Por d) existe $s \in X$ que es el máximo de A o bien el mínimo de B . Si es el mínimo de B , entonces es la menor cota superior de C , luego s es el supremo de C . Si s es el máximo de A , entonces existe un $c \in C$ tal que $s < c$, pero $c \in B$, luego c es una cota superior de C , luego c es el máximo, y en particular el supremo, de C . ■

Definición 1.35 Un conjunto totalmente ordenado X es *completo* si cumple cualquiera de las afirmaciones del teorema anterior.

Notemos que si X tiene máximo y mínimo la completitud equivale a que todo subconjunto de X tenga supremo e ínfimo, pues todo conjunto está acotado superior e inferiormente y \emptyset tiene al mínimo por supremo y al máximo por ínfimo.

En estos términos, el último teorema de la sección anterior afirma que un cuerpo ordenado arquimediano es completo como espacio métrico si y sólo si es completo como conjunto ordenado. Podemos precisar un poco más esta relación:

Teorema 1.36 *Todo cuerpo ordenado completo (como conjunto ordenado) es arquimediano y, por consiguiente, completo como espacio métrico.*

DEMOSTRACIÓN: Sea R un cuerpo ordenado completo. Si no es arquimediano, entonces \mathbb{N} está acotado superiormente, luego tiene supremo, digamos s . Por definición de supremo, $s - 1/2$ no es cota superior de \mathbb{N} , luego existe un $n \in \mathbb{N}$ tal que $s - 1/2 < n \leq s$, pero entonces $s < n + 1/2 < n + 1$, en contradicción con que s sea cota superior de \mathbb{N} . ■

Así pues, en lo sucesivo podemos hablar de cuerpos ordenados completos, sin necesidad de exigir que sean arquimedianos, pues sería redundante. La completitud de un cuerpo tiene muchas consecuencias algebraicas:

Teorema 1.37 *Si R es un cuerpo ordenado completo, para cada $x \in R$, $x \geq 0$ existe un único $y \in R$, $y \geq 0$ tal que $y^2 = x$.*

DEMOSTRACIÓN: Podemos suponer que $x > 0$. Consideremos el conjunto $A = \{u \in R \mid u > 0, u^2 < x\}$. Como R es arquimediano, existen números naturales $x < n < n^2$ y $1/x < m < m^2$, con lo que $1/m^2 < x$ y así $1/m \in A$ y n es una cota superior de A . Esto implica que A tiene supremo. Llamémoslo y . Claramente $y > 0$.

Supongamos que $x < y^2$. Tomemos un número natural n tal que $n > 1/y$ y $n > 2y/(y^2 - x)$. Así $2y/n < y^2 - x$ y en consecuencia

$$\left(y - \frac{1}{n}\right)^2 = y^2 - 2y\frac{1}{n} + \frac{1}{n^2} > y^2 - y^2 + x + \frac{1}{n^2} > x.$$

Así, si $u \in A$ tenemos que $u^2 < x < (y - 1/n)^2$, luego $u < y - 1/n$, pero esto significa que $y - 1/n$ es una cota superior de A , en contradicción con que y es el supremo.

Supongamos ahora que $y^2 < x$. Entonces tomamos un número natural n que cumpla $n > 4y/(x - y^2)$ y $n^2 > 2/(x - y^2)$. Así

$$\left(y + \frac{1}{n}\right)^2 = y^2 + 2y\frac{1}{n} + \frac{1}{n^2} < y^2 + \frac{x - y^2}{2} + \frac{x - y^2}{2} = y^2 + x - y^2 = x,$$

luego $y + 1/n \in A$, y esto supone de nuevo una contradicción. Por lo tanto ha de ser $y^2 = x$.

La unicidad es clara, pues si $z \in R$ cumple $z > 0$, $z \neq y$, entonces $z^2 < y^2$ o $z^2 > y^2$ según si $z < y$ o $z > y$. ■

Definición 1.38 Si R es un cuerpo ordenado completo y $x \in R$ cumple $x \geq 0$, llamaremos *raíz cuadrada positiva* de x al único $y \in R$ que cumple $y^2 = x$, y lo representaremos por $y = \sqrt{x}$.

Notemos que \sqrt{x} no es la única raíz cuadrada de x en R , pues $-\sqrt{x}$ es otra (distinta siempre que $x \neq 0$), y no hay más, porque el polinomio $t^2 - x$ sólo puede tener dos raíces en R .

La existencia de raíces cuadradas implica que en un cuerpo ordenado completo la relación de orden está determinada por la suma y el producto: se cumple $x \leq y$ si y sólo si existe un $z \in R$ tal que $y - x = z^2$.

Con esto estamos casi a punto de demostrar el resultado principal que estamos persiguiendo: salvo isomorfismo, existe a lo sumo un cuerpo ordenado completo. Pero antes hay que precisar qué significa “salvo isomorfismo”:

Definición 1.39 Una aplicación $f : M \rightarrow N$ entre dos espacios métricos es una *inmersión isométrica* si conserva las distancias, es decir, si para todos los puntos $x, y \in M$, se cumple $d(f(x), f(y)) = d(x, y)$.

Notemos que una inmersión isométrica es necesariamente inyectiva, por la primera condición de la definición de distancia. Una inmersión isométrica biyectiva entre dos espacios métricos es una *isometría*.

Dos espacios métricos son *isométricos* cuando existe una isometría entre ellos. Esto hace que ambos compartan las mismas propiedades definibles en términos de la distancia.

Una *inmersión isométrica* (resp. *isometría*) entre dos cuerpos métricos es un monomorfismo (resp. isomorfismo) de cuerpos $f : K \rightarrow L$ tal que para todo $x \in K$ se cumple $|f(x)| = |x|$.

Es claro que las isometrías de cuerpos métricos son isometrías de espacios métricos, y hacen que dos cuerpos métricos isométricos compartan las mismas propiedades definibles en términos de la métrica, la suma y el producto.

Un *isomorfismo de cuerpos ordenados* es un isomorfismo que además conserve el orden, es decir, que cumpla $x \leq y$ si y sólo si $f(x) \leq f(y)$. Es claro que todo isomorfismo de cuerpos ordenados arquimedianos es una isometría de cuerpos métricos, pues al conservarse el orden se conservan el valor absoluto y la distancia, que se definen a partir de ellos.

Teorema 1.40 Sean M y N dos espacios métricos completos, sea $M_0 \subset M$ un subconjunto denso, considerado como espacio métrico con la restricción de la distancia de M . Sea $f : M_0 \rightarrow N$ una inmersión isométrica. Entonces existe una única inmersión isométrica $F : M \rightarrow N$ que extiende a f (es decir, tal que $F|_{M_0} = f$). Si $f[M_0]$ es denso en N , entonces F es una isometría. Más aún, si M y N son cuerpos métricos, M_0 es un subcuerpo y f es una inmersión isométrica (resp. isometría) de cuerpos métricos, entonces F también lo es.

DEMOSTRACIÓN: Sea $x \in M$. El teorema 1.20 implica que existe una sucesión $\{x_n\}_{n=0}^{\infty}$ en M_0 convergente a x . En particular es una sucesión de Cauchy en M , pero es claro que esto es lo mismo que decir que $\{x_n\}_{n=0}^{\infty}$ es una sucesión de Cauchy en M_0 como espacio métrico, pues la definición de sucesión de Cauchy depende sólo de los valores que toma la distancia sobre los términos de la sucesión. Como f es una inmersión isométrica, resulta que la sucesión $\{f(x_n)\}_{n=0}^{\infty}$ es de Cauchy en $f[M_0]$, luego en N , y por la completitud de N converge en N . Definimos $F(x) = \lim_n f(x_n)$.

Para que esta definición sea correcta tenemos que comprobar que no depende de la sucesión utilizada para calcular $F(x)$. En efecto, si tomamos otra sucesión $\{y_n\}_{n=0}^{\infty}$ en M_0 convergente a x , es fácil ver que la sucesión $\{z_k\}_{k=0}^{\infty}$ dada por

$$z_k = \begin{cases} x_n & \text{si } k = 2n, \\ y_n & \text{si } k = 2n + 1, \end{cases}$$

también converge a x . Según hemos visto, entonces $\{f(z_k)\}_{k=0}^{\infty}$ converge a un cierto z en N , pero las sucesiones $\{f(x_n)\}_{n=0}^{\infty}$ y $\{f(y_n)\}_{n=0}^{\infty}$ son subsucesiones de esta sucesión, luego ambas convergen al mismo límite z .

Por lo tanto F está bien definida. Veamos ahora que es una inmersión isométrica, es decir, que cumple $d(F(x), F(y)) = d(x, y)$.

Para ello observamos que si $\lim_n x_n = x$ y $\lim_n y_n = y$, donde ambas sucesiones están en M_0 , entonces $\lim_n d(x_n, y_n) = d(x, y)$. En efecto,

$$\begin{aligned} |d(x, y) - d(x_n, y_n)| &\leq |d(x, y) - d(x_n, y)| + |d(x_n, y) - d(x_n, y_n)| \\ &\leq d(x_n, x) + d(y_n, y), \end{aligned}$$

luego, dado $\epsilon > 0$, podemos tomar un $m \in \mathbb{N}$ tal que si $n \geq m$ se cumple $d(x_n, x) < \epsilon/2$ y $d(y_n, y) < \epsilon/2$, con lo que $|d(x, y) - d(x_n, y_n)| < \epsilon$.

Igualmente, $\lim_n d(f(x_n), f(y_n)) = d(F(x), F(y))$, pero ambas sucesiones de distancias son la misma, luego los límites son el mismo.

Si $f[M_0]$ es denso en N entonces F es suprayectiva, pues dado $y \in N$, podemos tomar una sucesión en $f[M_0]$ convergente a y y pasarla a M_0 con f , con lo que obtenemos una sucesión de Cauchy convergente a un cierto $x \in M$ que claramente cumplirá $F(x) = y$.

La unicidad de F es clara, pues si G fuera otra inmersión isométrica tal que $G|_{M_0} = f$, entonces, dado $x \in M$, lo expresamos como $x = \lim_n x_n$, donde la sucesión está en M_0 , y entonces $d(G(x), G(x_n)) = d(x, x_n)$, de donde se sigue inmediatamente que

$$G(x) = \lim_n G(x_n) = \lim_n f(x_n) = F(x).$$

Supongamos ahora que M y N son cuerpos métricos, que M_0 es un subcuerpo y que f es una inmersión isométrica de cuerpos métricos. Entonces, dados $x = \lim_n x_n$, $y = \lim_n y_n$ en M , el teorema 1.14 nos da que $x + y = \lim_n (x_n + y_n)$, luego podemos calcular F con esta sucesión y entonces

$$F(x + y) = \lim_n (f(x_n) + f(y_n)) = \lim_n f(x_n) + \lim_n f(y_n) = F(x) + F(y),$$

e igualmente se razona con el producto. Por lo tanto F es un monomorfismo de cuerpos. Como $|x| = d(x, 0)$, el hecho de que F sea una inmersión isométrica implica que $|F(x)| = |x|$. ■

En particular:

Teorema 1.41 *Dos cuerpos ordenados completos cualesquiera son isomorfos.*

DEMOSTRACIÓN: Supongamos que R y R' son cuerpos ordenados completos. Entonces ambos contienen a \mathbb{Q} o, si queremos ser más precisos, existen subcuerpos densos $\mathbb{Q} \subset R$, $\mathbb{Q}' \subset R'$ y un isomorfismo de cuerpos ordenados $f: \mathbb{Q} \rightarrow \mathbb{Q}'$.

El teorema anterior nos da que el isomorfismo f se extiende a una única isometría $F: R \rightarrow R'$ de cuerpos métricos, pero vamos a probar que además conserva el orden, y por lo tanto un isomorfismo de cuerpos ordenados. Esto es una consecuencia inmediata del teorema 1.37: se cumple $x \leq y$ si y sólo si existe $z \in R$ tal que $y - x = z^2$, si y sólo si existe $z' \in R'$ tal que $F(y) - F(x) = z'^2$, si y sólo si $F(x) \leq F(y)$. ■

Definición 1.42 Llamaremos *cuerpo \mathbb{R} de los números reales* a cualquier cuerpo ordenado completo.

En la sección 2.2 de [G] se construye un cuerpo ordenado completo, y el teorema anterior prueba que dos cualesquiera de ellos son isomorfos, por lo que es indiferente trabajar con uno o con otro. En el Apéndice A presentamos otra construcción de \mathbb{R} basada en el concepto de sucesión de Cauchy y que tiene la ventaja de que se generaliza a espacios métricos y cuerpos métricos arbitrarios.

Así pues, \mathbb{R} es un cuerpo ordenado arquimediano caracterizado salvo isomorfismo por su completitud, que viene expresada por cualquiera de las condiciones siguientes:

- Todo subconjunto no vacío de \mathbb{R} acotado superiormente tiene supremo, (y si está acotado inferiormente tiene ínfimo).
- Una sucesión en \mathbb{R} es convergente si y sólo si es de Cauchy.
- Toda sucesión de intervalos encajados en \mathbb{R} tiene intersección no vacía.
- Un conjunto $I \subset \mathbb{R}$ es un intervalo si y sólo si cuando $a < c < b$ con $a, b \in I$ y $c \in \mathbb{R}$, entonces $c \in I$.
- Si $\mathbb{R} = A \cup B$ de modo que $A \neq \emptyset \neq B$ y todo elemento de A es menor que todo elemento de B , entonces, o bien A tiene máximo, o bien B tiene mínimo.

A éstas caracterizaciones de la completitud de \mathbb{R} podemos añadir una más, y es que todo desarrollo decimal es convergente en \mathbb{R} . Esto es consecuencia inmediata del teorema siguiente:¹¹

Teorema 1.43 *Todo cuerpo ordenado arquimediano es isomorfo a un subcuerpo de \mathbb{R} .*

¹¹Si R es un cuerpo ordenado arquimediano en el que todo desarrollo decimal converge, por el teorema siguiente podemos suponer que R es un subcuerpo de \mathbb{R} , pero todo $\alpha \in \mathbb{R}$ admite un desarrollo decimal que debe converger en R , luego $\alpha \in R$, y así $R = \mathbb{R}$.

DEMOSTRACIÓN: Basta aplicar la prueba del teorema 1.41 a un cuerpo ordenado arquimediano R y a \mathbb{R} . Como se indica allí, R contiene un subcuerpo denso \mathbb{Q}' isomorfo a \mathbb{Q} , y sólo hemos de observar que para extender un isomorfismo $f : \mathbb{Q}' \rightarrow \mathbb{Q}$ a una inmersión isométrica $F : R \rightarrow \mathbb{R}$ no hace falta que R sea completo, pues la primera parte de la prueba del teorema 1.40 no requiere que M sea completo. ■

Así pues, en las definiciones de valor absoluto y distancia, no perdemos generalidad si suponemos que el cuerpo ordenado arquimediano en el que toman valores es precisamente \mathbb{R} .

1.6 El cardinal del continuo

En esta última sección demostraremos que existe una diferencia puramente conjuntista entre el conjunto de los números racionales y el de los números reales, y es que el primero es numerable y el segundo no. El concepto de numerabilidad está expuesto en el apéndice B de [A1].

Teorema 1.44 \mathbb{Z} y \mathbb{Q} son conjuntos numerables.

DEMOSTRACIÓN: La aplicación $f : \{\pm 1\} \times \mathbb{N} \rightarrow \mathbb{Z}$ dada por $f(e, n) = en$ es claramente suprayectiva y $|\{\pm 1\} \times \mathbb{N}| = 2 \cdot \aleph_0 = \aleph_0$, luego \mathbb{Z} es numerable por el teorema [A1 B.4].

A su vez, la aplicación $g : \mathbb{Z} \times (\mathbb{N} \setminus \{0\}) \rightarrow \mathbb{Q}$ dada por $g(m, n) = m/n$ es suprayectiva, y los dos factores del dominio de g son numerables, luego concluimos igualmente que \mathbb{Q} es numerable. ■

Para probar que \mathbb{R} no es numerable nos basamos en lo siguiente:

Teorema 1.45 Si un conjunto A tiene al menos dos elementos y B es infinito numerable, entonces el conjunto A^B de todas las aplicaciones $f : B \rightarrow A$ no es numerable.

DEMOSTRACIÓN: Sea $A' = \{a_0, a_1\} \subset A$ un conjunto con dos elementos. Como $A'^B \subset A^B$, basta probar que el primer conjunto es no numerable o, equivalentemente, podemos suponer que $A = \{a_0, a_1\}$. Entonces, la aplicación $G : A^B \rightarrow \mathcal{P}B$ dada por $G(f) = f^{-1}[a_1]$ es biyectiva, y $\mathcal{P}B$ no es numerable por el teorema de Cantor, luego A^B tampoco. ■

Teorema 1.46 Todo intervalo no trivial de números reales es no numerable. En particular, \mathbb{R} no es numerable.

DEMOSTRACIÓN: Por intervalos triviales entendemos los de la forma $[a, a]$, que tienen un elemento, y los de la forma $]a, a[$, $]a, a]$, $[a, a[$, que son vacíos.

Sea F el conjunto de las aplicaciones $f : \mathbb{N} \rightarrow \{0, 1\}$, que por el teorema anterior es no numerable. Equivalentemente, F es el conjunto de todas las sucesiones de ceros y unos. Sea $h : F \rightarrow [0, 1]$ la aplicación que a cada $\{a_n\}_{n=0}^\infty \in F$ le asigna el número real cuya expresión decimal es $0.a_0a_1a_2\dots$

Como ninguna expresión de este tipo es finalmente constante igual a 9 (sólo hay cifras iguales a 0 o a 1) la aplicación h es inyectiva, luego $[0, 1]$ no puede ser numerable.

Si $[a, b]$ es cualquier intervalo cerrado con $a < b$, se comprueba fácilmente que la aplicación $g : [0, 1] \rightarrow [a, b]$ dada por $g(t) = ta + (1 - t)b$ es biyectiva, luego todos los intervalos cerrados no triviales son no numerables. Pero es claro que todo intervalo no trivial contiene un intervalo cerrado no trivial, luego todos son no numerables. ■

Por consiguiente, el conjunto de los números irracionales es no numerable (si fuera numerable, al unirlo a \mathbb{Q} , que es numerable, obtendríamos que \mathbb{R} es numerable). Más aún, todo intervalo contiene una cantidad no numerable de números irracionales. En particular, entre dos números reales cualesquiera existen infinitos irracionales.¹²

La frontera entre la numerabilidad de \mathbb{Q} y la no numerabilidad de \mathbb{R} puede localizarse mejor. Para ello probamos un resultado auxiliar:

Teorema 1.47 *Si A es un anillo conmutativo y unitario numerable y X es un conjunto numerable no vacío, el anillo de polinomios $A[X]$ es infinito numerable.*

DEMOSTRACIÓN: Consideremos primero el caso en que $X = \{x\}$. Sea $A_n[x]$ el conjunto de los polinomios de grado $\leq n$ y veamos por inducción que es numerable. Como $A_0[x] = A$, se cumple para $n = 0$. Si vale para n , consideramos la aplicación $f : A \times A_n[x] \rightarrow A_{n+1}[x]$ dada por $f(a, p(x)) = ax^{n+1} + p(x)$. Claramente es suprayectiva, luego $A_{n+1}[x]$ es numerable.

Entonces $A[x] = \bigcup_{n \in \mathbb{N}} A_n[x]$ es numerable por ser unión numerable de conjuntos numerables.

Usando que $A[x_1, \dots, x_{n+1}] = A[x_1, \dots, x_n][x]$, una simple inducción implica que el teorema vale siempre que el conjunto de indeterminadas X es finito. Por último, si X es infinito numerable, como en cada polinomio aparece un número finito de indeterminadas, se cumple que

$$A[X] = \bigcup_{Y \in \mathcal{P}^f X} A[Y]$$

luego $A[X]$ es numerable por ser unión numerable de conjuntos numerables. ■

Definición 1.48 Un número real se dice *algebraico* si es la raíz de un polinomio no nulo de $\mathbb{Q}[x]$. En caso contrario se dice *trascendente*.

Así, $\sqrt{2}$ es irracional, pero algebraico, porque es la raíz del polinomio $x^2 - 2$.

¹²Este hecho puede probarse también mediante un argumento mucho más elemental: dados dos números reales $\alpha < \beta$, se cumple que $\alpha/\sqrt{2} < \beta/\sqrt{2}$, luego existe un número racional $\alpha/\sqrt{2} < r < \beta/\sqrt{2}$, luego $\alpha < r\sqrt{2} < \beta$, y $r\sqrt{2}$ es irracional, ya que si fuera racional también lo sería $\sqrt{2}$, por ser cociente de números racionales.

Teorema 1.49 *El conjunto de los números reales algebraicos es numerable. En particular, todo intervalo no trivial contiene infinitos números trascendentes.*

DEMOSTRACIÓN: Para cada $p \in \mathbb{Q}[x] \setminus \{0\}$, sea R_p el conjunto de sus raíces en \mathbb{R} , que es finito, pues un polinomio no nulo tiene a lo sumo tantas raíces en un cuerpo como indica su grado. Entonces el conjunto de los números reales algebraicos es

$$\bigcup_{p \in \mathbb{Q}[x] \setminus \{0\}} R_p,$$

que es unión numerable de conjuntos finitos, luego es numerable. ■

Acabamos de probar que existen infinitos números reales trascendentes, pero la prueba es esencialmente no constructiva, hasta el punto de que no nos permite justificar que ninguno en concreto lo sea.

Definición 1.50 *Llamaremos $\mathfrak{c} = \mathcal{P}\mathbb{N}$ y diremos que un conjunto X tiene cardinal \mathfrak{c} , y lo representaremos por $|X| = \mathfrak{c}$ si existe $f : X \rightarrow \mathcal{P}\mathbb{N}$ biyectiva.*

El teorema de Cantor prueba que los conjuntos de cardinal \mathfrak{c} no son numerables. El teorema siguiente permite probar fácilmente que casi todos los conjuntos no numerables que nos vamos a encontrar en la práctica tendrán cardinal \mathfrak{c} :

Teorema 1.51 *Los conjuntos siguientes tienen cardinal \mathfrak{c} :*

$\mathcal{P}\mathbb{N}$, $\mathbb{N}^{\mathbb{N}}$, $\{1, 2, \dots, n\}^{\mathbb{N}}$ ($n \geq 2$), \mathbb{R} , $\mathbb{R}^{\mathbb{N}}$, todo intervalo no trivial en \mathbb{R} .

DEMOSTRACIÓN: $|\mathcal{P}\mathbb{N}| = \mathfrak{c}$ por definición. Llamemos $I_n = \{1, \dots, n\}$. En la prueba de 1.45 hemos visto que existe una aplicación $f : I_2^{\mathbb{N}} \rightarrow \mathcal{P}\mathbb{N}$ biyectiva, luego $f^{-1} : \mathcal{P}\mathbb{N} \rightarrow I_2^{\mathbb{N}}$ es inyectiva. Por otra parte, toda aplicación de \mathbb{N} en I_n es un subconjunto de $\mathbb{N} \times I_n$, es decir, $I_n^{\mathbb{N}} \subset \mathcal{P}(\mathbb{N} \times I_n)$. Ahora bien, existe una biyección de $\mathbb{N} \times I_n$ en \mathbb{N} , la cual induce una biyección $\mathcal{P}(\mathbb{N} \times I_n) \rightarrow \mathcal{P}\mathbb{N}$, la cual se restringe a su vez a una aplicación $I_n^{\mathbb{N}} \rightarrow \mathcal{P}\mathbb{N}$ inyectiva. Por el teorema de Cantor-Bernstein concluimos que existe $I_n^{\mathbb{N}} \rightarrow \mathcal{P}\mathbb{N}$ biyectiva, luego $|I_n^{\mathbb{N}}| = \mathfrak{c}$. Exactamente el mismo argumento prueba que $|\mathbb{N}^{\mathbb{N}}| = \mathfrak{c}$.

En la prueba del teorema 1.46 se ve que existe una aplicación $I_2^{\mathbb{N}} \rightarrow [0, 1]$ inyectiva, luego también existe una aplicación $\mathcal{P}\mathbb{N} \rightarrow [0, 1]$ inyectiva. Por otra parte, la aplicación $\mathbb{R} \rightarrow \mathcal{P}\mathbb{Q}$ que a cada $\alpha \in \mathbb{R}$ le asigna el conjunto $\{r \in \mathbb{Q} \mid r < \alpha\}$ es claramente inyectiva, y una biyección $\mathbb{Q} \rightarrow \mathbb{N}$ induce una biyección $\mathcal{P}\mathbb{Q} \rightarrow \mathcal{P}\mathbb{N}$, luego en total tenemos una aplicación inyectiva $\mathbb{R} \rightarrow \mathcal{P}\mathbb{N}$.

El teorema de Cantor-Bernstein nos da biyecciones $\mathcal{P}\mathbb{N} \rightarrow \mathbb{R}$ y $\mathcal{P}\mathbb{N} \rightarrow [0, 1]$. Con esto no sólo tenemos que $|\mathbb{R}| = \mathfrak{c}$, sino que la prueba de 1.46 muestra ahora que todos los intervalos $[a, b] \subset \mathbb{R}$ tienen cardinal \mathfrak{c} , y de ahí es fácil deducir que todos los intervalos, más aún, todo subconjunto de \mathbb{R} que contenga un intervalo, tiene cardinal \mathfrak{c} .

Una biyección $\mathbb{R} \rightarrow \mathbb{N}^{\mathbb{N}}$ induce una biyección $\mathbb{R}^{\mathbb{N}} \rightarrow (\mathbb{N}^{\mathbb{N}})^{\mathbb{N}}$, luego basta probar que el último conjunto tiene cardinal \mathfrak{c} , pero es fácil definir biyecciones $(\mathbb{N}^{\mathbb{N}})^{\mathbb{N}} \rightarrow \mathbb{N}^{\mathbb{N} \times \mathbb{N}} \rightarrow \mathbb{N}^{\mathbb{N}}$. ■

Para terminar de perfilar la situación probamos un último teorema:

Teorema 1.52 Si $A \subset B$, $|A| \leq \aleph_0$ y $|B| = \mathfrak{c}$, entonces $|B \setminus A| = \mathfrak{c}$.

DEMOSTRACIÓN: En primer lugar, $B \setminus A$ tiene que ser no numerable, porque en caso contrario $B = A \cup (B \setminus A)$ sería numerable. Por el teorema [Al B.8] podemos tomar $N \subset B \setminus A$ numerable. Sea $C = (B \setminus A) \setminus N$. Entonces

$$B = A \cup N \cup C, \quad B \setminus A = N \cup C,$$

y todas las uniones son disjuntas. Una biyección entre $A \cup N$ y N se extiende a una biyección $B \rightarrow B \setminus A$ sin más que dejar invariantes a los elementos de C . ■

Esto implica que $|\mathbb{R} \setminus \mathbb{Q}| = \mathfrak{c}$, y también que hay \mathfrak{c} números reales trascendentes, por ejemplo.

Por supuesto, es fácil encontrar conjuntos que tengan cardinal mayor que \mathfrak{c} , como pueden ser $\mathcal{P}\mathbb{R}$ o $\mathbb{R}^{\mathbb{R}}$.

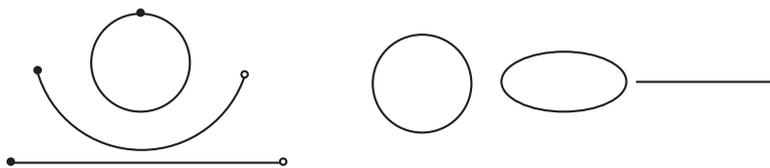
Capítulo II

Topología

La topología puede considerarse como la forma más abstracta de la geometría. El concepto principal que puede definirse a partir de la estructura topológica es el de aplicación continua, que viene a ser una transformación realizada sin cortes o saltos bruscos o, dicho de otro modo, que transforma puntos próximos en puntos próximos. Los resultados topológicos son aplicables tanto a la geometría propiamente dicha como a la descripción de otros muchos objetos más cercanos a la teoría de conjuntos general, si bien aquí nos centraremos en la vertiente geométrica. Al combinarla con el álgebra obtendremos el cálculo diferencial, que constituye la herramienta más potente para el estudio de la geometría.

2.1 Espacios topológicos

Según acabamos de indicar, una aplicación continua es una aplicación que transforma puntos próximos en puntos próximos. Nuestro objetivo ahora es definir una estructura matemática en la que esta afirmación pueda convertirse en una definición rigurosa. En primer lugar conviene reformularla así: una aplicación continua es una aplicación que transforma los puntos de alrededor de un punto dado en puntos de alrededor de su imagen. En efecto, si cortamos una circunferencia por un punto P para convertirla en un segmento, la transformación no es continua, pues los puntos de alrededor de P se transforman unos en los puntos de un extremo del segmento y otros en los puntos del otro extremo, luego no quedan todos alrededor del mismo punto. En cambio, podemos transformar continuamente (aunque no biyectivamente) una circunferencia en un segmento sin más que aplastarla.



Para llegar a la noción de “espacio topológico”, que es la estructura mínima respecto a la cual es posible dar una definición precisa de “continuidad”, partiremos de la formalización algebraica del concepto intuitivo de “espacio”. Ésta se lleva a cabo a través de la estructura de espacio afín euclídeo tridimensional sobre un cuerpo R que cumpla unos requisitos mínimos para que se cumplan las propiedades geométricas básicas, entre los que se encuentra que R sea un subcuerpo del cuerpo \mathbb{R} de los números reales. No obstante, así como la tridimensionalidad del espacio está en la base de nuestra intuición geométrica, desde un punto de vista algebraico es posible trabajar exactamente igual con espacios afines euclídeos de cualquier dimensión finita. Recordemos que un espacio afín euclídeo E tiene asociado un espacio vectorial \vec{E} de dimensión finita en el cual hay definido un producto vectorial $\vec{E} \times \vec{E} \rightarrow \mathbb{R}$. Éste a su vez permite definir la norma de un vector como $\|\vec{v}\| = \sqrt{\vec{v}\vec{v}}$ y el ángulo entre dos vectores, mediante la relación

$$\cos \widehat{\vec{v}\vec{w}} = \frac{\vec{v}\vec{w}}{\|\vec{v}\| \|\vec{w}\|}.$$

A su vez la norma permite definir la distancia entre dos puntos de E como $d(P, Q) = \|\overrightarrow{PQ}\|$, que puede verse a su vez como la longitud del segmento que une P con Q .

Estas estructuras son demasiado particulares y restrictivas desde el punto de vista topológico. La medida de ángulos es un sinsentido en topología, y la de longitudes tiene un interés secundario, pues no importan las medidas concretas, sino tan sólo la noción de proximidad. Ello nos lleva a una primera abstracción de estas ideas que conserva los mínimos elementos que van a tener alguna relevancia en este contexto, aunque sea marginal:

Definición 2.1 Sea \mathbb{K} un cuerpo métrico y E un espacio vectorial sobre \mathbb{K} . Una *norma* en E es una aplicación $\|\cdot\| : E \rightarrow [0, +\infty[$ que cumpla las propiedades siguientes:

- a) $\|v\| = 0$ si y sólo si $v = 0$.
- b) $\|v + w\| \leq \|v\| + \|w\|$.
- c) $\|\alpha v\| = |\alpha| \|v\|$,

para $v, w \in E$ y todo $\alpha \in \mathbb{K}$.

Un *espacio normado* es un par $(E, \|\cdot\|)$ en estas condiciones. En la práctica escribiremos E , sin indicar explícitamente la norma.

Es conocido [G 4.21] que la norma derivada de un producto escalar en el espacio vectorial asociado a un espacio afín euclídeo cumple estas propiedades, pero aquí no exigimos que la norma derive de ningún producto escalar, admitimos que el espacio E tenga dimensión infinita y admitimos una clase más amplia de cuerpos de escalares.

Veamos algunos ejemplos de normas en \mathbb{R}^n :

Teorema 2.2 \mathbb{R}^n es un espacio normado con cualquiera de estas normas:

$$\|x\|_1 = \sum_{i=1}^n |x_i|, \quad \|x\|_2 = \sqrt{\sum_{i=1}^n |x_i|^2}, \quad \|x\|_\infty = \max\{|x_i| \mid i = 1, \dots, n\}.$$

La única propiedad que no es inmediata es la desigualdad triangular para la segunda norma, pero dicha norma es precisamente la norma euclídea, es decir, la norma derivada del producto escalar en \mathbb{R}^n , y el hecho de que cumple las propiedades de la definición de norma es un caso particular de [G 4.21].

El hecho de que estas tres aplicaciones sean normas permite obtener un resultado más general:

Teorema 2.3 Sean E_1, \dots, E_n espacios normados. Entonces las aplicaciones siguientes son normas en $E = E_1 \times \dots \times E_n$.

$$\|x\|_1 = \sum_{i=1}^n \|x_i\|, \quad \|x\|_2 = \sqrt{\sum_{i=1}^n \|x_i\|^2}, \quad \|x\|_\infty = \max\{\|x_i\| \mid i = 1, \dots, n\}.$$

Además se cumplen las relaciones: $\|x\|_\infty \leq \|x\|_2 \leq \|x\|_1 \leq n\|x\|_\infty$.

DEMOSTRACIÓN: Tenemos $\|x\|_i = \|(\|x_1\|, \dots, \|x_n\|)\|_i$ para $i = 1, 2, \infty$. Usando el teorema anterior se ve inmediatamente que son normas.

$$\begin{aligned} \|x\|_\infty &= \sqrt{\|x\|_\infty^2} \leq \sqrt{\sum_{i=1}^n \|x_i\|^2} = \|x\|_2 \leq \sqrt{\sum_{i=1}^n \|x_i\|^2 + 2 \sum_{i < j} \|x_i\| \|x_j\|} \\ &= \sqrt{\left(\sum_{i=1}^n \|x_i\|\right)^2} = \|x\|_1 \leq \sum_{i=1}^n \|x\|_\infty = n\|x\|_\infty. \quad \blacksquare \end{aligned}$$

Observemos que un cuerpo métrico \mathbb{K} es un espacio normado tomando como norma el propio valor absoluto de \mathbb{K} , y aplicando el teorema anterior con todos los espacios iguales a \mathbb{K} obtenemos que el teorema 2.2 vale con \mathbb{K} en lugar de \mathbb{R} .

Ejercicio: Probar que en un espacio normado se cumple $|\|x\| - \|y\|| \leq \|x - y\|$.

El concepto de norma tiene interés en nuestro contexto porque nos permitirá relacionar los conceptos topológicos con los algebraicos, pero a la hora de aislar el concepto de espacio topológico, las normas contienen todavía información irrelevante. Obtenemos un grado más de abstracción pasando al concepto de espacio métrico que ya introdujimos en el capítulo anterior, pero recordamos aquí la definición:

Definición 2.4 Una *distancia* en un conjunto M es una aplicación

$$d : M \times M \longrightarrow \mathbb{R},$$

que cumpla las propiedades siguientes:

- a) $d(x, y) \geq 0$ y $d(x, y) = 0$ únicamente cuando $x = y$,
- b) $d(x, y) = d(y, x)$,
- c) $d(x, z) \leq d(x, y) + d(y, z)$,

de las cuales se deduce una cuarta: $|d(x, y) - d(x, z)| \leq d(y, z)$.

Un *espacio métrico* es un par (M, d) , donde M es un conjunto y d es una distancia en M .

Todo espacio normado E es un espacio métrico con la distancia definida por $d(x, y) = \|x - y\|$. Las propiedades de la definición de norma implican inmediatamente las de la definición de distancia. En particular en \mathbb{K}^n tenemos definidas tres distancias:

$$d_1(x, y) = \sum_{i=1}^n |x_i - y_i|, \quad d_2(x, y) = \sqrt{\sum_{i=1}^n |x_i - y_i|^2},$$

$$d_\infty(x, y) = \max\{|x_i - y_i| \mid 1 \leq i \leq n\}.$$

La geometría nos enseña que la distancia d_2 en \mathbb{R}^n es la distancia euclídea usual, es decir, la distancia que (al menos cuando $n = 1, 2, 3$) se corresponde con la noción intuitiva de distancia. Veremos que las otras distancias son formas alternativas de medir la proximidad entre puntos que, aunque no sean “la natural”, serán equivalentes a efectos topológicos (precisamente porque el valor exacto de la distancia será irrelevante, y lo único que importará es la proximidad o lejanía entre puntos).

Más en general, las fórmulas anteriores permiten definir distancias en cualquier producto finito de espacios métricos. La prueba del teorema siguiente es muy sencilla a partir de los teoremas 2.2 y 2.3.

Teorema 2.5 Sean M_1, \dots, M_n espacios métricos. Sea $M = M_1 \times \dots \times M_n$. Entonces las aplicaciones $d_1, d_2, d_\infty : M \times M \rightarrow [0, +\infty[$ definidas como sigue son distancias en M :

$$d_1(x, y) = \sum_{i=1}^n d(x_i, y_i),$$

$$d_2(x, y) = \sqrt{\sum_{i=1}^n d(x_i, y_i)^2},$$

$$d_\infty(x, y) = \max\{d(x_i, y_i) \mid 1 \leq i \leq n\}.$$

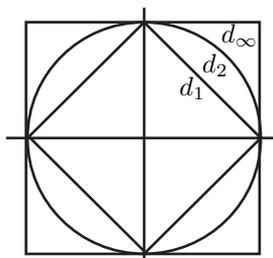
Además se cumplen las relaciones $d_\infty(x, y) \leq d_2(x, y) \leq d_1(x, y) \leq n d_\infty(x, y)$.

Con esto estamos casi a punto de llegar al concepto general de espacio topológico. La transición la haremos mediante el concepto siguiente:

Definición 2.6 Sea M un espacio métrico, $x \in M$ y $\epsilon > 0$ (en estos casos sobrentenderemos $\epsilon \in \mathbb{R}$). Definimos

$$B_\epsilon(x) = \{y \in M \mid d(x, y) < \epsilon\} \quad (\text{Bola abierta de centro } x \text{ y radio } \epsilon).$$

$$B'_\epsilon(x) = \{y \in M \mid d(x, y) \leq \epsilon\} \quad (\text{Bola cerrada de centro } x \text{ y radio } \epsilon).$$



La figura muestra las bolas de centro $(0, 0)$ y radio 1 para las tres métricas que hemos definido en \mathbb{R}^2 . Las bolas con otros centros son trasladadas de éstas, y las bolas de otros radios tienen el mismo aspecto, pero de mayor o menor tamaño.

Las bolas abiertas se diferencian de las cerradas en que las primeras no contienen los puntos del borde. El interés de las bolas reside en que una bola de centro un punto P contiene todos los puntos de alrededor de P , por pequeño que sea su radio.

Observemos que el concepto de “puntos de alrededor” es un tanto escurridizo: Según lo que acabamos de decir, ningún punto en particular (distinto de P) está alrededor de P , pues siempre podemos tomar una bola suficientemente pequeña como para que deje fuera a dicho punto. Esto significa que no podemos dar sentido a la afirmación “ Q es un punto de alrededor de P ”, pero lo importante es que sí tiene sentido decir “El conjunto A contiene a todos los puntos de alrededor de P ”. Esto sucede cuando A contiene una bola cualquiera de centro P , y entonces diremos que A es un entorno de P . Aunque el concepto de entorno podría tomarse como concepto topológico básico, lo cierto es que es más cómodo partir de un concepto “más regular”: diremos que un conjunto es abierto si es un entorno de todos sus puntos. Los conjuntos abiertos tienen las propiedades que recoge la definición siguiente:

Definición 2.7 Una *topología* en un conjunto X es una familia \mathcal{T} de subconjuntos de X a cuyos elementos llamaremos *abiertos*, tal que cumpla las propiedades siguientes:

- \emptyset y X son abiertos.
- La unión de cualquier familia de abiertos es un abierto.
- La intersección de dos abiertos es un abierto.

Un *espacio topológico* es un par (X, \mathcal{T}) , donde X es un conjunto y \mathcal{T} es una topología en X . En la práctica escribiremos simplemente X en lugar de (X, \mathcal{T}) .

Sea M un espacio métrico. Diremos que un conjunto $G \subset M$ es *abierto* si para todo $x \in G$ existe un $\epsilon > 0$ tal que $B_\epsilon(x) \subset G$. Es inmediato comprobar que los conjuntos abiertos así definidos forman una topología en M , a la que llamaremos *topología inducida por la métrica*. En lo sucesivo consideraremos siempre a los espacios métricos como espacios topológicos con esta topología.

En el párrafo previo a la definición de topología hemos definido “abierto” como un conjunto que es entorno de todos sus puntos. Puesto que formalmente hemos definido los espacios topológicos a partir del concepto de abierto, ahora hemos de definir el concepto de entorno.

Si X es un espacio topológico, $U \subset X$ y $x \in U$, diremos que U es un *entorno* de x si existe un abierto G tal que $x \in G \subset U$.

Es inmediato comprobar que, en un espacio métrico, U es entorno de x si y sólo si existe un $\epsilon > 0$ tal que $B_\epsilon(x) \subset U$, es decir, si y sólo si U contiene a todos los puntos de alrededor de x , tal y como habíamos afirmado.

Ejemplo El intervalo $I = [0, 1]$, visto como subconjunto de \mathbb{R} , es entorno de todos sus puntos excepto de sus extremos 0 y 1, pues si $0 < x < 1$ siempre podemos tomar $\epsilon = \min\{x, 1 - x\}$ y entonces $B_\epsilon(x) =]x - \epsilon, x + \epsilon[\subset I$. En cambio, I no contiene todos los puntos de alrededor de 1, pues toda bola de centro 1 contiene puntos a la derecha de 1 y ninguno de ellos está en I . El caso del 0 es similar. En particular I no es abierto. ■

El teorema siguiente recoge las propiedades básicas de los entornos. La prueba es inmediata.

Teorema 2.8 *Sea X un espacio topológico, $x \in X$ y E_x la familia de todos los entornos de x .*

- a) *Un conjunto $G \subset X$ es abierto si y sólo si es un entorno de todos sus puntos.*
- b1) *$X \in E_x$.*
- b2) *Si $U \in E_x$ y $U \subset V \subset X$ entonces $V \in E_x$.*
- b3) *Si $U, V \in E_x$ entonces $U \cap V \in E_x$.*

Puesto que los abiertos pueden definirse a partir de los entornos, es obvio que si dos topologías sobre un mismo conjunto tienen los mismos entornos entonces son iguales. Las desigualdades del teorema 2.5 implican que una bola para una de las tres distancias definidas en el producto M contiene otra bola del mismo centro para cualquiera de las otras distancias. De aquí se sigue que un subconjunto de M es entorno de un punto para una distancia si y sólo si lo es para las demás, y de aquí a su vez que las tres distancias definen la misma topología en el producto. En particular, si \mathbb{K} es un cuerpo métrico, las tres distancias que tenemos definidas sobre \mathbb{K}^n definen la misma topología, a la que llamaremos *topología usual* o *topología euclídea* en \mathbb{K}^n .

Ésta es una primera muestra del carácter auxiliar de las distancias en topología. Cuando queramos probar un resultado puramente topológico sobre \mathbb{R}^n podremos apoyarnos en la distancia que resulte más conveniente, sin que ello suponga una pérdida de generalidad. La distancia d_2 es la distancia euclídea y por lo tanto la más natural desde un punto de vista geométrico, pero las distancias d_1 y d_∞ son formalmente más sencillas y a menudo resultan más adecuadas.

Ejemplo Es fácil definir una distancia en \mathbb{R}^n que induzca una topología distinta de la usual. De hecho, si X es un conjunto cualquiera podemos considerar la distancia $d : X \times X \rightarrow \mathbb{R}$ dada por

$$d(x, y) = \begin{cases} 1 & \text{si } x \neq y, \\ 0 & \text{si } x = y \end{cases}$$

Es fácil ver que efectivamente es una distancia y para todo punto x se cumple que $B_1(x) = \{x\}$, luego $\{x\}$ es un entorno de x , luego es un abierto y, como toda unión de abiertos es abierta, de hecho todo subconjunto de X es abierto. La métrica d recibe el nombre de *métrica discreta* y la topología que induce es la *topología discreta*. Un espacio topológico cuya topología sea la discreta es un *espacio discreto*.

En un espacio discreto un punto no tiene más puntos a su alrededor que él mismo. Esta topología es la más adecuada para conjuntos como \mathbb{N} o \mathbb{Z} , pues, efectivamente, un número entero no tiene alrededor a ningún otro. ■

Las bolas abiertas de un espacio métrico son abiertas. Esto es fácil de ver intuitivamente, pero el mero hecho de que las hayamos llamado así no justifica que lo sean:

Teorema 2.9 *Las bolas abiertas de un espacio métrico son conjuntos abiertos.*

DEMOSTRACIÓN: Sea $B_\epsilon(x)$ una bola abierta y sea $y \in B_\epsilon(x)$. Entonces $d(x, y) < \epsilon$. Sea $0 < \delta < \epsilon - d(x, y)$. Basta probar que $B_\delta(y) \subset B_\epsilon(x)$. Ahora bien, si $z \in B_\delta(y)$, entonces $d(z, x) \leq d(z, y) + d(y, x) < \delta + d(x, y) < \epsilon$, luego en efecto $z \in B_\epsilon(x)$. ■

2.2 Bases y subbases

Hemos visto que la topología en un espacio métrico se define a partir de las bolas abiertas. El concepto de “bola abierta” no tiene sentido en un espacio topológico arbitrario en el que no tengamos dada una distancia, sin embargo hay otras familias de conjuntos que pueden representar un papel similar.

Definición 2.10 Sea X un espacio topológico. Diremos que una familia \mathcal{B} de abiertos de X (a los que llamaremos *abiertos básicos*) es una *base* de X si para todo abierto G de X y todo punto $x \in G$ existe un abierto $B \in \mathcal{B}$ tal que $x \in B \subset G$.

Si $x \in X$ diremos que una familia E de entornos (abiertos) de x (a los que llamaremos *entornos básicos* de x) es una *base de entornos* (abiertos) de x si todo entorno de x contiene un elemento de E .

En estos términos la propia definición de los abiertos métricos (junto con el hecho de que las bolas abiertas son realmente conjuntos abiertos) prueba que las bolas abiertas son una base de la topología métrica, y también es claro que

las bolas abiertas de centro un punto x forman una base de entornos abiertos de x . Pero estos conceptos son mucho más generales. Pensemos por ejemplo que otras bases de un espacio métrico son las bolas abiertas de radio menor que 1, las bolas abiertas de radio racional, etc. Cualquier base determina completamente la topología y en cada ocasión puede convenir trabajar con una base distinta.

Teorema 2.11 *Sea X un espacio topológico.*

- a) *Una familia de abiertos \mathcal{B} es una base de X si y sólo si todo abierto de X es unión de abiertos de \mathcal{B} .*
- b) *Si \mathcal{B} es una base de X y $x \in X$ entonces $\mathcal{B}_x = \{B \in \mathcal{B} \mid x \in B\}$ es una base de entornos abiertos de x .*
- c) *Si para cada punto $x \in X$ el conjunto E_x es una base de entornos abiertos de x entonces $\mathcal{B} = \bigcup_{x \in X} E_x$ es una base de X .*

DEMOSTRACIÓN: a) Si \mathcal{B} es una base de X y G es un abierto es claro que G es la unión de todos los abiertos de \mathcal{B} contenidos en G , pues una inclusión es obvia y si $x \in G$ existe un $B \in \mathcal{B}$ tal que $x \in B \subset G$, luego x está en la unión considerada. El recíproco es obvio.

b) Los elementos de \mathcal{B}_x son obviamente entornos de x y si U es un entorno de x entonces existe un abierto G tal que $x \in G \subset U$, y a su vez existe $B \in \mathcal{B}$ tal que $x \in B \subset G$, luego $B \in \mathcal{B}_x$ y $B \subset U$. Esto prueba que \mathcal{B}_x es una base de entornos abiertos de x .

c) Si G es un abierto de X y $x \in G$, entonces G es un entorno de x , luego existe un entorno básico $B \in E_x$ tal que $x \in B \subset G$ y ciertamente $B \in \mathcal{B}$. Como además los elementos de \mathcal{B} son abiertos, tenemos que \mathcal{B} es una base de X . ■

Una forma habitual de definir una topología en un conjunto es especificar una base o una base de entornos abiertos de cada punto. Por ejemplo, la topología métrica puede definirse como la topología que tiene por base a las bolas abiertas o como base de entornos de cada punto x a las bolas abiertas de centro x . No obstante, para que una familia de conjuntos pueda ser base de una topología ha de cumplir unas propiedades muy simples que es necesario comprobar. El teorema siguiente da cuenta de ellas.

Teorema 2.12 *Sea X un conjunto y \mathcal{B} una familia de subconjuntos de X que cumpla las propiedades siguientes:*

- a) $X = \bigcup_{B \in \mathcal{B}} B$,
- b) *Si $U, V \in \mathcal{B}$ y $x \in U \cap V$ entonces existe $W \in \mathcal{B}$ tal que $x \in W \subset U \cap V$.*

Entonces existe una única topología en X para la cual \mathcal{B} es una base.

DEMOSTRACIÓN: Definimos los abiertos de X como las uniones de elementos de \mathcal{B} . Basta comprobar que estos abiertos forman realmente una topología, pues ciertamente en tal caso \mathcal{B} será una base y la topología será única.

El conjunto vacío es abierto trivialmente (o si se prefiere, por definición). El conjunto X es abierto por la propiedad a).

La unión de abiertos es obviamente abierta (una unión de uniones de elementos de \mathcal{B} es al fin y al cabo una unión de elementos de \mathcal{B}).

Sean G_1 y G_2 abiertos y supongamos que $x \in G_1 \cap G_2$. Como G_1 es unión de elementos de \mathcal{B} existe un $U \in \mathcal{B}$ tal que $x \in U \subset G_1$. Similarmente $x \in V \subset G_2$ con $V \in \mathcal{B}$. Por la propiedad b) existe $W \in \mathcal{B}$ tal que $x \in W \subset U \cap V \subset G_1 \cap G_2$. Así pues, x está en la unión de los conjuntos $W \in \mathcal{B}$ tales que $W \subset G_1 \cap G_2$, y la otra inclusión es obvia, luego $G_1 \cap G_2$ es unión de elementos de \mathcal{B} . ■

El teorema siguiente nos da las condiciones que hemos de comprobar para definir una topología a partir de una familia de bases de entornos abiertos.

Teorema 2.13 *Sea X un conjunto y para cada $x \in X$ sea \mathcal{B}_x una familia no vacía de subconjuntos de X tal que:*

- a) *Si $U \in \mathcal{B}_x$, entonces $x \in U$.*
- b) *Si $U, V \in \mathcal{B}_x$, existe un $W \in \mathcal{B}_x$ tal que $W \subset U \cap V$.*
- c) *Si $x \in U \in \mathcal{B}_y$, existe un $V \in \mathcal{B}_x$ tal que $V \subset U$.*

Entonces existe una única topología para la cual cada \mathcal{B}_x es una base de entornos abiertos de x .

DEMOSTRACIÓN: Sea $\mathcal{B} = \bigcup_{x \in X} \mathcal{B}_x$. Veamos que \mathcal{B} cumple las condiciones del teorema anterior para ser base de una topología en X . Por la condición a) tenemos que $X = \bigcup_{B \in \mathcal{B}} B$.

Si $U, V \in \mathcal{B}$ y $x \in U \cap V$, entonces $U \in \mathcal{B}_y$ y $V \in \mathcal{B}_z$ para ciertos y, z . Existen $U', V' \in \mathcal{B}_x$ tales que $U' \subset U$ y $V' \subset V$ (por la condición c). Existe $W \in \mathcal{B}_x$ tal que $W \subset U' \cap V'$ (por la condición b). Así $x \in W \subset U \cap V$ con $W \in \mathcal{B}$.

Por lo tanto \mathcal{B} es la base de una topología en X para la que los elementos de cada \mathcal{B}_x son abiertos y, en particular, entornos de x . Si A es un entorno de x para dicha topología, existe un $U \in \mathcal{B}$ tal que $x \in U \subset A$. Por definición de \mathcal{B} , existe un $y \in X$ tal que $U \in \mathcal{B}_y$, y por c) existe un $V \in \mathcal{B}_x$ tal que $V \subset U \subset A$. Esto prueba que \mathcal{B}_x es una base de entornos de x .

Las bases de entornos determinan los entornos y por tanto la topología, es decir, se da la unicidad. ■

Ejemplo Consideremos dos conjuntos cualesquiera con la única condición de que no sean números reales, a los que llamaremos $+\infty$ y $-\infty$. Vamos a usar el teorema anterior para convertir en espacio topológico a $\overline{\mathbb{R}} = \mathbb{R} \cup \{-\infty, +\infty\}$. En primer lugar extendemos a $\overline{\mathbb{R}}$ el orden usual de \mathbb{R} estableciendo que $+\infty$ sea el máximo de $\overline{\mathbb{R}}$ y $-\infty$ sea el mínimo.

La idea es que queremos que “los puntos de alrededor” de $+\infty$ sean los números reales muy grandes, y “los puntos de alrededor” de $-\infty$ sean los números reales muy pequeños (muy grandes en valor absoluto, pero negativos).

Esto se concreta tomando como base de entornos abiertos de cada número real x al conjunto de los entornos abiertos de x en \mathbb{R} con la topología usual, la base de entornos abiertos de $+\infty$ está formada por los intervalos $]x, +\infty]$, donde x varía en \mathbb{R} , y la base de entornos abiertos de $-\infty$ la forman los intervalos $[-\infty, x[$, donde x varía en \mathbb{R} . Con esto estamos diciendo que un conjunto contiene a los alrededores de $+\infty$ si contiene a todos los números reales a partir de uno dado, y análogamente para $-\infty$.

Es fácil comprobar que las familias consideradas cumplen las propiedades del teorema anterior, luego definen una topología en \mathbb{R} .

Teniendo en cuenta que hemos definido los entornos abiertos de los números reales como los entornos abiertos que ya tienen en la topología usual, es inmediato que un subconjunto de \mathbb{R} es abierto en la topología usual de \mathbb{R} si y sólo si lo es en la topología que hemos definido en \mathbb{R} . ■

Hay un concepto análogo a los de base y base de entornos que es menos intuitivo, pero mucho más práctico a la hora de definir topologías. Se trata del concepto de subbase:

Definición 2.14 Sea X un espacio topológico. Una familia de abiertos S es una *subbase* de X si las intersecciones finitas de elementos de S forman una base de X .

Por ejemplo, es fácil ver que los intervalos abiertos $]a, b[$ forman una base de \mathbb{R} (son las bolas abiertas). Por consiguiente, los intervalos de la forma $] -\infty, a[$ y $]a, +\infty[$ forman una subbase de \mathbb{R} , pues son abiertos y entre sus intersecciones finitas se encuentran todos los intervalos $]a, b[$ (notar además que cualquier familia de abiertos que contenga a una base es una base).

La ventaja de las subbases consiste en que una familia no ha de cumplir ninguna propiedad en especial para ser subbase de una topología:

Teorema 2.15 Sea X un conjunto y S una familia de subconjuntos de X . Entonces existe una única topología en X de la cual S es subbase.

DEMOSTRACIÓN: Sea \mathcal{B} la familia de las intersecciones finitas de elementos de S . Entonces $X = \bigcap_{G \in \mathcal{B}} G$ y obviamente la intersección de dos intersecciones finitas de elementos de S es una intersección finita de elementos de S ; luego si $U, V \in \mathcal{B}$ también $U \cap V \in \mathcal{B}$, de donde se sigue que \mathcal{B} es la base de una topología en X , de la cual S es subbase. Claramente es única, pues \mathcal{B} es base de cualquier topología de la que S sea subbase. ■

2.3 Productos y subespacios

Hemos visto que el producto de una familia finita de espacios métricos es de nuevo un espacio métrico de forma natural (o mejor dicho, de tres formas distintas pero equivalentes desde un punto de vista topológico). Ahora veremos

que la topología del producto se puede definir directamente a partir de las topologías de los factores sin necesidad de considerar las distancias. Más aún, podemos definir el producto de cualquier familia de espacios topológicos, no necesariamente finita.

Definición 2.16 Sean $\{X_i\}_{i \in I}$ espacios topológicos. Consideremos su producto cartesiano $X = \prod_{i \in I} X_i$ y las proyecciones $p_i : X \rightarrow X_i$ que asignan a cada punto su coordenada i -ésima. Llamaremos *topología producto* en X a la que tiene por subbase a los conjuntos $p_i^{-1}[G]$, donde $i \in I$ y G es abierto en X_i .

Una base de la topología producto la forman los conjuntos de la forma $\bigcap_{i \in F} p_i^{-1}[G_i]$, donde F es un subconjunto finito de I y G_i es abierto en X_i .

Equivalentemente, la base está formada por los conjuntos $\prod_{i \in I} G_i$, donde cada G_i es abierto en X_i y $G_i = X_i$ salvo para un número finito de índices. Al conjunto de estos índices se le llama *soporte* del abierto básico $\prod_{i \in I} G_i$.

Si el número de factores es finito la restricción se vuelve vacía, de modo que un abierto básico en un producto $X_1 \times \cdots \times X_n$ es simplemente un conjunto de la forma $G_1 \times \cdots \times G_n$, donde cada G_i es abierto en X_i .

En lo sucesivo “casi todo i ” querrá decir “todo índice i salvo un número finito de ellos”.

Teorema 2.17 Sean $\{X_i\}_{i \in I}$ espacios topológicos, para cada i sea \mathcal{B}_i una base de X_i . Entonces los conjuntos de la forma $\prod_{i \in I} G_i$, donde cada G_i está en \mathcal{B}_i o es X_i (y casi todos son X_i) forman una base de $\prod_{i \in I} X_i$.

DEMOSTRACIÓN: Consideremos la topología \mathcal{T} en el producto que tiene por subbase a los conjuntos $p_i^{-1}[G_i]$ con G_i en \mathcal{B}_i (y, por consiguiente, tienen por base a los abiertos del enunciado). Como ciertamente estos conjuntos son abiertos para la topología producto, tenemos que todo abierto de \mathcal{T} lo es de la topología producto. Recíprocamente, un abierto subbásico de la topología producto es $p_i^{-1}[G_i]$, con G_i abierto en X_i . Entonces $G_i = \bigcup_{B \in A_i} B$, donde cada A_i es un subconjunto de \mathcal{B}_i . Por lo tanto $p_i^{-1}[G_i] = \bigcup_{B \in A_i} p_i^{-1}[B]$ es abierto de \mathcal{T} . Por consiguiente todo abierto de la topología producto lo es de \mathcal{T} y así ambas topologías coinciden. ■

En lo sucesivo, a pesar de que en el producto se puedan considerar otras bases, cuando digamos “abiertos básicos” nos referiremos a los abiertos indicados en el teorema anterior tomando como bases de los factores las propias topologías salvo que se esté considerando alguna base en concreto.

Tal y como anunciábamos, el producto de espacios topológicos generaliza al producto de espacios métricos (o de espacios normados). El teorema siguiente lo prueba.

Teorema 2.18 Si M_1, \dots, M_n son espacios métricos, entonces la topología inducida por las métricas de 2.5 en $M = M_1 \times \dots \times M_n$ es la topología producto.

DEMOSTRACIÓN: Como las tres métricas inducen la misma topología sólo es necesario considerar una de ellas, pero para la métrica d_∞ se cumple $B_\epsilon(x) = B_\epsilon(x_1) \times \dots \times B_\epsilon(x_n)$, luego la base inducida por la métrica es base de la topología producto. ■

La definición de topología producto es sin duda razonable para un número finito de factores. Sin embargo cuando tenemos infinitos factores hemos exigido una condición de finitud que no hemos justificado. En principio podríamos considerar en $\prod_{i \in I} X_i$ la topología que tiene por base a los productos $\prod_{i \in I} G_i$ con G_i abierto en X_i (sin ninguna restricción de finitud). Ciertamente estos conjuntos son base de una topología a la que se le llama *topología de cajas*, y el teorema siguiente muestra que no coincide con la topología producto que hemos definido. La topología producto resulta ser mucho más útil que la topología de cajas.

Teorema 2.19 Sea $\{X_i\}_{i \in I}$ una familia de espacios topológicos. Los únicos abiertos en $\prod_{i \in I} X_i$ de la forma $\prod_{i \in I} G_i \neq \emptyset$ son los abiertos básicos, es decir, los que además cumplen que cada G_i es abierto y $G_i = X_i$ para casi todo i .

DEMOSTRACIÓN: Supongamos que $\prod_{i \in I} G_i$ es un abierto no vacío. Consideremos un punto $x \in \prod_{i \in I} G_i$. Existirá un abierto básico $\prod_{i \in I} H_i$ tal que $x \in \prod_{i \in I} H_i \subset \prod_{i \in I} G_i$, luego para cada índice i se cumplirá $x_i \in H_i \subset G_i$, y como casi todo H_i es igual a X_i , tenemos que $G_i = X_i$ para casi todo i . Además tenemos que G_i es un entorno de x_i , pero dado cualquier elemento $a \in G_i$ siempre podemos formar un $x \in \prod_{i \in I} G_i$ tal que $x_i = a$, luego en realidad tenemos que G_i es un entorno de todos sus puntos, o sea, es abierto. ■

Nos ocupamos ahora de los subespacios de un espacio topológico. Es evidente que todo subconjunto N de un espacio métrico M es también un espacio métrico con la misma distancia restringida a $N \times N$. Por lo tanto tenemos una topología en M y otra en N . Vamos a ver que podemos obtener la topología de N directamente a partir de la de M , sin pasar por la métrica.

Teorema 2.20 Sea X un espacio topológico (con topología \mathcal{T}) y $A \subset X$. Definimos $\mathcal{T}_A = \{G \cap A \mid G \in \mathcal{T}\}$. Entonces \mathcal{T}_A es una topología en A llamada topología relativa a X (o topología inducida por X) en A . En lo sucesivo sobrentenderemos siempre que la topología de un subconjunto de un espacio X es la topología relativa.

DEMOSTRACIÓN: $A = X \cap A \in \mathcal{T}_A$, $\emptyset = \emptyset \cap A \in \mathcal{T}_A$.

Sea $C \subset \mathcal{T}_A$. Para cada $G \in C$ sea $U_G = \{U \in \mathcal{T} \mid U \cap A = G\} \neq \emptyset$ y sea V_G la unión de todos los abiertos de U_G .

De este modo V_G es un abierto en X y $V_G \cap A = G$.

$$\bigcup_{G \in \mathcal{C}} G = \bigcup_{G \in \mathcal{C}} V_G \cap A, \quad \text{y} \quad \bigcup_{G \in \mathcal{C}} V_G \in \mathcal{T}, \quad \text{luego} \quad \bigcup_{G \in \mathcal{C}} G \in \mathcal{T}_A.$$

Si $U, V \in \mathcal{T}_A$, $U = U' \cap A$ y $V = V' \cap A$ con $U', V' \in \mathcal{T}$. Entonces $U \cap V = U' \cap V' \cap A \in \mathcal{T}_A$, pues $U' \cap V' \in \mathcal{T}$. Así \mathcal{T}_A es una topología en A . ■

Ejemplo Consideremos $I = [0, 1] \subset \mathbb{R}$. Resulta que $]1/2, 1]$ es abierto en I , pues $]1/2, 1] =]1/2, 2[\cap I$ y $]1/2, 2[$ es abierto en \mathbb{R} . Sin embargo $]1/2, 1]$ no es abierto en \mathbb{R} porque no es entorno de 1. Intuitivamente, $]1/2, 1]$ no contiene a todos los puntos de alrededor de 1 en \mathbb{R} (faltan los que están a la derecha de 1), pero sí contiene a todos los puntos de alrededor de 1 en I . ■

La relación entre espacios y subespacios viene perfilada por los teoremas siguientes. El primero garantiza que la topología relativa no depende del espacio desde el que relativicemos.

Teorema 2.21 *Si X es un espacio topológico (con topología \mathcal{T}) y $A \subset B \subset X$, entonces $\mathcal{T}_A = (\mathcal{T}_B)_A$.*

DEMOSTRACIÓN: Si U es abierto en \mathcal{T}_A , entonces $U = V \cap A$ con $V \in \mathcal{T}$, luego $V \cap B \in \mathcal{T}_B$ y $U = V \cap A = (V \cap B) \cap A \in (\mathcal{T}_B)_A$.

Si $U \in (\mathcal{T}_B)_A$, entonces $U = V \cap A$ con $V \in \mathcal{T}_B$, luego $V = W \cap B$ con $W \in \mathcal{T}$. Así pues, $U = W \cap B \cap A = W \cap A \in \mathcal{T}_A$. Por lo tanto $\mathcal{T}_A = (\mathcal{T}_B)_A$. ■

Teorema 2.22 *Si \mathcal{B} es una base de un espacio X y $A \subset X$, entonces el conjunto $\{B \cap A \mid B \in \mathcal{B}\}$ es una base de A .*

DEMOSTRACIÓN: Sea U un abierto en A y $x \in U$. Existe un V abierto en X tal que $U = V \cap A$. Existe un $B \in \mathcal{B}$ tal que $x \in B \subset V$, luego $x \in B \cap A \subset V \cap A = U$. Por lo tanto la familia referida es base de A . ■

Similarmente se demuestra:

Teorema 2.23 *Si \mathcal{B}_x es una base de entornos (abiertos) de un punto x de un espacio X y $x \in A \subset X$, entonces $\{B \cap A \mid B \in \mathcal{B}_x\}$ es una base de entornos (abiertos) de x en A .*

Teorema 2.24 *Sea M un espacio métrico y sea $A \subset M$. Entonces $d' = d|_{A \times A}$ es una distancia en A y la topología que induce es la topología relativa.*

DEMOSTRACIÓN: Una base para la topología inducida por la métrica de A sería la formada por las bolas

$$B_\epsilon^{d'}(x) = \{a \in A \mid d'(x, a) < \epsilon\} = \{a \in M \mid d(x, a) < \epsilon\} \cap A = B_\epsilon^d(x) \cap A,$$

pero éstas son una base para la topología relativa por el teorema 2.22. ■

Teorema 2.25 Sea $\{X_i\}_{i \in I}$ una familia de espacios topológicos y para cada i sea $Y_i \subset X_i$. Entonces la topología inducida en $\prod_{i \in I} Y_i$ por $\prod_{i \in I} X_i$ es la misma que la topología producto de los $\{Y_i\}_{i \in I}$.

(La base obtenida por el teorema 2.22 a partir de la base usual de la topología producto es claramente la base usual de la topología producto.)

Ejercicio: Probar que la topología que hemos definido en \mathbb{R} induce en \mathbb{R} la topología euclídea.

2.4 Algunos conceptos topológicos

Dedicamos esta sección a desarrollar el lenguaje topológico, es decir, a introducir las características de un espacio y sus subconjuntos que pueden definirse a partir de su topología. Hasta ahora hemos visto únicamente los conceptos de abierto y entorno. Otro concepto importante es el dual conjuntista de “abierto”:

Definición 2.26 Diremos que un subconjunto de un espacio topológico es *cerrado* si su complementario es abierto.

Por ejemplo, un semiplano (sin su recta frontera) es un conjunto abierto, mientras que un semiplano con su frontera es cerrado, pues su complementario es el semiplano opuesto sin su borde, luego es abierto. Pronto veremos que la diferencia entre los conjuntos abiertos y los cerrados es precisamente que los primeros no contienen a los puntos de su borde y los segundos contienen todos los puntos de su borde. Es importante notar que un conjunto no tiene por qué ser ni abierto ni cerrado. Baste pensar en el intervalo $[0, 1[$.

Ejercicio: Sea $X = [0, 1] \cup]3, 4]$. Probar que $]3, 4]$ es a la vez abierto y cerrado en X .

Las propiedades de los cerrados se deducen inmediatamente de las de los abiertos:

Teorema 2.27 Sea X un espacio topológico. Entonces:

- \emptyset y X son cerrados.
- La intersección de cualquier familia de cerrados es un cerrado.
- La unión de dos cerrados es un cerrado.

Puesto que la unión de abiertos es abierta, al unir todos los abiertos contenidos en un conjunto dado obtenemos el mayor abierto contenido en él. Similarmente, al intersecar todos los cerrados que contienen a un conjunto dado obtenemos el menor cerrado que lo contiene:

Definición 2.28 Sea X un espacio topológico. Llamaremos *interior* de un conjunto $A \subset X$ al mayor abierto contenido en A . Lo representaremos por $\text{int } A$ o $\overset{\circ}{A}$. Llamaremos *clausura* de A al menor cerrado que contiene a A . Lo representaremos por $\text{cl } A$ o \bar{A} . Los puntos de $\overset{\circ}{A}$ se llaman *puntos interiores* de A , mientras que los de \bar{A} se llaman *puntos adherentes* de A .

Así pues, para todo conjunto A tenemos que $\overset{\circ}{A} \subset A \subset \overline{A}$. El concepto de punto interior es claro: un punto x es interior a un conjunto A si y sólo si A es un entorno de x . Por ejemplo, en un semiplano cerrado, los puntos interiores son los que no están en el borde. El teorema siguiente nos caracteriza los puntos adherentes.

Teorema 2.29 *Sea X un espacio topológico y A un subconjunto de X . Un punto x es adherente a A si y sólo si todo entorno de x corta a A .*

DEMOSTRACIÓN: Supongamos que x es adherente a A . Sea U un entorno de x . Existe un abierto G tal que $x \in G \subset U$. Basta probar que $G \cap A \neq \emptyset$. Ahora bien, en caso contrario $X \setminus G$ sería un cerrado que contiene a A , luego $\overline{A} \subset X \setminus G$, mientras que $x \in \overline{A} \cap G$.

Recíprocamente, si x tiene esta propiedad entonces $x \in \overline{A}$, ya que de lo contrario $X \setminus \overline{A}$ sería un entorno de x que no corta a A . ■

Vemos, pues, que, como su nombre indica, los puntos adherentes a un conjunto A son los que “están pegados” a A , en el sentido de que tienen alrededor puntos de A . Por ejemplo, es fácil ver que los puntos adherentes a un semiplano abierto son sus propios puntos más los de su borde. Veamos ahora que el concepto de “borde” corresponde a una noción topológica general:

Definición 2.30 *Sea X un espacio topológico y $A \subset X$. Llamaremos *frontera* de A al conjunto $\partial A = \overline{A} \cap \overline{X \setminus A}$.*

Así, los puntos frontera de un conjunto son aquellos que tienen alrededor puntos que están en A y puntos que no están en A . Esto es claramente una definición general del “borde” de un conjunto. Por ejemplo, la frontera de un triángulo la forman los puntos de sus lados.

Teorema 2.31 *Sea X un espacio topológico. Se cumple:*

- a) *Si $A \subset X$ entonces $\overset{\circ}{A} \subset A \subset \overline{A}$, además $\overset{\circ}{A}$ es abierto y \overline{A} es cerrado.*
- b) *Si $A \subset B \subset X$ y A es abierto entonces $A \subset \overset{\circ}{B}$.*
- c) *Si $A \subset B \subset X$ y B es cerrado entonces $\overline{A} \subset B$.*
- d) *Si $A \subset B \subset X$ entonces $\overset{\circ}{A} \subset \overset{\circ}{B}$ y $\overline{A} \subset \overline{B}$.*
- e) *Si $A, B \subset X$, entonces $\text{int}(A \cap B) = \text{int} A \cap \text{int} B$, $\overline{A \cup B} = \overline{A} \cup \overline{B}$.*
- f) *$A \subset X$ es abierto si y sólo si $A = \overset{\circ}{A}$, y es cerrado si y sólo si $A = \overline{A}$.*
- g) *Si $A \subset B \subset X$, entonces $\overline{A}^B = \overline{A}^X \cap B$.*
- h) *Si $A \subset X$, entonces $\text{int}(X \setminus A) = X \setminus \text{cl} A$ y $\text{cl}(X \setminus A) = X \setminus \text{int} A$.*

DEMOSTRACIÓN: Muchas de estas propiedades son inmediatas. Probaremos sólo algunas.

e) Claramente $A \cup B \subset \overline{A \cup B}$, y el segundo conjunto es cerrado, luego $\overline{A \cup B} \subset \overline{A \cup B}$. Por otra parte es claro que $\overline{A} \subset \overline{A \cup B}$ y $\overline{B} \subset \overline{A \cup B}$, luego tenemos la otra inclusión. La prueba con interiores es idéntica.

g) Observemos en primer lugar que los cerrados de B son exactamente las intersecciones con B de los cerrados de X . En efecto, si C es cerrado en X entonces $X \setminus C$ es abierto en X , luego $B \cap (X \setminus C) = B \setminus C$ es abierto en B , luego $B \setminus (B \setminus C) = B \cap C$ es cerrado en B . El recíproco es similar.

Por definición, \overline{A}^X es la intersección de todos los cerrados en X que contienen a A , luego $\overline{A}^X \cap B$ es la intersección de todas las intersecciones con B de los cerrados en X que contienen a A , pero estos son precisamente los cerrados de B que contienen a A , o sea, $\overline{A}^X \cap B$ es exactamente \overline{A}^B .

h) Tenemos que $A \subset \text{cl} A$, luego $X \setminus \text{cl} A \subset X \setminus A$ y el primero es abierto, luego $X \setminus \text{cl} A \subset \text{int}(X \setminus A)$.

Por otra parte $\text{int}(X \setminus A) \subset X \setminus A$, luego $A \subset X \setminus \text{int}(X \setminus A)$, y éste es cerrado, luego $\overline{A} \subset X \setminus \text{int}(X \setminus A)$ y $\text{int}(X \setminus A) \subset X \setminus \overline{A}$. ■

En la prueba de la propiedad g) hemos visto lo siguiente:

Teorema 2.32 Si X es un espacio topológico y $A \subset X$, los cerrados en la topología relativa de A son las intersecciones con A de los cerrados de X .

Conviene observar que el análogo a g) para interiores es falso. Es decir, no se cumple en general que $\overset{\circ}{A}^B = \overset{\circ}{A} \cap B$. Por ejemplo, es fácil ver que en \mathbb{R} se cumple $\overset{\circ}{\mathbb{N}} = \emptyset$, luego $\overset{\circ}{\mathbb{N}} \cap \mathbb{Z} = \emptyset$, mientras que $\overset{\circ}{\mathbb{N}}^{\mathbb{Z}} = \mathbb{N}$.

Ejercicio: Probar que en un espacio normado sobre \mathbb{R} la clausura de una bola abierta es la bola cerrada del mismo radio y el interior de una bola cerrada es la bola abierta. ¿Cuál es la frontera de ambas? Dar ejemplos que muestren la falsedad de estos hechos en un espacio métrico arbitrario.

Vamos a refinar el concepto de punto adherente. Hemos visto que los puntos adherentes a un conjunto A son aquellos que tienen alrededor puntos de A . Sucede entonces que todo punto $x \in A$ es trivialmente adherente, porque x es un punto de alrededor de x y está en A . Cuando eliminamos esta posibilidad trivial tenemos el concepto de punto de acumulación:

Definición 2.33 Sea X un espacio topológico y $A \subset X$. Diremos que un punto $x \in X$ es un *punto de acumulación* de A si todo entorno U de x cumple $(U \setminus \{x\}) \cap A \neq \emptyset$. El conjunto de puntos de acumulación de A se llama *conjunto derivado* de A y se representa por A' .

Ejemplo Consideremos el conjunto $A = \{1/n \mid n \in \mathbb{N} \setminus \{0\}\} \subset \mathbb{R}$. Es fácil ver que $\overline{A} = A \cup \{0\}$. Sin embargo, $A' = \{0\}$. En efecto, en general se cumple que $A' \subset \overline{A}$, pero ningún punto $1/n \in A$ es de acumulación, pues

$$\left] \frac{1}{n} - \frac{1}{n(n+1)}, \frac{1}{n} + \frac{1}{n(n+1)} \right[$$

es un entorno de $1/n$ que no corta a A salvo en este mismo punto. ■

Como ya hemos dicho, siempre es cierto que $A' \subset \bar{A}$. También es claro que un punto adherente que no esté en A ha de ser un punto de acumulación de A . En otras palabras, $\bar{A} = A \cup A'$. Los puntos de A pueden ser de acumulación o no serlo. Por ejemplo, todos los puntos de $[0, 1]$ son de acumulación, mientras que los puntos del ejemplo anterior no lo eran.

Definición 2.34 Sea X un espacio topológico y $A \subset X$. Los puntos de $A \setminus A'$ se llaman *puntos aislados* de A .

Un punto $x \in A$ es aislado si y sólo si tiene un entorno U tal que $U \cap A = \{x\}$. El entorno lo podemos tomar abierto, y entonces vemos que los puntos aislados de A son los puntos que son abiertos en la topología relativa. Vemos, pues, que un espacio es discreto si y sólo si todos sus puntos son aislados. Es el caso del ejemplo anterior.

Definición 2.35 Un subconjunto A de un espacio topológico X es *denso* si $\bar{A} = X$.

Aplicando la propiedad h) de 2.31 vemos que A es denso en X si y sólo si $X \setminus A$ tiene interior vacío, es decir, si y sólo si todo abierto de X corta a A . Esto significa que los puntos de A están “en todas partes”. Es claro que, en el caso de espacios métricos, los conjuntos densos son los mismos definidos en 1.19.

Por ejemplo, puesto que todo intervalo de números reales contiene números racionales e irracionales, es claro que \mathbb{Q} y $\mathbb{R} \setminus \mathbb{Q}$ son densos en \mathbb{R} . De aquí se sigue fácilmente que \mathbb{Q}^n y $(\mathbb{R} \setminus \mathbb{Q})^n$ son densos en \mathbb{R}^n .

Ejercicio: Probar que si A es abierto en un espacio X y D es denso en X entonces $A \cap D$ es denso en A .

Hay una propiedad que no cumplen todos los espacios topológicos, pero sí la práctica totalidad de espacios de interés.

Definición 2.36 Diremos que un espacio topológico X es un *espacio de Hausdorff* si para todo par de puntos distintos $x, y \in X$ existen abiertos disjuntos U y V tales que $x \in U, y \in V$ (se dice que los abiertos U y V *separan* a x e y).

Por ejemplo, si en un conjunto X con más de un punto consideramos la topología formada únicamente por los abiertos \emptyset y X (topología trivial) obtenemos un espacio que no es de Hausdorff. Se trata de un espacio patológico donde todo punto está alrededor de cualquier otro. Aunque la topología trivial es ciertamente la más patológica posible, lo cierto es que todas las topologías no de Hausdorff comparten con ella su patología, y rara vez resultan de interés. Veamos las propiedades de los espacios de Hausdorff:

Teorema 2.37 *Se cumplen las propiedades siguientes:*

- a) *En un espacio de Hausdorff, todo conjunto finito es cerrado.*
- b) *Todo espacio de Hausdorff finito es discreto.*

- c) Todo subespacio de un espacio de Hausdorff es un espacio de Hausdorff.
- d) El producto de una familia de espacios de Hausdorff es un espacio de Hausdorff.
- e) Todo espacio métrico es un espacio de Hausdorff.
- f) Un espacio X es de Hausdorff si y sólo si la diagonal $\Delta = \{(x, x) \mid x \in X\}$ es cerrada en $X \times X$.

DEMOSTRACIÓN: a) Basta probar que todo punto $\{x\}$ en un espacio de Hausdorff X es cerrado. Ahora bien, dado $y \in X \setminus \{x\}$, existen abiertos disjuntos U, V tales que $x \in U, y \in V$, luego $y \in V \subset X \setminus \{x\}$, lo que prueba que $X \setminus \{x\}$ es entorno de todos sus puntos, luego $\{x\}$ es cerrado.

b) En un espacio de Hausdorff finito todo subconjunto es cerrado, luego todo subconjunto es abierto, luego es discreto.

c) Si X es un espacio de Hausdorff y $A \subset X$, dados dos puntos $x, y \in A$, existen abiertos disjuntos U, V en X que separan a x e y , luego $U \cap A, V \cap A$ son abiertos disjuntos en A que separan a x e y .

d) Consideremos un producto de espacios de Hausdorff $\prod_{i \in I} X_i$ y dos de sus puntos x, y . Sea i_0 un índice tal que $x_{i_0} \neq y_{i_0}$. Existen abiertos U, V en X_{i_0} que separan a x_{i_0} e y_{i_0} . Entonces $p_{i_0}^{-1}[U]$ y $p_{i_0}^{-1}[V]$ son abiertos subbásicos disjuntos en el producto que separan a x e y .

e) Si X es un espacio métrico, dos de sus puntos x, y están separados por las bolas de centros x, y y radio $d(x, y)/2$.

f) La diagonal Δ es cerrada si y sólo si su complementario es abierto, si y sólo si para todo par $(x, y) \in X \times X$ con $x \neq y$ existe un abierto básico $U \times V$ en $X \times X$ tal que $(x, y) \in U \times V \subset X \times X \setminus \Delta$. Ahora bien, la condición $U \times V \subset X \times X \setminus \Delta$ equivale a $U \cap V = \emptyset$, luego la diagonal es cerrada si y sólo si X es Hausdorff. ■

Ejercicio: Probar que si un producto de espacios topológicos es un espacio de Hausdorff no vacío, entonces cada uno de los factores es un espacio de Hausdorff.

Terminamos la sección con algunas propiedades métricas, no topológicas, es decir, propiedades definidas a partir de la distancia en un espacio métrico y que no se pueden expresar en términos de su topología.

Definición 2.38 Un subconjunto A de un espacio métrico es *acotado* si existe un $M > 0$ tal que para todo par de puntos $x, y \in A$ se cumple $d(x, y) \leq M$. El *diámetro* de un conjunto acotado A es el supremo de las distancias $d(x, y)$ cuando (x, y) varía en $A \times A$.

Es fácil probar que todo subconjunto de un conjunto acotado es acotado, así como que la unión finita de conjuntos acotados está acotada. Sin embargo hemos de tener presente el hecho siguiente: dado un espacio métrico M , podemos

definir $d'(x, y) = \min\{1, d(x, y)\}$. Es fácil ver que d' es una distancia en M y las bolas de radio menor que 1 para d' coinciden con las bolas respecto a d . Como estas bolas forman una base de las respectivas topologías métricas, concluimos que ambas distancias definen la misma topología. Sin embargo, respecto a d' todos los conjuntos están acotados. Esto prueba que el concepto de acotación no es topológico.

Ejercicio: Calcular el diámetro de una bola abierta en \mathbb{R}^n y en un espacio con la métrica discreta.

Definición 2.39 Si M es un espacio métrico, $A \neq \emptyset$ un subconjunto de M y $x \in M$, definimos la distancia de x a A como $d(x, A) = \inf\{d(x, y) \mid y \in A\}$.

Es evidente que si $x \in A$ entonces $d(x, A) = 0$, pues entre las distancias cuyo ínfimo determinan $d(x, A)$ se encuentra $d(x, x) = 0$. Sin embargo los puntos que cumplen $d(x, A) = 0$ no están necesariamente en A .

Teorema 2.40 Si M es un espacio métrico y $A \subset M$, entonces un punto x cumple $d(x, A) = 0$ si y sólo si x es adherente a A .

DEMOSTRACIÓN: Si $d(x, A) = 0$, para probar que es adherente basta ver que toda bola abierta de centro x corta a A . Dado $\epsilon > 0$ tenemos que $d(x, A) < \epsilon$, lo que significa que existe un $y \in A$ tal que $d(x, y) < \epsilon$, es decir, $y \in B_\epsilon(x) \cap A$. El recíproco se prueba igualmente. ■

En el caso de espacios normados podemos hacer algunas afirmaciones adicionales. La prueba del teorema siguiente es inmediata.

Teorema 2.41 Sea E un espacio normado y $A \subset E$. Las afirmaciones siguientes son equivalentes:

- a) A es acotado.
- b) Existe un $M > 0$ tal que $\|x\| \leq M$ para todo $x \in A$.
- c) Existe un $M > 0$ tal que $A \subset B_M(0)$.

Adición de un punto infinito a un espacio métrico Hemos visto cómo dotar de estructura de espacio topológico a $\mathbb{R} \cup \{\pm\infty\}$ de modo que los puntos de “alrededor” de $+\infty$ sean los números positivos grandes, y los puntos de “alrededor” de $-\infty$ sean los números negativos grandes. Esto depende de la estructura de orden de \mathbb{R} , pero si renunciamos a distinguir entre $+\infty$ y $-\infty$ podemos hacer algo similar en un espacio métrico cualquiera.

En efecto, si M es un espacio métrico, tomamos un conjunto cualquiera $\infty \notin M$ y definimos $M^\infty = M \cup \{\infty\}$. Dotamos a M^∞ de estructura de espacio topológico tomando como abiertos los abiertos de M más los conjuntos de la forma $A \cup \{\infty\}$, donde A es abierto en M y $M \setminus A$ es acotado.

Observemos que estos conjuntos forman realmente una topología. Por una parte, \emptyset y M^∞ son trivialmente abiertos. Si \mathcal{F} es una familia de abiertos de M^∞ , vamos a ver que la unión $V = \bigcup \mathcal{F}$ es abierta. Si ∞ no está en ningún elemento de \mathcal{F} , entonces V es abierto porque M es un espacio topológico: si $\infty \in U \in \mathcal{F}$, entonces, llamando $V' = V \setminus \{\infty\}$, tenemos que $V = V' \cup \{\infty\}$, $V' = \bigcup_{W \in \mathcal{F}} (W \setminus \{\infty\})$ es abierto en M y $M \setminus V' \subset M \setminus U$ es acotado, luego V también es abierto en este caso.

Finalmente, si U_1, \dots, U_n son abiertos en M^∞ y $V = U_1 \cap \dots \cap U_n$, entonces

$$V \setminus \{\infty\} = (U_1 \setminus \{\infty\}) \cap \dots \cap (U_n \setminus \{\infty\})$$

es abierto en M , por ser intersección de abiertos. Si $\infty \notin V$, directamente $V = V \setminus \{\infty\}$ es abierto en M^∞ , y si $\infty \in V$, entonces cada $U'_i = U_i \setminus \{\infty\}$ cumple que $M \setminus U'_i$ es acotado, luego

$$\bigcup_{i=1}^n (M \setminus U'_i) = M \setminus \bigcap_{i=1}^n U'_i = M \setminus (V \setminus \{\infty\})$$

es acotado, por ser unión finita de acotados.

Es inmediato que la topología inducida en M desde M^∞ es la que ya teníamos en M , y también es claro que M^∞ es un espacio de Hausdorff, porque dos puntos de M tienen entornos disjuntos en M , luego en M^∞ , y si $m \in M$, entonces $U = B_1(m)$ y $V = M^\infty \setminus B'_1(m)$ son entornos disjuntos de m e ∞ .

Si M está acotado, entonces $\{\infty\}$ es abierto en M^∞ , por lo que ∞ es un punto aislado y la construcción resulta trivial: el único punto “alrededor de ∞ ” es el propio ∞ , pero si M no está acotado, entonces los puntos alrededor de ∞ son ∞ y los puntos que están “lejos” de cualquier punto de M , en el sentido de estar fuera de cualquier bola de centro dicho punto.

Por ejemplo, si M es un espacio normado, es claro que una base de entornos de ∞ la forman los conjuntos

$$\{\infty\} \cup \{x \in M \mid \|x\| > K\}, \quad K \in \mathbb{R}. \quad \blacksquare$$

2.5 Continuidad

Finalmente estamos en condiciones de formalizar la idea de función continua como función f que envía los alrededores de un punto a los alrededores de su imagen. No exigimos que las imágenes de los puntos de alrededor de un punto x sean todos los puntos de alrededor de $f(x)$. Por ejemplo, sea S la circunferencia unidad en \mathbb{R}^2 y f la aplicación dada por

$$\begin{aligned} [0, 1] &\longrightarrow S \\ x &\longmapsto (\cos 2\pi x, \operatorname{sen} 2\pi x) \end{aligned}$$

Lo que hace f es “pegar” los extremos del intervalo en un mismo punto $(1, 0)$. Queremos que esta aplicación sea continua, y vemos que $[0, 1/4]$ es un entorno de 0 que se transforma en “medio” entorno de $(1, 0)$, en los puntos de alrededor de $(1, 0)$ contenidos en el semiplano $y > 0$. Esto no debe ser, pues, un obstáculo a la continuidad.

Pedir que los puntos de alrededor de x sean enviados a puntos de alrededor de $f(x)$ (no necesariamente todos) es pedir que todo entorno U de $f(x)$ contenga a las imágenes de los puntos de alrededor de x , es decir, las imágenes de un entorno de x , es decir, que $f^{-1}[U]$ contenga un entorno de x , pero esto equivale a que él mismo lo sea. Así pues:

Definición 2.42 Una aplicación $f : X \rightarrow Y$ entre dos espacios topológicos es *continua* en un punto $x \in X$ si para todo entorno U de $f(x)$ se cumple que $f^{-1}[U]$ es un entorno de x . Diremos que f es *continua* si lo es en todos los puntos de X .

Observemos que en esta definición podemos sustituir “entorno” por “entorno básico”. En un espacio métrico podemos considerar concretamente bolas abiertas, y entonces la definición se particulariza como sigue:

Teorema 2.43 Una aplicación $f : M \rightarrow N$ entre dos espacios métricos es continua en un punto $x \in M$ si y sólo si para todo $\epsilon > 0$ existe un $\delta > 0$ tal que si $d(x, x') < \delta$ entonces $d(f(x), f(x')) < \epsilon$.

Veamos varias caracterizaciones de la continuidad.

Teorema 2.44 Sea $f : X \rightarrow Y$ una aplicación entre espacios topológicos. Las afirmaciones siguientes son equivalentes:

- a) f es continua.
- b) Para todo abierto (básico) G de Y se cumple que $f^{-1}[G]$ es abierto en X .
- c) Para todo cerrado C de Y se cumple que $f^{-1}[C]$ es cerrado en X .
- d) Para todo $A \subset X$ se cumple $f[\overline{A}] \subset \overline{f[A]}$.
- e) Para todo $B \subset Y$ se cumple $\overline{f^{-1}[B]} \subset f^{-1}[\overline{B}]$.
- f) Para todo $B \subset Y$ se cumple $f^{-1}[\text{int } B] \subset \text{int } f^{-1}[B]$.

DEMOSTRACIÓN: a) \leftrightarrow b). Si f es continua y $x \in f^{-1}[G]$ entonces $f(x) \in G$, luego G es un entorno de $f(x)$, luego por definición de continuidad $f^{-1}[G]$ es un entorno de x , luego $f^{-1}[G]$ es abierto. Es claro que b) \rightarrow a).

Evidentemente b) \leftrightarrow c).

a) \rightarrow d). Si $x \in f[\overline{A}]$ entonces $x = f(y)$, con $y \in \overline{A}$. Si E es un entorno de x , por definición de continuidad $f^{-1}[E]$ es un entorno de y , luego $f^{-1}[E] \cap A \neq \emptyset$, de donde $E \cap f[A] \neq \emptyset$, lo que prueba que $x \in \overline{f[A]}$.

d) \rightarrow e). Tenemos que $f^{-1}[B] \subset X$, luego $f[\overline{f^{-1}[B]}] \subset \overline{f[f^{-1}[B]]} \subset \overline{B}$, luego $\overline{f^{-1}[B]} \subset f^{-1}[\overline{B}]$.

e) \rightarrow f). En efecto:

$$\begin{aligned} f^{-1}[\text{int } B] &= f^{-1}[Y \setminus (Y \setminus \text{int } B)] = X \setminus f^{-1}[Y \setminus \text{int } B] \\ &= X \setminus f^{-1}[\overline{Y \setminus B}] \subset X \setminus \overline{f^{-1}[Y \setminus B]} = X \setminus \overline{X \setminus f^{-1}[B]} \\ &= X \setminus (X \setminus \text{int } f^{-1}[B]) = \text{int } f^{-1}[B]. \end{aligned}$$

f) \rightarrow b). Si B es abierto en Y , entonces

$$f^{-1}[B] = f^{-1}[\text{int } B] \subset \text{int } f^{-1}[B] \subset f^{-1}[B],$$

luego $f^{-1}[B] = \text{int } f^{-1}[B]$ que es, por lo tanto, abierto. \blacksquare

Ahora vamos a probar una serie de resultados generales que nos permitirán reconocer en muchos casos la continuidad de una aplicación de forma inmediata. De la propia definición de continuidad se sigue inmediatamente:

Teorema 2.45 *Si $f : X \rightarrow Y$ es continua en un punto x y $g : Y \rightarrow Z$ es continua en $f(x)$, entonces $f \circ g$ es continua en x . En particular, la composición de aplicaciones continuas es una aplicación continua.*

Otro hecho básico es que la continuidad depende sólo de la topología en la imagen y no de la del espacio de llegada.

Teorema 2.46 *Sea $f : X \rightarrow Y$ una aplicación entre espacios topológicos. Entonces f es continua en un punto $x \in X$ como aplicación $f : X \rightarrow Y$ si y sólo si lo es como aplicación $f : X \rightarrow f[X]$.*

DEMOSTRACIÓN: Un entorno de $f(x)$ en $f[X]$ es $U \cap f[X]$, donde U es un entorno de $f(x)$ en Y , pero $f^{-1}[U \cap f[X]] = f^{-1}[U]$, luego es indistinto considerar entornos en $f[X]$ o en Y . \blacksquare

Teniendo en cuenta que la aplicación identidad en un conjunto es obviamente continua, de los teoremas anteriores se deduce inmediatamente el que sigue:

Teorema 2.47 *Si X es un espacio topológico y $A \subset X$, entonces la inclusión $i : A \rightarrow X$ dada por $i(x) = x$ es continua. Por tanto, si $f : X \rightarrow Y$ es continua en un punto $x \in A$, la restricción $f|_A = i \circ f$ es continua en x .*

En particular la restricción de una aplicación continua a un subconjunto es también continua. El recíproco no es cierto, pero se cumple lo siguiente:

Teorema 2.48 *Dada una aplicación $f : X \rightarrow Y$, si A es un entorno de un punto $x \in X$ y $f|_A$ es continua en x , entonces f es continua en x .*

DEMOSTRACIÓN: Si U es un entorno de $f(x)$ en Y , entonces $(f|_A)^{-1}[U] = f^{-1}[U] \cap A$ es un entorno de x en A , luego existe un entorno G de x en X de manera que $x \in G \cap A = f^{-1}[U] \cap A$, luego en particular $x \in G \cap A \subset f^{-1}[U]$, y $G \cap A$ es un entorno de x en X , luego $f^{-1}[U]$ también lo es. ■

Esto significa que la continuidad es una propiedad local, es decir, el que una función sea continua o no en un punto es un hecho que sólo depende del comportamiento de la función en un entorno del punto. En particular, si cubrimos un espacio topológico por una familia de abiertos, para probar que una aplicación es continua basta ver que lo es su restricción a cada uno de los abiertos. Esto es cierto también si cubrimos el espacio con cerrados a condición de que sean un número finito:

Teorema 2.49 *Sea $f : X \rightarrow Y$ una aplicación entre espacios topológicos. Sean C_1, \dots, C_n subconjuntos cerrados de X tales que $X = C_1 \cup \dots \cup C_n$. Entonces f es continua si y sólo si cada $f|_{C_i}$ es continua.*

DEMOSTRACIÓN: Una implicación es obvia. Si las restricciones son continuas, entonces dado un cerrado C de Y , se cumple que

$$f^{-1}[C] = (f^{-1}[C] \cap C_1) \cup \dots \cup (f^{-1}[C] \cap C_n) = (f|_{C_1})^{-1}[C] \cup \dots \cup (f|_{C_n})^{-1}[C].$$

Ahora, cada $(f|_{C_i})^{-1}[C]$ es cerrado en C_i , luego es la intersección con C_i de un cerrado de X , luego es la intersección de dos cerrados en X , luego es cerrado en X . Así pues, $f^{-1}[C]$ es la unión de un número finito de cerrados de X , luego es cerrado en X . Esto prueba que f es continua. ■

Teorema 2.50 *Si $\{X_i\}_{i \in I}$ es una familia de espacios topológicos, las proyecciones $p_i : \prod_{i \in I} X_i \rightarrow X_i$ son funciones continuas.*

DEMOSTRACIÓN: Las antiimágenes de abiertos en X_i son abiertos básicos del producto. ■

Ejercicio: Probar que la topología producto es la menor topología que hace continuas a las proyecciones.

Teorema 2.51 *Si $\{X_i\}_{i \in I}$ es una familia de espacios topológicos y X es un espacio topológico, entonces una aplicación $f : X \rightarrow \prod_{i \in I} X_i$ es continua si y sólo si lo son todas las funciones $f_i = f \circ p_i$.*

DEMOSTRACIÓN: Si f es continua las funciones $f \circ p_i$ también lo son por ser composición de funciones continuas.

Si cada $f \circ p_i$ es continua, sea $A = \prod_{i \in I} A_i$ un abierto básico del producto.

Sean i_1, \dots, i_n los índices tales que $A_{i_j} \neq X_{i_j}$. Entonces $f^{-1}[p_{i_j}^{-1}[A_{i_j}]] = (f \circ p_{i_j})^{-1}[A_{i_j}]$ es abierto en X , pero

$$A = \bigcap_{j=1}^n p_{i_j}^{-1}[A_{i_j}] \quad \text{y} \quad f^{-1}[A] = \bigcap_{j=1}^n f^{-1}[p_{i_j}^{-1}[A_{i_j}]]$$

es abierto en X . ■

Así, por ejemplo, para probar que la aplicación $f : \mathbb{R} \rightarrow \mathbb{R}^2$ dada por $f(x) = (x + 1, x^2)$ es continua, basta probar que lo son las aplicaciones $x + 1$ y x^2 .

Definición 2.52 Sean E y F espacios normados sobre un cuerpo métrico \mathbb{K} . Una aplicación $f : E \rightarrow F$ tiene la *propiedad de Lipschitz* si existe un $M > 0$ tal que para todos los vectores $v, w \in E$ se cumple que $\|f(v) - f(w)\| \leq M\|v - w\|$.

Teorema 2.53 Las aplicaciones con la propiedad de Lipschitz son continuas.

DEMOSTRACIÓN: Sea $f : E \rightarrow F$ una aplicación con la propiedad de Lipschitz con constante M . Vamos a aplicar el teorema 2.43. Dado $\epsilon > 0$ tomamos $\delta = \epsilon/M$. Así, si $\|v - w\| < \delta$, entonces $\|f(v) - f(w)\| \leq M\|v - w\| < \epsilon$. ■

Por ejemplo, es fácil ver que si E es un espacio normado entonces la norma $\|\cdot\| : E \rightarrow \mathbb{R}$ tiene la propiedad de Lipschitz con constante $M = 1$, luego es una aplicación continua. Un ejemplo menos trivial es el de la suma:

Teorema 2.54 Sea E un espacio normado. Entonces la suma $+: E \times E \rightarrow E$ tiene la propiedad de Lipschitz, luego es continua.

DEMOSTRACIÓN: Consideraremos a $E \times E$ como espacio normado con la norma $\|\cdot\|_1$. Entonces, si $(u, v), (a, b) \in E \times E$, tenemos que

$$\begin{aligned} \|(u + v) - (a + b)\| &= \|(u - a) + (v - b)\| \leq \|u - a\| + \|v - b\| \\ &= \|(u - a, v - b)\|_1 = \|(u, v) - (a, b)\|_1. \end{aligned}$$

■

El producto no cumple la propiedad de Lipschitz, pero aun así es continuo.

Teorema 2.55 Sea E un espacio normado sobre un cuerpo métrico \mathbb{K} . El producto $\cdot : \mathbb{K} \times E \rightarrow E$ es una aplicación continua.

DEMOSTRACIÓN: Veamos que el producto es continuo en un punto (λ, x) de $\mathbb{K} \times E$. Usaremos la norma $\|\cdot\|_\infty$ en $\mathbb{K} \times E$. Dado $\epsilon > 0$, sea $(\lambda', x') \in \mathbb{K} \times E$.

$$\|\lambda'x' - \lambda x\| = \|\lambda'(x' - x) + (\lambda' - \lambda)x\| \leq |\lambda'| \|x' - x\| + |\lambda' - \lambda| \|x\|.$$

Tomemos ahora $0 < \delta < 1$ que cumpla además

$$\delta < \frac{\epsilon}{|\lambda| + \|x\| + 1} < \epsilon.$$

Así si $\|(\lambda', x') - (\lambda, x)\|_\infty < \delta$, entonces $|\lambda' - \lambda| < \delta$ y $\|x' - x\| < \delta$. En particular $|\lambda'| \leq |\lambda| + \delta < |\lambda| + 1$. Así

$$\|\lambda'x' - \lambda x\| < \frac{|\lambda'| \epsilon}{|\lambda| + \|x\| + 1} + \frac{\|x\| \epsilon}{|\lambda| + \|x\| + 1} < \epsilon. \quad \blacksquare$$

Definición 2.56 Un *espacio vectorial topológico* es un par (V, \mathcal{T}) , donde V es un espacio vectorial sobre un cuerpo métrico \mathbb{K} y \mathcal{T} es una topología de Hausdorff sobre V de modo que las aplicaciones $+$: $V \times V \rightarrow V$ y \cdot : $\mathbb{K} \times V \rightarrow V$ son continuas.

Hemos demostrado que todo espacio normado es un espacio vectorial topológico.

En general, si X es un conjunto y V es un espacio vectorial sobre un cuerpo K , el conjunto V^X de todas las aplicaciones de X en V es un espacio vectorial con las operaciones dadas por $(f + g)(x) = f(x) + g(x)$ y $(\alpha f)(x) = \alpha f(x)$.

Si X e Y son espacios topológicos, llamaremos $C(X, Y)$ al conjunto de todas las aplicaciones continuas de X en Y . Es frecuente abreviar $C(X) = C(X, \mathbb{R})$.

Si X es un espacio topológico y V un espacio vectorial topológico, entonces $C(X, V)$ resulta ser un subespacio vectorial de V^X , pues la suma de dos funciones continuas f y g puede obtenerse como composición de la aplicación $h : X \rightarrow V \times V$ dada por $h(x) = (f(x), g(x))$, que es continua, y la suma en V , que también lo es, luego $f + g$ es continua.

Similarmente se prueba que si $\alpha \in \mathbb{K}$ y $f \in C(X, V)$, entonces $\alpha f \in C(X, V)$.

El espacio \mathbb{K}^X tiene también estructura de anillo conmutativo y unitario con el producto dado por $(fg)(x) = f(x)g(x)$. Como el producto \cdot : $\mathbb{K} \times \mathbb{K} \rightarrow \mathbb{K}$ es continuo (tomando $E = \mathbb{K}$ en el teorema anterior), resulta que $C(X, \mathbb{K})$ es un subanillo de \mathbb{K}^X .

En general, una cuádrupla $(A, +, \cdot, \circ)$, donde la terna $(A, +, \cdot)$ es un espacio vectorial sobre un cuerpo K , la terna $(A, +, \circ)$ es un anillo unitario y además $\alpha(ab) = (\alpha a)b = a(\alpha b)$ para todo $\alpha \in K$, y todos los $a, b \in A$, se llama una K -álgebra.

Tenemos, pues, que si X es un espacio topológico, entonces $C(X, \mathbb{K})$ es una \mathbb{K} -álgebra de funciones. El espacio $C(X, \mathbb{K})$ contiene un subcuerpo isomorfo a \mathbb{K} , a saber, el espacio de las funciones constantes. Cuando no haya confusión identificaremos las funciones constantes con los elementos de \mathbb{K} , de modo que, por ejemplo, 2 representará a la función que toma el valor 2 sobre todos los puntos de X .

Más aún, si $p(x_1, \dots, x_n) \in \mathbb{K}[x_1, \dots, x_n]$, el polinomio p determina una única función evaluación $p : \mathbb{K}^n \rightarrow \mathbb{K}$, de modo que los polinomios constantes se corresponden con las funciones constantes y la indeterminada x_i con la proyección en la i -ésima coordenada. Como todo polinomio es combinación de sumas y productos de constantes e indeterminadas, podemos concluir que $\mathbb{K}[x_1, \dots, x_n] \subset C(\mathbb{K}^n, \mathbb{K})$.

Por ejemplo, la aplicación $f : \mathbb{R}^3 \rightarrow \mathbb{R}^2$ dada por $f(x, y, z) = (3x - 2yz, xyz)$ es claramente continua.

Teorema 2.57 Si \mathbb{K} es un cuerpo métrico, la aplicación $h : \mathbb{K} \setminus \{0\} \rightarrow \mathbb{K}$ dada por $h(x) = 1/x$ es continua.

DEMOSTRACIÓN: Sea $x \in \mathbb{K} \setminus \{0\}$. Sea $\delta = |x|/2$. Si $y \in B_\delta(x)$ entonces $|x| - |y| \leq ||x| - |y|| \leq |x - y| < \delta$, luego $|y| > |x| - \delta = \delta$.

Dado $\epsilon > 0$, sea $\delta' < \delta$, $\delta' < \delta|x|\epsilon$. Así, si $|y - x| < \delta'$ tenemos

$$\left| \frac{1}{y} - \frac{1}{x} \right| = \frac{|y - x|}{|y||x|} < \frac{\delta|x|\epsilon}{\delta|x|} = \epsilon. \quad \blacksquare$$

Consecuentemente, si $f \in C(X, \mathbb{K})$ y f no se anula en X , podemos definir la función $1/f \in C(X, \mathbb{K})$ mediante $(1/f)(x) = 1/f(x)$.

Por ejemplo, la función $f : \mathbb{R} \setminus \{1\} \rightarrow \mathbb{R}$ dada por

$$f(x) = \frac{x^2}{x-1}$$

es obviamente continua.

Teorema 2.58 *La función $\sqrt{\cdot} : [0, +\infty[\rightarrow [0, +\infty[$ es continua.*

DEMOSTRACIÓN: Basta notar que

$$\sqrt{x} \in]a, b[\Leftrightarrow a < \sqrt{x} < b \Leftrightarrow a^2 < x < b^2 \Leftrightarrow x \in]a^2, b^2[.$$

Esto significa que la antiimagen del intervalo $]a, b[$ es el intervalo $]a^2, b^2[$. Similarmente se ve que la antiimagen de un intervalo $[0, b[$ es el intervalo $[0, b^2[$ y, como estos intervalos constituyen una base de $[0, +\infty[$, tenemos que $\sqrt{\cdot}$ es una función continua. \blacksquare

Teorema 2.59 *Sea M un espacio métrico y $A \neq \emptyset$ un subconjunto de M . Entonces las aplicaciones $d : M \times M \rightarrow \mathbb{R}$ y $d(\cdot, A) : M \rightarrow \mathbb{R}$ son continuas.*

DEMOSTRACIÓN: Consideremos en $M \times M$ la distancia d_∞ . Dado un par $(x, y) \in M \times M$, y un $\epsilon > 0$, basta probar que si (x', y') dista de (x, y) menos de $\epsilon/2$, es decir, si se cumple $d(x, x') < \epsilon/2$ y $d(y, y') < \epsilon/2$, entonces tenemos $|d(x, y) - d(x', y')| < \epsilon$. En efecto, en tal caso

$$d(x, y) \leq d(x, x') + d(x', y') + d(y', y) < d(x', y') + \epsilon,$$

luego $d(x, y) - d(x', y') < \epsilon$, e igualmente se llega a $d(x', y') - d(x, y) < \epsilon$, luego efectivamente $|d(x, y) - d(x', y')| < \epsilon$.

Para probar la continuidad de $d(\cdot, A)$ en un punto x observamos que

$$|d(x, A) - d(y, A)| \leq d(x, y).$$

En efecto, para todo $z \in A$ se cumple $d(x, z) \leq d(x, y) + d(y, z)$. De aquí se sigue claramente $d(x, A) \leq d(x, y) + d(y, A)$, y tomando el ínfimo en z vemos que $d(x, A) \leq d(x, y) + d(y, A)$. Similarmente se prueba $d(y, A) \leq d(x, y) + d(x, A)$, de donde se sigue la relación con valores absolutos. A su vez esta relación implica que si $d(x, y) < \epsilon$, entonces $|d(x, A) - d(y, A)| < \epsilon$, lo que expresa la continuidad de la aplicación. \blacksquare

Para terminar con las propiedades generales de las aplicaciones continuas probaremos un hecho de interés teórico:

Teorema 2.60 Sean $f, g : X \rightarrow Y$ aplicaciones continuas, $D \subset X$ un conjunto denso tal que $f|_D = g|_D$ y supongamos que Y es un espacio de Hausdorff. Entonces $f = g$.

DEMOSTRACIÓN: Sea $h : X \rightarrow Y \times Y$ dada por $h(x) = (f(x), g(x))$. Claramente es continua y $h[D] \subset \Delta$, donde

$$\Delta = \{(y, y) \mid y \in Y\}$$

es un cerrado en $Y \times Y$ (por 2.37). Por lo tanto $h[X] = h[\overline{D}] \subset \overline{h[D]} \subset \Delta$, de donde $f = g$. ■

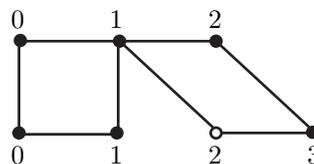
Definición 2.61 Una aplicación biyectiva $f : X \rightarrow Y$ entre dos espacios topológicos es un *homeomorfismo* si f y su inversa son ambas continuas. Dos espacios topológicos son *homeomorfos* si existe un homeomorfismo entre ellos.

Un homeomorfismo induce una biyección entre los abiertos de los dos espacios, por lo que ambos son topológicamente indistinguibles. Es claro que cualquier propiedad definida exclusivamente a partir de la topología se conserva por homeomorfismos, luego dos espacios homeomorfos tienen las mismas propiedades topológicas.

Es importante notar que no toda biyección continua es un homeomorfismo. Por ejemplo, la identidad $I : \mathbb{R} \rightarrow \mathbb{R}$, cuando en el primer espacio consideramos la topología discreta y en el segundo la euclídea es biyectiva y continua, pero no un homeomorfismo. Veamos un ejemplo más intuitivo:

Ejemplo Sea $f : [0, 1] \cup]2, 3] \rightarrow [0, 2]$ la aplicación dada por

$$f(x) = \begin{cases} x & \text{si } 0 \leq x \leq 1 \\ x - 1 & \text{si } 2 < x \leq 3 \end{cases}$$



Es fácil ver que f es biyectiva, y es continua porque sus restricciones a los abiertos $[0, 1]$ y $]2, 3]$ de su dominio son ambas continuas (son polinomios). Sin embargo no es un homeomorfismo. La aplicación f^{-1} es continua en todos los puntos de $[0, 2]$ excepto en $x = 1$. En efecto, $[0, 1]$ es un entorno de $f^{-1}(1) = 1$, pero la antiimagen de este intervalo es el mismo $[0, 1]$, que no es un entorno de 1 en $[0, 2]$. En los demás puntos es continua, pues f^{-1} restringida a los abiertos de su dominio $[0, 1[$ y $]1, 2]$ es polinómica.

Lo que sucede es que, a pesar de ser biyectiva, f está “pegando” los intervalos $[0, 1]$ y $]2, 3]$ en el intervalo $[0, 2]$, por lo que f^{-1} “corta” éste por el punto 1. En general, si una aplicación continua es una aplicación que no corta, aunque puede pegar, un homeomorfismo es una aplicación que no corta ni pega. Para que no pegue ha de ser biyectiva, pero acabamos de ver que esto no es suficiente. El ejemplo de la topología discreta se interpreta igual: los puntos de \mathbb{R} con la topología discreta están todos “separados” entre sí, luego al pasar a \mathbb{R} con la topología euclídea los estamos “pegando”, aunque no identifiquemos puntos. ■

Definición 2.62 Una aplicación $f : X \rightarrow Y$ entre dos espacios topológicos es *abierto* si para todo abierto A de X , se cumple que $f[A]$ es abierto en Y .

De este modo, un homeomorfismo es una biyección continua y abierta. La propiedad de ser abierta no es muy intuitiva y tampoco es muy frecuente (salvo en el caso de los homeomorfismos). Sin embargo es de destacar el hecho siguiente:

Teorema 2.63 *Las proyecciones de un espacio producto en cada uno de sus factores son aplicaciones abiertas.*

DEMOSTRACIÓN: Basta ver que la proyección de un abierto básico es un abierto, pero los abiertos básicos son productos de abiertos, y las proyecciones son sus factores. ■

Definición 2.64 La *gráfica* de una aplicación $f : X \rightarrow Y$ es el conjunto¹

$$\{(x, f(x)) \in X \times Y \mid x \in X\}.$$

En los casos de funciones $f : A \subset \mathbb{R}^m \rightarrow \mathbb{R}^n$, donde $m + n \leq 3$ la gráfica de f tiene una interpretación en el plano o espacio euclídeo intuitivos, y esta representación permite reconocer rápidamente las características de f . El resultado más importante por lo que a la continuidad se refiere es el siguiente:

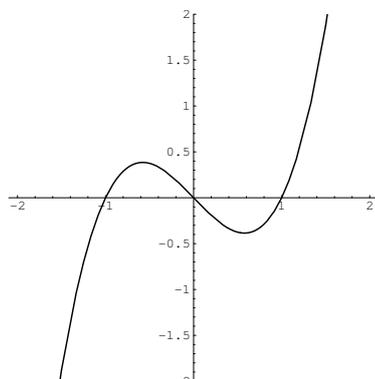
Teorema 2.65 *Si $f : X \rightarrow Y$ es una aplicación continua, entonces la aplicación $x \mapsto (x, f(x))$ es un homeomorfismo entre X y la gráfica de f .*

DEMOSTRACIÓN: La aplicación es obviamente biyectiva y continua. Su inversa es la restricción a la gráfica de la proyección sobre el primer factor de $X \times Y$, luego también es continua. ■

Ejemplo La gráfica del polinomio $x^3 - x$ es la que se muestra en la página siguiente.² Vemos que es una línea ondulada. Podemos considerarla como una imagen “típica” de espacio homeomorfo a \mathbb{R} . Se obtiene deformando la recta “elásticamente”, sin cortarla ni pegarla. Pero notemos que la aplicación f no es un homeomorfismo, la gráfica muestra cómo transforma a \mathbb{R} en su imagen: ésta resulta de “aplstar” la curva sobre el eje vertical, con lo que \mathbb{R} se “pliega” sobre sí mismo, de modo que parte de sus puntos se superponen tres a tres. ■

¹Observemos que desde el punto de vista de la teoría de conjuntos la gráfica de una función f es la propia función f como conjunto.

²El lector debería preguntarse cómo se sabe que la gráfica de f tiene esta forma y no otra. Más adelante veremos cómo determinar analíticamente las características de la gráfica de una función, pero de momento nos bastará con lo siguiente: Para obtener una figura como la anterior, basta programar a un ordenador para que calcule la función considerada sobre los suficientes números racionales, digamos sobre los números de la forma $k/100$, donde k varía entre -200 y 200 , y dibuje un pequeño cuadrado con coordenadas $(x, f(x))$. El resultado es una gráfica como la que hemos mostrado. El proceso sólo involucra la aritmética de los números racionales, que no tiene ninguna dificultad.

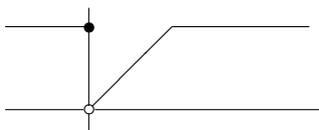


Veamos ahora un ejemplo de gráfica de una función discontinua.

Ejemplo Consideremos la función $f : \mathbb{R} \rightarrow [0, 1]$ dada por

$$f(x) = \begin{cases} 1 & \text{si } x \leq 0 \\ x & \text{si } 0 < x < 1 \\ 1 & \text{si } x \geq 1 \end{cases}$$

Su gráfica es la siguiente:



Es claro que f es continua en todo punto distinto de 0. En efecto, su restricción al abierto $]-\infty, 0[$ es constante, luego continua, su restricción al abierto $]0, +\infty[$ también es continua, pues este abierto es a su vez unión de dos cerrados, $]0, 1]$ y $[1, +\infty[$, en los cuales f es continua, pues en el primero es el polinomio x y en el segundo es constante. Es fácil ver que f no es continua en 0.

La gráfica de f no es homeomorfa a \mathbb{R} , aunque no estamos en condiciones de probarlo ahora. Es fácil ver que $x \mapsto (x, f(x))$ no es un homeomorfismo, pero esto no prueba que no exista otra biyección que sí lo sea. De todos modos, intuitivamente vemos que la gráfica está formada por dos piezas, por lo que para transformar \mathbb{R} en la gráfica es necesario “cortar” por algún punto. ■

Ejemplo Los homeomorfismos no pueden “cortar” ni “pegar”, pero sí “estirar” arbitrariamente un conjunto. Por ejemplo, la homotecia $f(x) = ax$, donde $a > 0$ transforma el intervalo $[-1, 1]$ en el intervalo $[-a, a]$. Las traslaciones $x \mapsto x + b$ son claramente homeomorfismos de \mathbb{R} , luego combinando homotecias y traslaciones podemos construir un homeomorfismo entre cualquier par de intervalos cerrados de la forma $[a, b]$. También es claro que cualquier par de intervalos abiertos $]a, b[$ son homeomorfos entre sí, al igual que cualquier par de intervalos semiabiertos $[a, b[$ y $]a, b]$. En la sección siguiente veremos que también es

posible “estirar infinitamente” un intervalo, de modo que, por ejemplo, $]0, 1[$ es homeomorfo a \mathbb{R} . ■

La topología euclídea Sea \mathbb{K} un cuerpo métrico. Entonces toda aplicación lineal $f : \mathbb{K}^n \rightarrow \mathbb{K}^m$ es continua, pues cada una de sus funciones coordenadas es un polinomio. En particular todo automorfismo de \mathbb{K}^n es un homeomorfismo. Si V es cualquier \mathbb{K} -espacio vectorial de dimensión finita n , existe un isomorfismo $f : V \rightarrow \mathbb{K}^n$. Podemos considerar la topología en V formada por los conjuntos $f^{-1}[G]$, donde G es abierto en \mathbb{K}^n para la topología euclídea. Es claro que V es así un espacio topológico homeomorfo a \mathbb{K}^n . Además, la topología en V no depende de la elección de f , pues si $g : V \rightarrow \mathbb{K}^n$ es cualquier otro isomorfismo y G es un abierto en \mathbb{K}^n , entonces $g^{-1}[G] = f^{-1}[(g^{-1} \circ f)[G]]$ y $(g^{-1} \circ f)[G]$ es un abierto en \mathbb{K}^n , porque $g^{-1} \circ f$ es un isomorfismo y por consiguiente un homeomorfismo.

Más aún, la aplicación $\| \cdot \| : V \rightarrow \mathbb{R}$ dada por $\|v\| = \|f(v)\|$ es claramente una norma en V que induce la topología que acabamos de definir. Esto justifica las definiciones siguientes:

Definición 2.66 Un *isomorfismo topológico* entre dos espacios vectoriales topológicos es una aplicación entre ambos que sea a la vez isomorfismo y homeomorfismo. Si V es un \mathbb{K} -espacio vectorial de dimensión finita n sobre un cuerpo métrico \mathbb{K} , llamaremos *topología euclídea* en V a la única topología respecto a la cual todos los isomorfismos de V en \mathbb{K}^n son topológicos.

Si $h : V \rightarrow W$ es una aplicación lineal entre dos \mathbb{K} -espacios vectoriales de dimensión finita, entonces es continua, pues considerando isomorfismos (topológicos) $f : V \rightarrow \mathbb{K}^m$ y $g : W \rightarrow \mathbb{K}^n$ tenemos que la aplicación $h' = f^{-1} \circ h \circ g : \mathbb{K}^m \rightarrow \mathbb{K}^n$ es lineal, luego continua, luego $h = f \circ h' \circ g^{-1}$ también es continua.

Si W es un subespacio de V , entonces la restricción a W de la topología euclídea en V es la topología euclídea. En efecto, para probarlo descomponemos $V = W \oplus W'$ y observamos que la identidad de W con la topología euclídea a W con la topología inducida es continua por ser lineal y su inversa es continua por ser la restricción de la proyección de V en W , que también es lineal. Por lo tanto se trata de un homeomorfismo y ambas topologías coinciden.

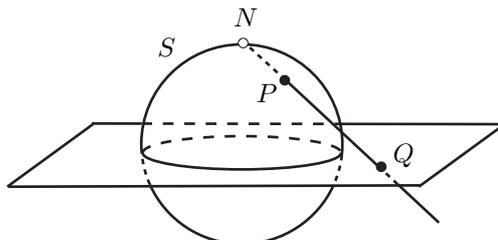
Similarmente se prueba que la topología euclídea en un producto de espacios vectoriales coincide con el producto de las topologías euclídeas.

Todo subespacio W de un espacio vectorial de dimensión finita V es cerrado. Para probarlo observamos que W puede expresarse como el núcleo de una aplicación lineal de W en \mathbb{K}^n , es decir, como la antiimagen de 0 por una aplicación continua y, como $\{0\}$ es cerrado, su antiimagen también.

Todos estos hechos se trasladan sin dificultad a los espacios afines sobre \mathbb{K} . La topología euclídea en un espacio afín se define de modo que todas las afinidades son continuas y las variedades afines son cerradas.

Ejemplo: La proyección estereográfica Consideremos la esfera en \mathbb{R}^{n+1} de centro 0 y radio 1, es decir, el conjunto $S = \{x \in \mathbb{R}^{n+1} \mid \|x\|^2 = 1\}$. Sea $N = (0, \dots, 0, 1)$ el “polo Norte” de S .

La *proyección estereográfica* es la biyección entre $S \setminus \{N\}$ y \mathbb{R}^n que a cada punto $P \in S$ le asigna el punto Q donde la recta NP corta al plano $x_{n+1} = 0$ (que podemos identificar con \mathbb{R}^n eliminando la última coordenada de sus puntos).



Si P tiene coordenadas x , la recta NP está formada por los puntos

$$(0, \dots, 0, 1) + \lambda(x_1, \dots, x_n, x_{n+1} - 1),$$

con $\lambda \in \mathbb{R}$. El valor de λ que anula la última coordenada es el que cumple $1 + \lambda(x_{n+1} - 1) = 0$, o sea, $\lambda = 1/(1 - x_{n+1})$, luego la proyección de P es el punto

$$f(x) = \left(\frac{x_1}{1 - x_{n+1}}, \dots, \frac{x_n}{1 - x_{n+1}} \right).$$

Similarmente se calcula la proyección inversa, que es

$$g(y) = \left(\frac{2y_1}{\|y\|^2 + 1}, \dots, \frac{2y_n}{\|y\|^2 + 1}, \frac{\|y\|^2 - 1}{\|y\|^2 + 1} \right).$$

Es evidente que tanto f como g son continuas, luego la proyección estereográfica es un homeomorfismo. Así pues, el plano es homeomorfo a una esfera menos un punto.

Más aún, llamemos $X = \mathbb{R}^n$ y consideremos $X^\infty = \mathbb{R}^n \cup \{\infty\}$ con la topología que hemos definido al final de la sección anterior. Vamos a probar que si extendemos $f : S \rightarrow X^\infty$ mediante $f(N) = \infty$ obtenemos un homeomorfismo, con lo cual tenemos una interpretación geométrica de X^∞ : al añadirle a \mathbb{R}^n un punto infinito, el resultado es, desde el punto de vista topológico, lo mismo que una esfera de la misma dimensión.

Sólo falta probar que f es continua en N y que f^{-1} es continua en ∞ . Para ello basta probar que f biyecta una base de entornos de N con una base de entornos de ∞ .

Consideremos, concretamente, la base de entornos de N formada por las bolas abiertas de radio menor que 1. Un entorno básico de N está formado por los puntos $x \in S$ tales que $\sqrt{x_1^2 + \dots + x_n^2 + (x_{n+1} - 1)^2} < \epsilon < 1$, y como $\|x\|^2 = 1$, esto equivale a $x_{n+1} > 1 - \epsilon^2/2$. Puesto que ϵ es arbitrario, los entornos básicos están formados por los puntos $x \in S$ tales que $1 - x_{n+1} < \epsilon$, para $\epsilon > 0$. Calculemos la imagen de uno de estos entornos.

Para ello notamos que si $x \neq N$,

$$\|f(x)\| = \sqrt{\frac{1+x_{n+1}}{1-x_{n+1}}},$$

y que si $x_{n+1} < x'_{n+1}$ entonces $x_{n+1} - x_{n+1}x'_{n+1} < x'_{n+1} - x_{n+1}x'_{n+1}$, luego

$$\frac{x_{n+1}}{1-x_{n+1}} < \frac{x'_{n+1}}{1-x'_{n+1}}, \quad \sqrt{1 + \frac{2x_{n+1}}{1-x_{n+1}}} < \sqrt{1 + \frac{2x'_{n+1}}{1-x'_{n+1}}},$$

$$\sqrt{\frac{1+x_{n+1}}{1-x_{n+1}}} < \sqrt{\frac{1+x'_{n+1}}{1-x'_{n+1}}}.$$

En otras palabras, cuanto mayor es x_{n+1} , mayor es la distancia a 0 de $f(x)$. Por consiguiente, los puntos tales que $x_{n+1} > 1 - \epsilon$ se corresponden con los puntos tales que

$$\|f(x)\| > \sqrt{\frac{1+1-\epsilon}{1-(1-\epsilon)}} = \sqrt{\frac{2}{\epsilon} - 1}.$$

Ahora bien, cualquier $K > 0$ es de esta forma para algún ϵ , concretamente $\epsilon = 2/(K^2 + 1)$. Concluimos que las imágenes de los entornos básicos de N son los conjuntos de la forma

$$\{u \in X \mid \|u\| > K\} \cup \{\infty\},$$

pero sabemos que estos conjuntos forman una base de entornos de ∞ en X^∞ .

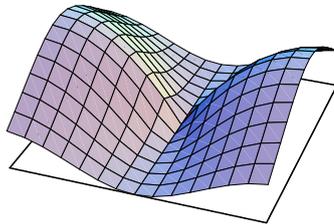
En particular tenemos que \mathbb{R}^∞ es homeomorfo a una circunferencia. ■

2.6 Límites de funciones

El concepto de límite es, junto al de continuidad, uno de los conceptos más importantes a los que la estructura topológica sirve de soporte. Para comprender su contenido podemos considerar la función $f : \mathbb{R}^2 \setminus \{(0,0)\} \rightarrow \mathbb{R}$ dada por

$$f(x,y) = x^2 \left(\frac{2}{\sqrt{x^2 + y^2}} - 1 \right).$$

Esta expresión no tiene sentido cuando $(x,y) = (0,0)$, por lo que es natural preguntarse si la gráfica de f mostrará alguna particularidad que explique por qué no puede ser calculada en este punto. He aquí dicha gráfica:



Su aspecto es el de una superficie homeomorfa a \mathbb{R}^2 , pero sabemos que no está definida en $(0, 0)$. Sin embargo, la gráfica hace pensar que $f(0, 0) = 0$. La explicación es que la función $f : \mathbb{R}^2 \rightarrow \mathbb{R}$ definida mediante

$$f(x, y) = \begin{cases} x^2 \left(\frac{2}{\sqrt{x^2 + y^2}} - 1 \right) & \text{si } (x, y) \neq (0, 0), \\ 0 & \text{si } (x, y) = (0, 0), \end{cases}$$

es continua, aunque esto no es evidente y, de hecho, no estamos en condiciones de probarlo ahora. Si queremos expresar la situación en términos de la expresión original de f , no definida en $(0, 0)$, habremos de decir que f transforma los puntos de alrededor de $(0, 0)$ en puntos de alrededor de 0, o también, que los valores que toma f son más cercanos a 0 cuanto más cercanos a $(0, 0)$ son los puntos que consideramos. Todo esto tiene sentido aunque la función f no esté definida en $(0, 0)$.

Comenzamos introduciendo el concepto topológico que permite expresar con rigor las ideas que acabamos de introducir:

Definición 2.67 Sean X, Y espacios topológicos, $A \subset X$ y $f : A \rightarrow Y$. Sea $a \in A$ y $b \in Y$. Diremos que f converge a b cuando x tiende a a si para todo entorno V de b existe un entorno U de a tal que si $x \in U \cap A$ y $x \neq a$, entonces $f(x) \in V$.

La interpretación es clara: los puntos $f(x)$ están alrededor de b [= en un entorno arbitrario V de b] siempre que x está alrededor de a [= en un entorno adecuado U de a , que dependerá de V], o más simplemente: si f envía los puntos de alrededor de a a los alrededores de b .

Si Y es un espacio de Hausdorff, una función converge a lo sumo a un único punto b para cada punto a . En efecto, si f converge a dos puntos b y b' cuando x tiende a a , podríamos tomar entornos disjuntos V y V' de b y b' , para los cuales existirían entornos U y U' de a de modo que si $x \in U \cap U' \cap A$, entonces $f(x) \in V \cap V'$, contradicción.

Por ello, si se da la convergencia, diremos que b es el *límite* cuando x tiende a a de $f(x)$, y lo representaremos por

$$b = \lim_{x \rightarrow a} f(x).$$

No exigimos que la función f esté definida en a . Tan sólo que a sea un punto de acumulación del dominio de f o, de lo contrario, no existirían puntos x a los que aplicar la definición y f convergería trivialmente a todos los puntos de Y .

En estos términos, lo que afirmábamos antes es que existe

$$\lim_{(x,y) \rightarrow (0,0)} x^2 \left(\frac{2}{\sqrt{x^2 + y^2}} - 1 \right) = 0,$$

de modo que los valores que toma esta expresión se acercan más a 0 cuanto más se acercan las variables al punto $(0, 0)$. Todavía no podemos probarlo.

Por supuesto, es posible que la función f esté definida en a , pero esto es irrelevante, pues en la definición de límite aparecen sólo puntos $x \neq a$, lo que significa que el límite es independiente de $f(a)$. En otras palabras, si modificáramos el valor de $f(a)$, la existencia del límite y su valor concreto no se alterarían.

También es obvio que la existencia o no de límite sólo depende del comportamiento de la función en un entorno del punto. En otras palabras, que si dos funciones coinciden en un entorno de un punto a (salvo quizá en a) entonces ambas tienen límite en a o ninguna lo tiene y, si lo tienen, éstos coinciden.

La relación entre los límites y la continuidad es la siguiente:

Teorema 2.68 Sean X, Y espacios topológicos y $f : X \rightarrow Y$. Sea $a \in X'$. Entonces f es continua en a si y sólo si existe $\lim_{x \rightarrow a} f(x) = f(a)$.

DEMOSTRACIÓN: Si el límite vale $f(a)$, entonces la definición de límite se cumple trivialmente cuando $x = a$ y lo que queda es la definición de continuidad en a . Recíprocamente, la definición de continuidad de f en a implica trivialmente la definición de límite con $b = f(a)$. ■

Mediante este teorema podemos deducir las propiedades algebraicas de los límites a partir de las propiedades correspondientes de las funciones continuas.

Teorema 2.69 Sea \mathbb{K} un cuerpo métrico, sean $f, g : A \subset X \rightarrow \mathbb{K}$ dos funciones definidas sobre un espacio topológico X y sea $a \in A'$. Si existen

$$\lim_{x \rightarrow a} f(x) \quad \text{y} \quad \lim_{x \rightarrow a} g(x)$$

entonces también existen

$$\begin{aligned} \lim_{x \rightarrow a} (f(x) + g(x)) &= \lim_{x \rightarrow a} f(x) + \lim_{x \rightarrow a} g(x), \\ \lim_{x \rightarrow a} f(x)g(x) &= \left(\lim_{x \rightarrow a} f(x) \right) \left(\lim_{x \rightarrow a} g(x) \right), \\ \lim_{x \rightarrow a} \frac{f(x)}{g(x)} &= \frac{\lim_{x \rightarrow a} f(x)}{\lim_{x \rightarrow a} g(x)}, \end{aligned}$$

(suponiendo, además, en el tercer caso que $\lim_{x \rightarrow a} g(x) \neq 0$.)

DEMOSTRACIÓN: La prueba es la misma en todos los casos. Veamos la primera. Consideramos la función f' que coincide con f en todos los puntos salvo en a , donde toma el valor del límite. Definimos igualmente g' . Entonces el teorema anterior implica que f' y g' son continuas en a , luego $f' + g'$ también lo es, luego por el teorema anterior existe

$$\lim_{x \rightarrow a} (f'(x) + g'(x)) = f'(a) + g'(a),$$

pero como $f' + g'$ coincide con $f + g$ salvo en a , el límite de $f' + g'$ coincide con el de $f + g$, y tenemos la relación buscada. ■

Es claro que el resultado sobre suma de límites es válido cuando las funciones toman valores en cualquier espacio vectorial topológico. Además en tal caso al multiplicar una función por un escalar el límite se multiplica por el mismo (para el caso de \mathbb{K} esto es un caso particular de la segunda igualdad). La misma técnica nos da inmediatamente:

Teorema 2.70 *Sea $f : A \subset X \rightarrow X_1 \times \cdots \times X_n$ una aplicación entre espacios topológicos, que será de la forma $f(x) = (f_1(x), \dots, f_n(x))$, para ciertas funciones $f_i : X \rightarrow X_i$. Sea $a \in A'$. Entonces existe $\lim_{x \rightarrow a} f(x)$ si y sólo si existen todos los límites $\lim_{x \rightarrow a} f_i(x)$ y en tal caso*

$$\lim_{x \rightarrow a} f(x) = \left(\lim_{x \rightarrow a} f_1(x), \dots, \lim_{x \rightarrow a} f_n(x) \right).$$

Para calcular el límite que tenemos pendiente necesitamos un hecho que a menudo resulta útil. Diremos que una función f con valores en un espacio métrico está *acotada* si su imagen es un conjunto acotado.

Teorema 2.71 *Sean $f, g : A \subset X \rightarrow \mathbb{K}$ dos funciones definidas de un espacio topológico X en un cuerpo métrico \mathbb{K} y sea $a \in A'$. Si existe $\lim_{x \rightarrow a} f(x) = 0$ y f está acotada, entonces existe $\lim_{x \rightarrow a} f(x)g(x) = 0$.*

DEMOSTRACIÓN: Si g está acotada, existe un $M > 0$ tal que $|g(x)| \leq M$ para todo $x \in A$. Entonces $|f(x)g(x)| \leq M|f(x)|$. El hecho de que f tienda a 0, substituyendo los entornos en \mathbb{K} de la definición general por bolas abiertas, queda así:

Para todo $\epsilon > 0$ existe un entorno U de a tal que si $x \in U \cap A$ y $x \neq a$, entonces $|f(x)| < \epsilon$.

Tomamos ahora $\epsilon > 0$ y aplicamos este hecho a ϵ/M , con lo que si $x \in U \cap A$ y $x \neq a$, tenemos que $|f(x)g(x)| \leq M|f(x)| < \epsilon$, y esto significa que fg tiende a 0. ■

Ejemplo Se cumple

$$\lim_{(x,y) \rightarrow (0,0)} x^2 \left(\frac{2}{\sqrt{x^2 + y^2}} - 1 \right) = 0.$$

En efecto, basta probar que

$$\lim_{(x,y) \rightarrow (0,0)} \frac{2x^2}{\sqrt{x^2 + y^2}} = 0,$$

pues el otro sumando, x^2 tiende obviamente a 0, por continuidad. Factorizamos

$$x \frac{2x}{\sqrt{x^2 + y^2}}.$$

El primer factor tiende a 0 y el segundo está acotado, pues se comprueba fácilmente que

$$-1 \leq \frac{x}{\sqrt{x^2 + y^2}} \leq 1.$$

Ahora basta aplicar el teorema anterior. ■

El hecho de que la composición de funciones continuas es continua se traduce ahora en el hecho siguiente:

Teorema 2.72 Sea $f : X \rightarrow Y$ y $g : Y \rightarrow Z$, sea $a \in X'$ y supongamos que existe

$$\lim_{x \rightarrow a} f(x) = b$$

y que g es continua en b . Entonces

$$g\left(\lim_{x \rightarrow a} f(x)\right) = \lim_{x \rightarrow a} g(f(x)).$$

DEMOSTRACIÓN: Sea U un entorno de $g(b)$. Entonces $g^{-1}[U]$ es un entorno de b . Existe un entorno V de a tal que si $x \in V$, $x \neq a$, entonces $f(x) \in g^{-1}[U]$, luego $g(f(x)) \in U$. Por lo tanto se cumple la definición de límite. ■

Por ejemplo, la continuidad de la raíz cuadrada implica que

$$\lim_{(x,y) \rightarrow (0,0)} |x| \sqrt{\frac{2}{\sqrt{x^2 + y^2}} - 1} = 0.$$

Ejercicio: Dar un ejemplo que muestre la falsedad de la afirmación siguiente: Dadas dos funciones $f, g : \mathbb{R} \rightarrow \mathbb{R}$, si existen $\lim_{x \rightarrow a} f(x) = b$, $\lim_{x \rightarrow b} g(x) = c$ entonces existe también $\lim_{x \rightarrow a} g(f(x)) = c$. Probar que es cierta si $f(x) \neq b$ para $x \neq a$.

Límites restringidos El límite de una función en un punto depende del dominio de ésta, por eso es importante lo que ocurre al restringir una función a dominios menores. Por ejemplo pensemos en la función $f : \mathbb{R} \rightarrow \mathbb{R}$ dada por

$$f(x) = \begin{cases} -1 & \text{si } x < 0 \\ 1 & \text{si } x \geq 0 \end{cases}$$

Es claro que no tiene límite en 0, pero sí lo tienen las funciones $f|_{]-\infty, 0[}$ y $f|_{]0, +\infty[}$. Si consideramos sólo la parte de la izquierda de la función, resulta que es constante igual a -1 , de donde su límite es -1 . Igualmente el límite de la parte derecha es 1. Por ello definimos:

Definición 2.73 Sea $f : A \subset X \rightarrow Y$, $B \subset A$ y $a \in B'$. Definimos

$$\lim_{\substack{x \rightarrow a \\ B}} f(x) = \lim_{x \rightarrow a} f|_B(x).$$

Para el caso de funciones definidas sobre subconjuntos de \mathbb{R} se define

$$\lim_{x \rightarrow a^-} f(x) = \lim_{\substack{x \rightarrow a \\]-\infty, a[}} f(x), \quad \lim_{x \rightarrow a^+} f(x) = \lim_{\substack{x \rightarrow a \\]a, +\infty[}} f(x),$$

y se llaman, respectivamente, límite por la izquierda y límite por la derecha de f en a . También se les llama *límites laterales*. Su utilidad se debe al teorema siguiente:

Teorema 2.74 *Sea A un subconjunto de un espacio X y $f : A \rightarrow Y$. Supongamos que $A = B_1 \cup \dots \cup B_n$ y que a es un punto de acumulación de todos estos conjuntos. Entonces existe $\lim_{x \rightarrow a} f(x)$ si y sólo si existen los límites $\lim_{x \rightarrow a} f(x)$ para B_i $i = 1, \dots, n$ y todos coinciden. En tal caso $f(x)$ es igual al límite común.*

DEMOSTRACIÓN: Si existen tales límites y todos valen L , sea V un entorno de L . Por definición existen entornos U_i de a tales que si $x \in U_i \cap (A \cap B_i)$ y $x \neq a$, entonces $f(x) \in V$. Si U es la intersección de los conjuntos U_i , entonces U es un entorno de a y si $x \in U \cap A$, $x \neq a$, existirá un i tal que $x \in B_i$, luego $f(x) \in V$, es decir, $\lim_{x \rightarrow a} f(x) = L$. El recíproco es más sencillo. ■

Ejemplo Se cumple

$$\lim_{(x,y) \rightarrow (0,0)} x \sqrt{\frac{2}{\sqrt{x^2 + y^2}} - 1} = 0.$$

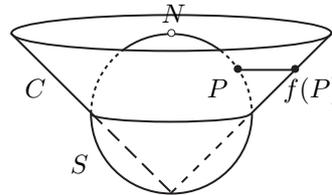
Para comprobarlo calculamos los límites

$$\begin{aligned} \lim_{\substack{(x,y) \rightarrow (0,0) \\ x \leq 0}} x \sqrt{\frac{2}{\sqrt{x^2 + y^2}} - 1} &= - \lim_{\substack{(x,y) \rightarrow (0,0) \\ x \leq 0}} |x| \sqrt{\frac{2}{\sqrt{x^2 + y^2}} - 1} = 0, \\ \lim_{\substack{(x,y) \rightarrow (0,0) \\ x \geq 0}} x \sqrt{\frac{2}{\sqrt{x^2 + y^2}} - 1} &= \lim_{\substack{(x,y) \rightarrow (0,0) \\ x \geq 0}} |x| \sqrt{\frac{2}{\sqrt{x^2 + y^2}} - 1} = 0, \end{aligned}$$

y aplicamos el teorema anterior. ■

Ejemplo: la proyección cónica Veamos ahora un caso en el que nos aparece de forma natural el cálculo de un límite.

Consideremos de nuevo la esfera S y el cono C formado al unir el polo sur $(0, 0 - 1)$ con los puntos del ecuador de S . Vamos a probar que $S \setminus \{N\}$ es homeomorfo a una porción de C mediante la aplicación que traslada radialmente cada punto.



Un punto (x, y, z) de C está en el segmento que une el punto $(0, 0, -1)$ con un punto $(a, b, 0)$, donde $a^2 + b^2 = 1$. Por tanto será de la forma

$$(x, y, z) = (0, 0, -1) + \lambda(a, b, 1), \quad \text{con } \lambda \in \mathbb{R}.$$

Entonces $\lambda = z + 1$, luego $x = a(z + 1)$, $y = b(z + 1)$. De aquí se sigue que

$$z + 1 = \pm \sqrt{x^2 + y^2}.$$

Recíprocamente, si un punto cumple esta ecuación, si $z = -1$ ha de ser $(0, 0, -1)$, que está en C , y si $z \neq -1$, entonces los valores

$$\lambda = z + 1, \quad a = \frac{x}{z + 1}, \quad b = \frac{y}{z + 1},$$

permiten expresar a (x, y, z) en la forma paramétrica anterior, luego se trata de un punto de C . Nos interesa sólo la porción de C formada por los puntos con altura entre -1 y 1 , por lo que definimos

$$C = \{(x, y, z) \in \mathbb{R}^3 \mid z + 1 = \sqrt{x^2 + y^2}, -1 \leq z < 1\}.$$

Notar que los puntos con $z = 1$ no están en C . El homeomorfismo que buscamos ha de transformar cada punto (x, y, z) de $S \setminus \{N\}$ en un punto $(\lambda x, \lambda y, z)$, donde $\lambda \geq 0$ es el adecuado para llegar a C . La condición es

$$\sqrt{(\lambda x)^2 + (\lambda y)^2} = z + 1$$

y, teniendo en cuenta que los puntos de la esfera cumplen $\sqrt{x^2 + y^2} = \sqrt{1 - z^2}$,

$$\lambda = \frac{1 + z}{\sqrt{x^2 + y^2}} = \sqrt{\frac{1 + z}{1 - z}}.$$

Por lo tanto $f : S \setminus \{N\} \rightarrow C$ será la aplicación dada por

$$f(x, y, z) = \left(\sqrt{\frac{1 + z}{1 - z}} x, \sqrt{\frac{1 + z}{1 - z}} y, z \right).$$

La inversa se calcula sin dificultad a partir de esta expresión:

$$g(x, y, z) = \left(\sqrt{\frac{1 - z}{1 + z}} x, \sqrt{\frac{1 - z}{1 + z}} y, z \right), \quad \text{si } (x, y, z) \neq (0, 0, -1),$$

y por supuesto $g(0, 0, -1) = (0, 0, -1)$.

Obviamente f es continua. Lo mismo vale para g en todos los puntos salvo en $(0, 0, -1)$. Para probar la continuidad en este punto hacemos uso de la igualdad $1 + z = \sqrt{x^2 + y^2}$, que cumplen todos los puntos de C , la cual nos permite expresar g como

$$g(x, y, z) = \left(x \sqrt{\frac{2}{\sqrt{x^2 + y^2}} - 1}, y \sqrt{\frac{2}{\sqrt{x^2 + y^2}} - 1}, z \right).$$

Basta probar que

$$\lim_{(x,y,z) \rightarrow (0,0,-1)} \left(x \sqrt{\frac{2}{\sqrt{x^2+y^2}} - 1}, y \sqrt{\frac{2}{\sqrt{x^2+y^2}} - 1}, z \right) = (0, 0, -1),$$

pero los dos primeros límites son el que hemos calculado como ejemplo a lo largo de esta sección.

Así pues, f es un homeomorfismo entre $S \setminus \{N\}$ y C . Ahora bien, es inmediato comprobar que la proyección sobre el plano XY es un homeomorfismo entre C y la bola abierta de centro 0 y radio 2 (la aplicación inversa es $(x, y) \mapsto (x, y, -1 + \sqrt{x^2 + y^2})$, claramente continua), con lo que la composición

$$h(x, y, z) = \left(\sqrt{\frac{1+z}{1-z}} x, \sqrt{\frac{1+z}{1-z}} y \right).$$

es un homeomorfismo entre $S \setminus \{N\}$ y dicha bola. Más aún, si componemos la inversa de la proyección estereográfica con esta aplicación obtenemos un homeomorfismo entre \mathbb{R}^2 y el disco abierto. Es fácil ver que la composición es

$$t(x, y) = \left(\frac{2x \sqrt{x^2 + y^2}}{x^2 + y^2 + 1}, \frac{2y \sqrt{x^2 + y^2}}{x^2 + y^2 + 1} \right).$$

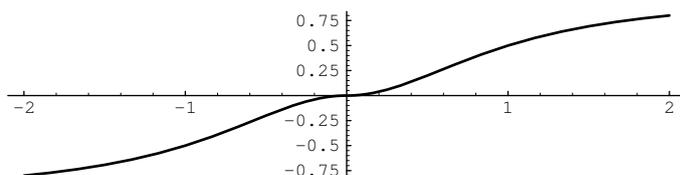
Si quitamos los doses obtenemos un homeomorfismo $t : \mathbb{R}^2 \rightarrow D$, donde D es la bola abierta (euclídea) de centro 0 y radio 1. Concretamente:

$$t(x, y) = \left(\frac{x \sqrt{x^2 + y^2}}{x^2 + y^2 + 1}, \frac{y \sqrt{x^2 + y^2}}{x^2 + y^2 + 1} \right).$$

Si restringimos esta aplicación a \mathbb{R} , es decir, a los puntos $(x, 0)$ obtenemos un homeomorfismo entre \mathbb{R} y el intervalo $] -1, 1[$. Concretamente

$$t(x) = \frac{x|x|}{x^2 + 1}.$$

He aquí su gráfica:



Si lo restringimos a $]0, +\infty[$ obtenemos un homeomorfismo entre este intervalo y $]0, 1[$. A saber:

$$t(x) = \frac{x^2}{x^2 + 1}.$$

A partir de aquí es fácil ver que dos intervalos abiertos cualesquiera (acotados o no) son homeomorfos. ■

Límites infinitos Consideremos ahora límites de funciones $f : X \rightarrow \mathbb{K}^\infty$, donde \mathbb{K} es un cuerpo métrico al que le hemos añadido un punto infinito, de modo que una base de entornos de ∞ la forman los conjuntos

$$\{\infty\} \cup \{\alpha \in \mathbb{K} \mid |\alpha| > K\}, \quad K \in \mathbb{R}.$$

Consideraremos también funciones $f : X \rightarrow \overline{\mathbb{R}} = \mathbb{R} \cup \{\pm\infty\}$, donde los entornos básicos de $+\infty$ son los intervalos $]K, +\infty]$ y los de $-\infty$ los intervalos $[-\infty, K[$, con $K \in \mathbb{R}$.

Es claro que

$$\lim_{x \rightarrow a} f(x) = +\infty$$

significa que para todo $K > 0$ existe un entorno U de a tal que si $x \in U$, $x \neq a$ está en el dominio de f , entonces $f(x) > K$. Similarmente ocurre con $-\infty$. Que el límite valga ∞ significa que para todo $K > 0$ existe un entorno U de a tal que si $x \in U$, $x \neq a$ está en el dominio de f , entonces $|f(x)| > K$.

Es fácil probar que los teoremas del estilo de “el límite de la suma es la suma de los límites” etc. valen también en los casos siguientes:

$$\begin{aligned} +\infty + (+\infty) &= +\infty, & -\infty + (-\infty) &= -\infty, \\ +\infty + a &= +\infty, & -\infty + a &= -\infty, & \infty + a &= \infty, \\ (+\infty)(+\infty) &= +\infty, & (-\infty)(-\infty) &= +\infty, & \infty \infty &= \infty. \\ \text{si } a > 0, & \quad (+\infty)a &= +\infty, & \quad (-\infty)a &= -\infty; \\ \text{si } a < 0, & \quad (+\infty)a &= -\infty, & \quad (-\infty)a &= +\infty. \\ \text{Si } a \neq 0, & \quad \infty a &= \infty, & \quad \frac{a}{\pm\infty} = \frac{a}{\infty} &= 0, & \quad \frac{a}{0} &= \infty. \end{aligned}$$

Por ejemplo, la igualdad $+\infty + (+\infty) = +\infty$ es una forma de expresar que la suma de dos funciones $f, g : A \subset X \rightarrow \overline{\mathbb{R}}$ que tienden a $+\infty$ en un punto $a \in A'$, tiende también a $+\infty$. La prueba es sencilla: dado un $K > 0$ existen entornos U y V de a tales que si $x \in U \cap A'$ entonces $f(x) > K/2$ y si $x \in V \cap A'$ entonces $g(x) > K/2$, con lo que si $x \in U \cap V \cap A'$ entonces $f(x) + g(x) > K$, luego se cumple la definición de límite. Similarmente se prueban todas las demás.

Cálculo de límites Estudiamos ahora los límites más simples y frecuentes. Comencemos con los límites en infinito de los polinomios en un cuerpo métrico \mathbb{K} (que será \mathbb{R} cuando distingamos entre $+\infty$ y $-\infty$). Obviamente,

$$\lim_{x \rightarrow \pm\infty} x = \pm\infty.$$

Aplicando n veces la regla $(+\infty)(+\infty) = +\infty$ vemos que si $n > 0$ entonces

$$\lim_{x \rightarrow +\infty} x^n = +\infty.$$

El límite en $-\infty$ será obviamente $(-1)^n \infty$. Si $n < 0$ usamos la regla $1/\pm\infty = 0$ y si $n = 0$ tenemos que x^0 es la constante 1, luego en total:

$$\lim_{x \rightarrow +\infty} x^n = \begin{cases} +\infty & \text{si } n > 0, \\ 1 & \text{si } n = 0, \\ 0 & \text{si } n < 0. \end{cases}$$

El límite en $-\infty$ difiere sólo en el primer caso, de forma obvia.

Para calcular el límite de un polinomio multiplicamos y dividimos por la potencia de x de mayor grado:

$$\lim_{x \rightarrow +\infty} 8x^5 - 3x^4 + x^3 + 2x - 5 = \lim_{x \rightarrow +\infty} x^5 \left(8 - \frac{3}{x} + \frac{1}{x^2} + \frac{2}{x^4} - \frac{5}{x^5} \right) = +\infty.$$

En general, si $p(x) \in \mathbb{R}[x]$ no es constante, se cumple $\lim_{x \rightarrow \pm\infty} p(x) = \pm\infty$, donde el signo depende en forma obvia del signo del coeficiente director de p y de la paridad del grado si el límite es en $-\infty$. Similarmente se razona que en \mathbb{K}^∞ todos los polinomios no constantes tienen límite ∞ .

En particular vemos que, a diferencia de casos como $1/\infty$ o $+\infty + \infty$, no hay una regla general para calcular un límite de la forma $+\infty - \infty$. En efecto, los tres límites siguientes son de este tipo y cada uno da un resultado distinto. Por ello se dice que $+\infty - \infty$ es una *indeterminación*.

$$\lim_{x \rightarrow +\infty} x^5 - x^2 = +\infty, \quad \lim_{x \rightarrow +\infty} x^2 - x^5 = -\infty, \quad \lim_{x \rightarrow +\infty} (x+2) - (x+1) = 1.$$

Nos ocupamos ahora de los límites de fracciones algebraicas (cocientes de polinomios). El ejemplo siguiente ilustra el caso general:

$$\lim_{x \rightarrow +\infty} \frac{6x^4 - 3x^3 + x + 1}{2x^4 + x^2 - 2x + 2} = \lim_{x \rightarrow +\infty} \frac{6 - \frac{3}{x} + \frac{1}{x^3} + \frac{1}{x^4}}{2 + \frac{1}{x^2} - \frac{2}{x^3} + \frac{2}{x^4}} = \frac{6}{2} = 3.$$

Es claro que el límite en $\pm\infty$ del cociente de dos polinomios del mismo grado es el cociente de los términos directores. Si los grados son distintos, al dividir entre la potencia de mayor grado obtenemos 0 si el grado del denominador es mayor y $\pm\infty$ si es mayor el del numerador, donde el signo se calcula de forma obvia. Esto vale igualmente (salvo lo dicho de los signos) para límites en \mathbb{K}^∞ . Vemos, pues, que el caso ∞/∞ es también una indeterminación.

Por ejemplo, antes hemos probado que la función

$$t(x) = \frac{x|x|}{x^2 + 1}$$

es un homeomorfismo entre \mathbb{R} y el intervalo $] -1, 1[$. Ahora es claro que

$$\lim_{x \rightarrow \pm\infty} t(x) = \pm 1,$$

luego si definimos $t(\pm\infty) = \pm 1$ tenemos una biyección continua $t: \overline{\mathbb{R}} \rightarrow [-1, 1]$. En el capítulo siguiente veremos que es un homeomorfismo.

Ejercicio: Calcular t^{-1} y probar que es continua.

2.7 Convergencia de sucesiones

En 1.10 hemos definido la convergencia de una sucesión en un espacio métrico. Ahora estamos en condiciones de probar que la convergencia de sucesiones es en realidad un concepto topológico y no métrico, es decir, que depende únicamente de la topología del espacio y no de la distancia. De hecho, el concepto de límite de una sucesión es un caso particular del concepto de límite de funciones que hemos estudiado en la sección anterior.

Definición 2.75 Sea X un espacio topológico, $l \in X$ y $\{x_n\}_{n=0}^\infty$ una sucesión en X . Diremos que $\{x_n\}_{n=0}^\infty$ converge a l si está finalmente en todo entorno de l , o sea, si para cada entorno V de l existe un $m \in \mathbb{N}$ tal que para todo $n \geq m$ se cumple $x_n \in V$.

Si X es un espacio de Hausdorff y $\{x_n\}_{n=0}^\infty$ converge en X , entonces el punto al cual converge es único, pues puntos distintos tienen entornos disjuntos, y una sucesión no puede estar finalmente en dos conjuntos disjuntos. Representaremos por

$$\lim_n x_n$$

al único límite de la sucesión $\{x_n\}_{n=0}^\infty$ en un espacio de Hausdorff, cuando éste exista.

Ahora observamos que si X es un espacio métrico, la definición 2.75 es equivalente a 1.10. En general, notemos que en la definición anterior no hace falta considerar todos los entornos V de l , sino que basta considerar los de una base de entornos de l . En el caso de un espacio métrico basta considerar las bolas $B_\epsilon(l)$, de modo que “para todo entorno V de l ” puede sustituirse por “para todo $\epsilon > 0$ ” y la condición $x_n \in V$ equivale a $d(x_n, l) < \epsilon$.

Por otra parte, 2.75 es un caso particular de 2.67. Para probarlo basta recordar que una sucesión $\{x_n\}_{n=0}^\infty$ en un espacio topológico X no es más que una aplicación $x : \mathbb{N} \rightarrow X$, y que podemos considerar a \mathbb{N} como subespacio de \mathbb{N}^∞ . Una base de entornos de ∞ en este espacio la forman los conjuntos

$$\{\infty\} \cup \{n \in \mathbb{N} \mid n \geq m\}, \quad m \in \mathbb{N},$$

y así, que $x : \mathbb{N} \rightarrow X$ converja a l cuando n tiende a ∞ en el sentido de 2.67 es precisamente lo que exige 2.75. Por lo tanto, en un espacio de Hausdorff tenemos que

$$\lim_n x_n = \lim_{n \rightarrow +\infty} x(n).$$

Sabiendo esto, todos los resultados generales para límites de funciones valen para sucesiones, por ejemplo, el límite de la suma de dos sucesiones de números reales es la suma de los límites, etc. Veamos ahora algunos hechos específicos sobre sucesiones:

Una sucesión converge si y sólo si converge finalmente, es decir, $\{x_n\}_{n=0}^\infty$ converge a l si y sólo si la sucesión $\{x_n\}_{n=k}^\infty$ converge a l . En particular, las sucesiones finalmente constantes convergen.

(Esto es consecuencia de que el límite de una función en un punto — el punto ∞ en este caso— depende sólo del comportamiento de la función en un entorno del punto).

Si $A \subset X$, $\{x_n\}_{n=0}^\infty \subset A$ y $l \in A$, entonces $\{x_n\}_{n=0}^\infty$ converge a l en A si y sólo si converge a l en X .

(Pues los entornos de l en A son las intersecciones con A de los entornos de l en X y, como la sucesión está en A , es equivalente que esté finalmente en un entorno de l en A o que esté finalmente en un entorno de l en X .)

Este hecho nos dice que la convergencia depende exclusivamente de la sucesión y de su límite, y no del espacio en el que los consideremos (siempre que no cambiemos de topología, claro está). Sin embargo, también de aquí se desprende que una sucesión convergente deja de serlo si eliminamos su límite. Por ejemplo, la sucesión $\{1/n\}_{n=0}^\infty$ no converge en el espacio $]0, 1]$, pues si convergiera a un punto l , tendríamos que en \mathbb{R} debería converger a la vez a l y a 0 .

Si una sucesión converge a un punto l , entonces todas sus subsucesiones convergen a l .

Pues si $\{x_n\}_{n=0}^\infty$ converge a l y $\{x_{n_k}\}_{k=0}^\infty$ es una subsucesión, para todo entorno U de l existe un n_0 tal que si $n \geq n_0$ se cumple $x_n \in U$, pero si $k \geq n_0$ entonces $n_k \geq k \geq n_0$, luego $x_{n_k} \in U$.

Ejercicio: Demostrar que una sucesión $\{x^n\}_{n=0}^\infty$ converge en un espacio producto $\prod_{i \in I} X_i$ a un punto x si y sólo si las sucesiones $\{x_i^n\}_{n=0}^\infty$ convergen a x_i para todo $i \in I$.

No podemos continuar nuestro estudio de las sucesiones sin introducir un nuevo concepto. Sucede que las sucesiones no se comportan adecuadamente en todos los espacios topológicos, sino sólo en aquellos que cumplen la siguiente propiedad adicional, por lo demás muy frecuente:

Definición 2.76 Un espacio X cumple el *primer axioma de numerabilidad* (1AN) si para todo punto $x \in X$ existe una base de entornos de x numerable, es decir, de la forma $\{V_n\}_{n=0}^\infty$.

Esta propiedad la tienen todos los espacios métricos, pues si x es un punto de un espacio métrico, una base de entornos de x la forman los conjuntos $\{B_{1/n}(x)\}_{n=1}^\infty$.

En $\overline{\mathbb{R}}$, una base de entornos de $+\infty$ es $\{]n, +\infty]\}_{n=0}^\infty$ y una base de entornos de $-\infty$ es $\{[-\infty, n[]\}_{n=0}^\infty$, luego este espacio también cumple 1AN.

Si M es un espacio métrico, el espacio M^∞ también es 1AN, pues, fijado $m_0 \in M$, una base de entornos de ∞ la forman los conjuntos

$$E_n = \{m \in M \mid d(m, m_0) > n\} \cup \{\infty\}, \quad \text{para } n = 1, 2, \dots$$

Así pues, todos los espacios topológicos que estamos considerando son 1AN.

Observemos que si X es un espacio que cumple 1AN y $x \in X$, podemos tomar una base de entornos de x de la forma $\{V_n\}_{n=0}^\infty$ que cumpla además

$$V_0 \supset V_1 \supset V_2 \supset V_3 \supset \dots$$

En efecto, si una base dada $\{W_n\}_{n=0}^\infty$ no cumple esta condición, la sustituimos por $V_n = W_0 \cap \dots \cap W_n$ y así tenemos las inclusiones indicadas.

Consideremos ahora esta sucesión:

$$1, \quad -1, \quad 1, \quad -1, \quad \dots$$

Es obvio que no es convergente, pero sin duda hay dos puntos “especiales” para esta sucesión, el 1 y el -1 . Quizá el lector se sienta tentado de afirmar que se trata de una sucesión con dos límites, pero nuestra definición de límite no admite esa posibilidad. Vamos a dar una definición que describa esta situación.

Definición 2.77 Un punto x de un espacio topológico X es un *punto adherente* de una sucesión $\{a_n\}_{n=0}^\infty$ en X si para cada entorno V de x y cada número natural n existe un número natural $m \geq n$ tal que $a_m \in V$.

Es decir, si la sucesión contiene puntos arbitrariamente lejanos en cualquier entorno de x o, si se prefiere, si la sucesión contiene infinitos puntos en cada entorno de x .

En estos términos, la sucesión $\{(-1)^n\}_{n=0}^\infty$ que hemos puesto como ejemplo tiene exactamente dos puntos adherentes, 1 y -1 .

Teorema 2.78 Sea X un espacio 1AN. Un punto $x \in X$ es un punto adherente de una sucesión $\{a_n\}_{n=0}^\infty$ si y sólo si ésta contiene una subsucesión que converge a x .

DEMOSTRACIÓN: Si x es un punto adherente de la sucesión, sea $\{V_n\}_{n=0}^\infty$ una base decreciente de entornos de x . Existe un punto $a_{n_0} \in V_0$. Existe un natural $n_1 \geq n_0 + 1$ tal que $a_{n_1} \in V_1$. Existe un natural $n_2 \geq n_1 + 1$ tal que $a_{n_2} \in V_2$. De este modo obtenemos una subsucesión $\{a_{n_k}\}_{k=0}^\infty$ tal que cada $a_{n_k} \in V_k$ y, como la sucesión de entornos es decreciente, en realidad V_k contiene todos los términos de la subsucesión posteriores a a_{n_k} , luego la subsucesión que hemos obtenido está finalmente en cada entorno V_k , es decir, converge a x .

Recíprocamente, si hay una subsucesión que converge a x , dicha subsucesión está finalmente en cualquier entorno de x , luego dicho entorno contiene infinitos términos de la sucesión dada. ■

En particular vemos que una sucesión convergente no tiene más punto adherente que su límite.

La continuidad de funciones puede caracterizarse en términos de sucesiones:

Teorema 2.79 Sea $f : X \rightarrow Y$ una aplicación entre espacios topológicos y supongamos que X cumple 1AN. Sea $x \in X$. Entonces f es continua en x si y sólo si para cada sucesión $\{a_n\}_{n=0}^\infty \subset X$ tal que $\lim_n a_n = x$, se cumple $\lim_n f(a_n) = f(x)$.

DEMOSTRACIÓN: Supongamos que f es continua en x . Sea V un entorno de $f(x)$. Entonces $f^{-1}[V]$ es un entorno de x y la sucesión $\{a_n\}_{n=0}^\infty$ está finalmente en $f^{-1}[V]$, luego $\{f(a_n)\}_{n=0}^\infty$ está finalmente en V , o sea, $\lim_n f(a_n) = f(x)$.

Recíprocamente, supongamos que f no es continua en x . Entonces existe un entorno V de $f(x)$ tal que $f^{-1}[V]$ no es entorno de x . Sea $\{V_n\}_{n=0}^\infty$ una base decreciente de entornos de x . Para cada natural n , no puede ocurrir que $V_n \subset f^{-1}[V]$, luego existe³ un punto $a_n \in V_n$ tal que $f(a_n) \notin V$.

Como la base es decreciente, todos los términos posteriores a a_n están en V_n , luego la sucesión $\{a_n\}_{n=0}^\infty$ está finalmente en cada V_n , con lo que converge a x . Sin embargo la sucesión $\{f(a_n)\}_{n=0}^\infty$ no tiene ningún término en V , luego no converge a $f(x)$. ■

Los puntos adherentes se caracterizan por sucesiones:

Teorema 2.80 Sea X un espacio topológico que cumpla 1AN. Sea $A \subset X$. Entonces \bar{A} está formado por los límites de las sucesiones convergentes contenidas en A .

DEMOSTRACIÓN: Si l es el límite de una sucesión contenida en A , entonces todo entorno de l contiene puntos de la sucesión, es decir, puntos de A , luego $l \in \bar{A}$.

Recíprocamente, si $x \in \bar{A}$, tomamos una base decreciente $\{V_n\}_{n=0}^\infty$ de entornos abiertos de x . Como $V_n \cap A \neq \emptyset$, existe un $a_n \in V_n \cap A$ y la sucesión $\{a_n\}_{n=0}^\infty$ así construida converge a x , y está contenida en A . ■

En particular, un conjunto A es cerrado si y sólo si el límite de toda sucesión convergente contenida en A , está en A , es decir, si no es posible “salir” de A mediante sucesiones.

2.8 Sucesiones y series numéricas

Ahora estamos en condiciones de continuar con más medios el estudio de la convergencia de sucesiones y series que iniciamos en la sección 1.2. Por ejemplo, allí vimos que si \mathbb{K} es un cuerpo métrico y $r \in \mathbb{K}$ cumple $|r| < 1$, entonces $\lim_n r^n = 0$. Ahora podemos extender este resultado. Consideremos primero el caso en que $r \in \mathbb{R}$ y $r > 1$. Entonces $s = r - 1 > 0$ y por el teorema del binomio de Newton, para $n \geq 1$, se cumple

$$r^n = (1 + s)^n = 1 + ns + \sum_{k=2}^n \binom{n}{k} s^k > ns.$$

³En este punto, al igual que en muchas situaciones similares, se usa el axioma de elección numerable.

Sabemos que $\lim_n ns = +\infty$. Del hecho de que $\{ns\}_{n=0}^\infty$ esté finalmente en cada entorno básico de $+\infty$, de la forma $]M, +\infty]$, se sigue claramente que lo mismo le sucede a $\{r^n\}_{n=0}^\infty$, luego si $r > 1$ concluimos que $\lim_n r^n = +\infty$.

De aquí obtenemos un argumento alternativo para el caso $0 < r < 1$. Basta observar que

$$\lim_n r^n = \lim_n \frac{1}{(1/r)^n} = \frac{1}{\infty} = 0.$$

Ahora consideramos el caso general en que \mathbb{K} es un cuerpo métrico y $r \in \mathbb{K}$. Entonces

$$\lim_n r^n = \begin{cases} 0 & \text{si } |r| < 1, \\ \infty & \text{si } |r| > 1. \end{cases}$$

En efecto, si $|r| > 1$, entonces $\lim_n |r^n| = \lim_n |r|^n = +\infty$, de donde es fácil deducir a partir de las meras definiciones que $\lim_n r^n = \infty$.

Si $|r| < 1$, entonces $\lim_n |r^n| = 0$, de donde también se sigue que $\lim_n r^n = 0$.

Puede probarse que si $|r| = 1$ el límite no existe salvo en el caso $r = 1$. Por ejemplo, ya hemos visto que $\lim_n (-1)^n$ no existe, pues se trata de la sucesión $1, -1, 1, -1, 1, -1, \dots$

Ahora observamos que las sucesiones monótonas siempre convergen:

Teorema 2.81 *Toda sucesión monótona creciente en \mathbb{R} converge a su supremo en $\overline{\mathbb{R}}$, y toda sucesión monótona decreciente converge a su ínfimo en $\overline{\mathbb{R}}$.*

DEMOSTRACIÓN: Por supremo de una sucesión $\{a_n\}_{n=0}^\infty$ entendemos el supremo del conjunto $\{a_n \mid n \in \mathbb{N}\}$. Sea s este supremo y supongamos que es finito. Entonces un entorno básico de s es de la forma $]s - \epsilon, s + \epsilon[$, para un $\epsilon > 0$. Como $s - \epsilon$ no es una cota superior de la sucesión, existe un natural m tal que $s - \epsilon < a_m \leq s$ y, por la monotonía, $s - \epsilon < a_n \leq s$ para todo $n \geq m$, es decir, que la sucesión está finalmente en el entorno. Una ligera modificación nos da el mismo resultado si $s = +\infty$ o si la sucesión es decreciente. ■

Ejemplo: Raíces anidadas Sea $a > 0$ y consideremos la sucesión definida recurrentemente como

$$x_0 = 0, \quad x_{n+1} = \sqrt{a + x_n}.$$

Así, por ejemplo,

$$x_5 = \sqrt{a + \sqrt{a + \sqrt{a + \sqrt{a + \sqrt{a}}}}}$$

Supongamos que existe $\lim_n x_n = c$. Entonces, como $x_{n+1}^2 = a + x_n$, el miembro izquierdo de la igualdad converge a c^2 , mientras que el derecho converge

a $a + c$, luego tiene que ser $c^2 = a + c$, lo que implica (teniendo en cuenta que tiene que ser $c \geq 0$) que $c = \frac{1 + \sqrt{1 + 4a}}{2}$.

Vamos a ver que, en efecto, la sucesión converge. Para ello probamos que es monótona creciente. Claramente $x_0 < x_1$ y, si se cumple que $x_n < x_{n+1}$, entonces $a + x_n < a + x_{n+1}$, luego $\sqrt{a + x_n} < \sqrt{a + x_{n+1}}$, luego $x_{n+1} < x_{n+2}$.

Por otra parte, vamos a probar que la sucesión está acotada superiormente por su presunto límite c . En efecto, es obvio que $x_0 < c$ y, si $x_n < c$, entonces $a + x_n < a + c$, luego $x_{n+1} = \sqrt{a + x_n} < \sqrt{a + c} = c$, donde la última igualdad se sigue que de $a + c = c^2$.

El teorema anterior demuestra que la sucesión converge, luego converge a c . Se suele expresar esto con la notación:

$$\sqrt{a + \sqrt{a + \sqrt{a + \sqrt{a + \sqrt{a + \dots}}}}} = \frac{1 + \sqrt{1 + 4a}}{2}.$$

En particular tenemos una expresión para el número áureo como límite de raíces anidadas:

$$\sqrt{1 + \sqrt{1 + \sqrt{1 + \sqrt{1 + \sqrt{1 + \dots}}}}} = \frac{1 + \sqrt{5}}{2}. \quad \blacksquare$$

Consideramos ahora series numéricas, es decir, sucesiones de la forma

$$\sum_{n=0}^{\infty} a_n = \left\{ \sum_{n=0}^k a_n \right\}_{k=0}^{\infty},$$

donde a_n varía en un cuerpo métrico \mathbb{K} . En 1.23 demostramos la convergencia de las series geométricas $\sum_{n=0}^{\infty} r^n$ con $|r| < 1$, cuya suma es $1/(1 - r)$.

A la hora de determinar si una serie $\sum_{n=0}^{\infty} a_n$ es convergente o divergente, un hecho básico es que para que converja su término general (es decir, la sucesión $\{a_n\}_{n=0}^{\infty}$) debe tender a 0. En efecto:

Teorema 2.82 *Sea $\{a_n\}_{n=0}^{\infty}$ una sucesión en un cuerpo métrico \mathbb{K} . Si la serie $\sum_{n=0}^{\infty} a_n$ es convergente, entonces $\lim_n a_n = 0$.*

DEMOSTRACIÓN: Sea $S_k = \sum_{n=0}^k a_n$. Que la serie converja a un número L significa por definición que existe $\lim_k S_k = L$. En tal caso también existe $\lim_k S_{k+1} = L$, pues es el límite de una subsucesión, y entonces

$$\lim_k a_k = \lim_k (S_{k+1} - S_k) = L - L = 0. \quad \blacksquare$$

Así, si no sumamos cada vez cantidades más pequeñas la serie no puede converger. El recíproco es tentador, pero falso. Basta considerar la serie determinada por la más sencilla de las sucesiones que tienden a 0:

Ejemplo La serie $\sum_{n=1}^{\infty} \frac{1}{n}$ es divergente en \mathbb{R} .

En efecto, observemos que

$$\begin{aligned} S_1 &= 1, \\ S_2 &= 1 + \frac{1}{2}, \\ S_4 &= S_2 + \frac{1}{3} + \frac{1}{4} > S_2 + \frac{2}{4} = 1 + \frac{1}{2} + \frac{1}{2}, \\ S_8 &= S_4 + \frac{1}{5} + \frac{1}{6} + \frac{1}{7} + \frac{1}{8} > S_4 + \frac{4}{8} > 1 + \frac{1}{2} + \frac{1}{2} + \frac{1}{2}. \end{aligned}$$

En general $S_{2^n} > 1 + \frac{n}{2}$, y esta sucesión tiende a $+\infty$, luego las sumas parciales no están acotadas y la serie diverge. ■

Se conocen muchos criterios para determinar el carácter convergente o divergente de una serie. Por ejemplo, el siguiente es aplicable a *series de términos positivos*, es decir, a series de números reales ≥ 0 :

Teorema 2.83 (Criterio de D'Alembert) Sea $\{a_n\}_{n=0}^{\infty}$ una sucesión de números reales positivos tal que exista $\lim_n \frac{a_{n+1}}{a_n} = L \in \mathbb{R}$. Entonces:

- a) Si $L < 1$ la serie $\sum_{n=0}^{\infty} a_n$ es convergente.
 b) Si $L > 1$ la serie $\sum_{n=0}^{\infty} a_n$ es divergente.

DEMOSTRACIÓN: a) Sea $\epsilon > 0$ tal que $L + \epsilon < 1$. Entonces $]-\infty, L + \epsilon[$ es un entorno de L , y por definición de límite existe un natural n_0 tal que si $n \geq n_0$ entonces $\frac{a_{n+1}}{a_n} < L + \epsilon$ o, lo que es lo mismo, $a_{n+1} < a_n(L + \epsilon)$. Así pues,

$$\begin{aligned} a_{n_0+1} &< a_{n_0}(L + \epsilon), \\ a_{n_0+2} &< a_{n_0+1}(L + \epsilon) < a_{n_0}(L + \epsilon)^2, \\ a_{n_0+3} &< a_{n_0+2}(L + \epsilon) < a_{n_0}(L + \epsilon)^3, \dots \end{aligned}$$

En general, si $n \geq n_0$, se cumple que

$$a_n < a_{n_0}(L + \epsilon)^{n-n_0} = \frac{a_{n_0}}{(L + \epsilon)^{n_0}}(L + \epsilon)^n.$$

De aquí que si $k > n_0$,

$$\begin{aligned} \sum_{n=0}^k a_n &\leq \sum_{n=0}^{n_0} a_n + \frac{a_{n_0}}{(L + \epsilon)^{n_0}} \sum_{n=n_0+1}^k (L + \epsilon)^n \\ &< \sum_{n=0}^{n_0} a_n + \frac{a_{n_0}}{(L + \epsilon)^{n_0}} \sum_{n=n_0+1}^{\infty} (L + \epsilon)^n < +\infty, \end{aligned}$$

donde la última serie converge porque es geométrica de razón $L + \epsilon < 1$. La serie $\sum_{n=0}^{\infty} a_n$ es de términos positivos y sus sumas parciales están acotadas, luego converge al supremo de dichas sumas parciales.

b) Sea ahora $\epsilon > 0$ tal que $1 < L - \epsilon$. Igual que en el apartado anterior existe un natural n_0 tal que si $n \geq n_0$ entonces $a_{n+1} > a_n(L - \epsilon)$, de donde se deduce igualmente que si $n \geq n_0$ entonces

$$a_n > \frac{a_{n_0}}{(L - \epsilon)^{n_0}} (L - \epsilon)^n,$$

y por consiguiente, para $k > n_0$,

$$\sum_{n=0}^k a_n \geq \sum_{n=0}^{n_0} a_n + \frac{a_{n_0}}{(L - \epsilon)^{n_0}} \sum_{n=n_0+1}^k (L - \epsilon)^n,$$

pero ahora la última serie diverge, pues es geométrica de razón mayor que 1, luego sus sumas parciales no están acotadas, y las de la primera serie tampoco. Así pues, ésta es divergente. ■

En el caso de que el límite L exista y valga 1 no es posible asegurar nada, hay casos en los que esto ocurre y la serie converge y casos en los que diverge.

Ejemplo Si $x > 0$, la serie $\sum_{n=0}^{\infty} \frac{x^n}{n!}$ es convergente, pues por el criterio de D'Alembert,

$$L = \lim_n \frac{x^{n+1}}{(n+1)!} : \frac{x^n}{n!} = \lim_n \frac{x}{n+1} = 0 < 1.$$

Incidentalmente, esto prueba también que $\lim_n \frac{x^n}{n!} = 0$.

Notemos que hemos determinado el carácter convergente de la serie, pero no hemos dicho nada sobre el cálculo efectivo de su límite. La razón es que no hay nada que decir. Por ejemplo, en el caso más simple, $x = 1$, el número

$$e = \sum_{n=0}^{\infty} \frac{1}{n!}$$

es un número "nuevo", en el sentido de que no es racional, ni la raíz cuadrada de un número racional, ni en general expresable en términos de otros números ya conocidos. El único sentido en que podemos "calcularlo" es en el de obtener aproximaciones racionales sumando términos de la serie. El resultado es

$$e = 2.7182818284590452353602874 \dots$$

■

Veamos ahora un criterio válido para las llamadas *series alternadas*, es decir, para series de números reales en las que los términos consecutivos tienen signos opuestos:

Teorema 2.84 (Criterio de Leibniz) Sea $\{a_n\}_{n=0}^{\infty}$ una sucesión de números reales positivos decreciente y convergente a 0. Entonces la serie $\sum_{n=0}^{\infty} (-1)^n a_n$ es convergente.

DEMOSTRACIÓN: Consideremos primero las sumas parciales pares. Por ejemplo:

$$S_6 = (a_0 - a_1) + (a_2 - a_3) + (a_4 - a_5) + a_6.$$

Teniendo en cuenta que la sucesión es decreciente, los sumandos así agrupados son todos mayores o iguales que 0, luego en general $S_{2n} \geq 0$.

Por otra parte, $S_8 = S_6 + (-a_7 + a_8) \leq S_6$, luego la sucesión $\{S_{2n}\}_{n=0}^{\infty}$ es monótona decreciente y acotada inferiormente por 0. Por lo tanto converge a un número L .

Ahora, $S_{2n+1} = S_{2n} + a_{2n+1}$, luego existe $\lim_n S_{2n+1} = L + 0$.

Es fácil comprobar que si las dos subsucesiones $\{S_{2n}\}_{n=0}^{\infty}$ y $\{S_{2n+1}\}_{n=0}^{\infty}$ convergen a un mismo número L , entonces toda la sucesión $\{S_n\}_{n=0}^{\infty}$ converge a L , es decir, la serie converge. ■

Por ejemplo, la serie $\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n}$ es convergente. De nuevo no tenemos más medio para calcular su límite que aproximarlos por una suma parcial. En realidad, cuando hacemos esto necesitamos saber cuál es el error cometido, para determinar el número de cifras decimales correctas. Por ejemplo, las primeras sumas parciales de esta serie son:

1	1	11	0.73654	21	0.71639	31	0.70901
2	0.50000	12	0.65321	22	0.67093	32	0.67776
3	0.83333	13	0.73013	23	0.71441	33	0.70806
4	0.58333	14	0.65870	24	0.67274	34	0.67865
5	0.78333	15	0.72537	25	0.71274	35	0.70722
6	0.61666	16	0.66287	26	0.67428	36	0.67945
7	0.75952	17	0.72169	27	0.71132	37	0.70647
8	0.63452	18	0.66613	28	0.67560	38	0.68016
9	0.74563	19	0.71877	29	0.71009	39	0.70580
10	0.64563	20	0.66877	30	0.67675	40	0.68080

El límite vale 0.6931471805599453094172321... Por lo tanto, si al calcular la décima suma afirmáramos que el límite vale 0.6456... estaríamos cometiendo un grave error. En realidad sólo la primera cifra decimal es exacta (se dice que un número decimal finito aproxima a otro con n cifras exactas si las n primeras cifras de ambos números coinciden). En general, al aproximar una serie por una suma parcial, o al aproximar cualquier límite de una sucesión por uno de sus términos, no sabemos cuál es el error cometido, ni en particular cuántas de las cifras de la aproximación son exactas. Sin embargo, en el caso de las series alternadas es fácil saberlo pues, según hemos visto, las sumas pares son decrecientes (siempre están sobre el límite) y las impares son crecientes (siempre están bajo el límite), luego las últimas sumas de la tabla nos dicen que el límite

se encuentra entre 0.68 y 0.71 y por lo tanto no sabemos ninguna de sus cifras con exactitud (salvo el 0). Yendo más lejos tenemos:

$$S_{1000} = 0.69264 \quad \text{y} \quad S_{1001} = 0.69364,$$

lo que nos permite afirmar que el límite está entre 0.692 y 0.694, es decir, es de la forma 0.69... y ya tenemos dos cifras exactas.

En la práctica no nos ocuparemos del problema técnico de determinar las cifras exactas que nos proporciona una aproximación dada, y consideraremos que un número es “conocido” si tenemos una sucesión que converge a él. Todas las aproximaciones que daremos (calculadas con ordenador con técnicas de análisis numérico) tendrán todas sus cifras exactas.

Estudiamos ahora las series sobre un cuerpo métrico completo \mathbb{K} . En primer lugar, la condición de Cauchy para una serie en \mathbb{K} puede expresarse como sigue:

Teorema 2.85 Si \mathbb{K} es un cuerpo métrico completo, una serie $\sum_{n=0}^{\infty} a_n$ en \mathbb{K} es convergente si y sólo si para todo $\epsilon > 0$ existe un natural n_0 tal que si $n_0 \leq m \leq p$, entonces $\left| \sum_{n=m}^p a_n \right| < \epsilon$.

DEMOSTRACIÓN: Se trata de la condición de Cauchy, pues $\sum_{n=m}^p a_n$ es la diferencia entre la suma parcial p -ésima menos la suma parcial $m-1$ -ésima, y su valor absoluto es la distancia entre ambas. ■

De aquí se sigue un hecho importantísimo.

Teorema 2.86 Sea $\sum_{n=0}^{\infty} a_n$ una serie en un cuerpo métrico completo \mathbb{K} . Si la serie $\sum_{n=0}^{\infty} |a_n|$ es convergente, entonces la serie $\sum_{n=0}^{\infty} a_n$ también lo es.

DEMOSTRACIÓN: Dado $\epsilon > 0$ existe un número natural n_0 de manera que si $n_0 \leq m \leq p$, entonces $\sum_{n=m}^p |a_n| < \epsilon$.

Ahora bien, $\left| \sum_{n=m}^p a_n \right| \leq \sum_{n=m}^p |a_n| < \epsilon$, luego la serie sin valores absolutos también es de Cauchy, luego converge. ■

Definición 2.87 Una serie $\sum_{n=0}^{\infty} a_n$ en un cuerpo métrico completo \mathbb{K} es *absolutamente convergente* si la serie $\sum_{n=0}^{\infty} |a_n|$ es convergente en \mathbb{R} .

Hemos probado que toda serie absolutamente convergente es convergente. Las series convergentes que no son absolutamente convergentes se llaman series *condicionalmente convergentes*.

Un ejemplo de serie condicionalmente convergente es

$$1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \frac{1}{5} - \frac{1}{6} + \frac{1}{7} - \frac{1}{8} + \cdots = 0.693147 \dots$$

La convergencia absoluta de una serie es esencial para ciertas cuestiones. Por ejemplo, una consecuencia inmediata de las propiedades de las sumas finitas y de los límites de sucesiones es que

$$\sum_{n=0}^{\infty} a_n + \sum_{n=0}^{\infty} b_n = \sum_{n=0}^{\infty} (a_n + b_n), \quad a \sum_{n=0}^{\infty} a_n = \sum_{n=0}^{\infty} aa_n,$$

entendiendo que si las series de la izquierda convergen, las de la derecha también lo hacen y se da la igualdad. Un resultado análogo para producto de series ya no es tan sencillo:

Definición 2.88 Sean $\sum_{n=0}^{\infty} a_n$ y $\sum_{n=0}^{\infty} b_n$ series en un cuerpo métrico completo \mathbb{K} .

Llamaremos *producto de Cauchy* de estas series a la serie $\sum_{n=0}^{\infty} \left(\sum_{k=0}^n a_k \cdot b_{n-k} \right)$.

La intención es que la serie que acabamos de definir converja al producto de las dos series de partida, pero esto no ocurre necesariamente si al menos una de ellas no converge absolutamente.

Teorema 2.89 Si $\sum_{n=0}^{\infty} a_n$ y $\sum_{n=0}^{\infty} b_n$ son dos series convergentes en un cuerpo métrico completo \mathbb{K} , al menos una de las cuales converge absolutamente, entonces

$$\sum_{n=0}^{\infty} \left(\sum_{k=0}^n a_k \cdot b_{n-k} \right) = \left(\sum_{n=0}^{\infty} a_n \right) \left(\sum_{n=0}^{\infty} b_n \right).$$

DEMOSTRACIÓN: Supongamos que la serie que converge absolutamente es $\sum_{n=0}^{\infty} a_n$ y definamos

$$A = \sum_{n=0}^{\infty} a_n, \quad B = \sum_{n=0}^{\infty} b_n, \quad c_n = \sum_{k=0}^n a_k \cdot b_{n-k}, \quad C_n = \sum_{k=0}^n c_k,$$

$$A_n = \sum_{k=0}^n a_k, \quad B_n = \sum_{k=0}^n b_k, \quad \beta_n = B_n - B.$$

Ahora,

$$\begin{aligned} C_n &= c_0 + \cdots + c_n = a_0 b_0 + (a_0 b_1 + a_1 b_0) + \cdots + (a_0 b_n + \cdots + a_n b_0) \\ &= a_0 B_n + \cdots + a_n B_0 = a_0 (B + \beta_n) + \cdots + a_n (B + \beta_0) \\ &= A_n B + (a_0 \beta_n + \cdots + a_n \beta_0) \end{aligned}$$

El teorema quedará probado si vemos que $a_0 \beta_n + \cdots + a_n \beta_0$ tiende a 0.

Sea $\epsilon > 0$. Sea $K = \sum_{n=0}^{\infty} |a_n|$. Sea $M = \sup\{|\beta_n| \mid n \geq 0\}$ (la sucesión β_n tiende a 0, luego está acotada). Existe un número natural n_0 tal que si $n \geq n_0$, entonces $|\beta_n| < \epsilon/2K$ y $\sum_{k=n_0+1}^n |a_k| < \epsilon/2M$. En consecuencia, si $n \geq 2n_0$,

$$\begin{aligned} |a_0\beta_n + \cdots + a_n\beta_0| &\leq \sum_{k=0}^n |a_k\beta_{n-k}| = \sum_{k=0}^{n_0} |a_k\beta_{n-k}| + \sum_{k=n_0+1}^n |a_k\beta_{n-k}| \\ &< \frac{\epsilon}{2K} \sum_{k=0}^{n_0} |a_k| + M \sum_{k=n_0+1}^n |a_k| \leq \frac{\epsilon}{2K} K + \frac{\epsilon}{2M} M = \epsilon. \end{aligned}$$

■

Si ninguna de las series converge absolutamente el resultado no tiene por qué cumplirse.

Ejemplo Consideremos la serie

$$\sum_{n=0}^{\infty} \frac{(-1)^n}{\sqrt{n+1}}.$$

La serie converge por el criterio de Leibniz. El producto de Cauchy de esta serie por sí misma tiene término general

$$c_n = (-1)^n \sum_{k=0}^n \frac{1}{\sqrt{(n-k+1)(k+1)}}.$$

Cuando $n \geq k$ tenemos

$$(n-k+1)(k+1) = \left(\frac{n}{2}+1\right)^2 - \left(\frac{n}{2}-k\right)^2 \leq \left(\frac{n}{2}+1\right)^2,$$

luego

$$\frac{1}{\sqrt{(n-k+1)(k+1)}} \geq \frac{1}{\frac{n}{2}+1} = \frac{2}{n+2}.$$

Por consiguiente

$$|c_n| \geq \sum_{k=0}^n \frac{2}{n+2} = \frac{2(n+1)}{n+2}.$$

Esta expresión converge a 2, luego c_n no converge a 0 y el producto de Cauchy no converge. El teorema anterior prueba, pues, que la serie dada es condicionalmente convergente. ■

Otro punto en el que la convergencia absoluta resulta crucial es en el de la reordenación de los términos de una serie.

Dada una serie convergente $\sum_{n=0}^{\infty} a_n$ y una biyección $\sigma : \mathbb{N} \rightarrow \mathbb{N}$, podemos considerar la serie $\sum_{n=0}^{\infty} a_{\sigma(n)}$ y estudiar su convergencia. De nuevo el resultado natural exige que la serie converja absolutamente:

Teorema 2.90 Si $\sum_{n=0}^{\infty} a_n$ es una serie absolutamente convergente en un cuerpo métrico completo \mathbb{K} y $\sigma : \mathbb{N} \rightarrow \mathbb{N}$ es una aplicación biyectiva, entonces la serie $\sum_{n=0}^{\infty} a_{\sigma(n)}$ es absolutamente convergente y tiene la misma suma.

DEMOSTRACIÓN: Una serie es absolutamente convergente si y sólo si las sumas parciales de sus módulos forman un conjunto acotado. Toda suma parcial de los módulos de la reordenación está mayorada por una suma parcial de los módulos de la serie original (tomando los sumandos necesarios para incluir todos los que aparecen en la suma dada). Por tanto las sumas parciales de los módulos de la reordenación están acotadas y la serie converge absolutamente.

Sea $\epsilon > 0$. Existe un número natural n_0 tal que si $n \geq n_0$

$$\left| \sum_{k=0}^n a_k - \sum_{k=0}^{\infty} a_k \right| = \left| \sum_{k=n+1}^{\infty} a_k \right| \leq \sum_{k=n+1}^{\infty} |a_k| = \sum_{k=0}^{\infty} |a_k| - \sum_{k=0}^n |a_k| < \frac{\epsilon}{2}.$$

Sea $m_0 \geq n_0$ tal que $\{0, 1, \dots, n_0\} \subset \{\sigma(0), \sigma(1), \dots, \sigma(m_0)\}$. Entonces si $n \geq m_0$,

$$\begin{aligned} \left| \sum_{k=0}^n a_{\sigma(k)} - \sum_{k=0}^{\infty} a_k \right| &\leq \left| \sum_{k=0}^n a_{\sigma(k)} - \sum_{k=0}^{n_0} a_k \right| + \left| \sum_{k=0}^{n_0} a_k - \sum_{k=0}^{\infty} a_k \right| \\ &< \sum_{k=n_0+1}^{\infty} |a_k| + \frac{\epsilon}{2} < \epsilon. \end{aligned}$$

Por lo tanto $\sum_{k=0}^{\infty} a_{\sigma(k)} = \sum_{k=0}^{\infty} a_k$. ■

Nota Observemos que si I es cualquier conjunto infinito numerable y $\{a_i\}_{i \in I}$ es cualquier familia de elementos de \mathbb{K} , la serie $\sum_{n=0}^{\infty} a_{i_n}$ que resulta de fijar una enumeración i_0, i_1, \dots de I es absolutamente convergente si y sólo si las sumas $\sum_{i \in F} |a_i|$, con $F \subset I$ finito, están acotadas.

En efecto, si las sumas están acotadas, la serie $\sum_{n=0}^{\infty} |a_{i_n}|$ tiene las sumas parciales acotadas, luego es convergente. Recíprocamente, si la serie de los valores absolutos es convergente, su suma acota a todas las sumas finitas.

En estas condiciones, la suma infinita $\sum_{i \in I} a_i$ puede definirse como la suma de la serie determinada por cualquier enumeración del conjunto I o, alternativamente, como el único número $S \in \mathbb{K}$ que cumple la condición siguiente (que no involucra la elección de ninguna enumeración):

Para todo $\epsilon > 0$ existe $F_0 \subset I$ finito tal que, para todo $F_0 \subset F \subset I$, se cumple $\left| S - \sum_{i \in F} a_i \right| < \epsilon$.

En efecto, si la serie converge a S al enumerar I , basta fijar una enumeración y tomar como $F_0 = \{i_0, \dots, i_{n_0}\}$ los primeros términos de la sucesión, de modo que

$$\left| \sum_{n=n_0+1}^{\infty} a_{i_n} \right| < \frac{\epsilon}{2}, \quad \sum_{n=n_0+1}^{\infty} |a_{i_n}| < \frac{\epsilon}{2}.$$

Entonces, si $F_0 \subset F \subset I$, con F finito:

$$\left| S - \sum_{i \in F} a_i \right| \leq \left| \sum_{n=0}^{\infty} a_{n_i} - \sum_{n=0}^{n_0} a_{i_n} \right| + \left| \sum_{i \in F_0} a_i - \sum_{i \in F} a_i \right| < \frac{\epsilon}{2} + \sum_{i \in F \setminus F_0} |a_i| < \epsilon.$$

Recíprocamente, si se cumple esta condición y fijamos una enumeración de I , para cada $\epsilon > 0$ tomamos F_0 según la hipótesis y elegimos un n_0 tal que $F_0 \subset \{i_0, \dots, i_{n_0}\}$. Así, para todo $n_1 \geq n_0$ se cumple que

$$\left| S - \sum_{n=0}^{n_1} a_i \right| < \epsilon,$$

luego la serie converge a S . ■

Las series absolutamente convergentes se pueden manipular exactamente igual que si fueran sumas finitas. El siguiente teorema justifica cualquier operación razonable entre ellas.

Teorema 2.91 *Sea $\{a_i\}_{i \in I}$ una familia numerable de elementos de un cuerpo métrico completo \mathbb{K} . Sea $I = \bigcup_{n=0}^{\infty} I_n$ una división de I en partes disjuntas. Entonces $\sum_{i \in I} a_i$ es (absolutamente) convergente si y sólo si lo son las series $\sum_{i \in I_n} a_i$ y $\sum_{n=0}^{\infty} \sum_{i \in I_n} |a_i|$. Además en tal caso*

$$\sum_{i \in I} a_i = \sum_{n=0}^{\infty} \sum_{i \in I_n} a_i.$$

DEMOSTRACIÓN: Si $\sum_{i \in I} a_i$ es absolutamente convergente, sus sumas parciales en módulo están acotadas, pero toda suma parcial en módulo de cada $\sum_{i \in I_n} |a_i|$ lo es también de la primera, luego éstas están acotadas, o sea, las series $\sum_{i \in I_n} a_i$ convergen absolutamente. Dado cualquier natural k , tomamos para cada $n \leq k$ un conjunto finito $F_n \subset I_n$ tal que $\sum_{i \in I_n} |a_i| - \sum_{i \in F_n} |a_i| < 1/(k+1)$. Entonces

$$\sum_{n=0}^k \sum_{i \in I_n} |a_i| < \sum_{n=0}^k \sum_{i \in F_n} |a_i| + 1 \leq \sum_{i \in I} |a_i| + 1,$$

luego las sumas parciales están acotadas y así todas las series convergen absolutamente.

Supongamos ahora que las series $\sum_{i \in I_n} |a_i|$ y $\sum_{n=0}^{\infty} \sum_{i \in I_n} |a_i|$ convergen absolutamente. Si $F \subset I$ es finito, para un cierto k suficientemente grande se cumple

$$\sum_{i \in F} |a_i| = \sum_{n=0}^k \sum_{i \in I_n \cap F} |a_i| \leq \sum_{n=0}^k \sum_{i \in I_n} |a_i| \leq \sum_{n=0}^{\infty} \sum_{i \in I_n} |a_i|,$$

luego las sumas parciales de $\sum_{i \in I} |a_i|$ están acotadas y la serie converge absolutamente.

Ahora supongamos la convergencia de todas las series y probemos la igualdad de las sumas. Notemos que la serie

$$\sum_{n=0}^{\infty} \sum_{i \in I_n} a_i$$

es convergente porque es absolutamente convergente. Sea $\epsilon > 0$. Existe un $F_* \subset I$ finito tal que si $F_* \subset F \subset I$ es finito, entonces

$$\left| \sum_{i \in I} a_i - \sum_{i \in F} a_i \right| < \frac{\epsilon}{4}.$$

Existe un número natural n_0 tal que $F_* \subset \bigcup_{n=0}^{n_0} I_n$ y

$$\left| \sum_{n=n_0+1}^{\infty} \sum_{i \in I_n} a_i \right| < \frac{\epsilon}{4}.$$

Para cada $n \leq n_0$ existe un conjunto finito $F_n \subset I_n$ tal que si $F_n \subset F \subset I_n$, entonces

$$\left| \sum_{i \in I_n} a_i - \sum_{i \in F} a_i \right| < \frac{\epsilon}{2(n_0 + 1)}.$$

Sea $F = F_* \cup \bigcup_{n=0}^{n_0} F_n$, de modo que

$$\begin{aligned} \left| \sum_{n=0}^{\infty} \sum_{i \in I_n} a_i - \sum_{i \in I} a_i \right| &\leq \left| \sum_{n=n_0+1}^{\infty} \sum_{i \in I_n} a_i \right| + \left| \sum_{n=0}^{n_0} \sum_{i \in I_n} a_i - \sum_{i \in F} a_i \right| \\ &+ \left| \sum_{i \in F} a_i - \sum_{i \in I} a_i \right| \\ &< \frac{\epsilon}{4} + \left| \sum_{n=0}^{n_0} \left(\sum_{i \in I_n} a_i - \sum_{i \in I_n \cap F} a_i \right) \right| + \frac{\epsilon}{4} < \frac{\epsilon}{2} + \frac{\epsilon}{2} = \epsilon. \end{aligned}$$

Por lo tanto ambas sumas coinciden. ■

Una aplicación destacada del teorema anterior es esta variante de 2.89:

Teorema 2.92 Si $\sum_{i \in I} a_i$ y $\sum_{j \in J} b_j$ son dos series (absolutamente) convergentes en un cuerpo métrico completo \mathbb{K} , entonces

$$\sum_{(i,j) \in I \times J} a_i b_j = \left(\sum_{i \in I} a_i \right) \left(\sum_{j \in J} b_j \right).$$

DEMOSTRACIÓN: Es claro que las series $\sum_{j \in J} |a_i b_j|$ son convergentes, al igual que lo es

$$\sum_{i \in I} \sum_{j \in J} |a_i b_j| = \sum_{i \in I} |a_i| \left(\sum_{j \in J} |b_j| \right),$$

pues en ambos casos estamos multiplicando por una constante los términos de una serie convergente. Por el teorema anterior, aplicado a $I \times J = \bigcup_{i \in I} \{i\} \times J$, concluimos que

$$\sum_{(i,j) \in I \times J} a_i b_j = \sum_{i \in I} \left(\sum_{j \in J} a_i b_j \right) = \sum_{i \in I} a_i \left(\sum_{j \in J} b_j \right) = \left(\sum_{i \in I} a_i \right) \left(\sum_{j \in J} b_j \right). \quad \blacksquare$$

Ejemplo Consideremos la serie condicionalmente convergente

$$S = \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n}.$$

Entonces

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{2n} = \frac{S}{2}.$$

Ahora consideremos la serie $\sum_{n=1}^{\infty} a_n$ cuyos términos impares son ceros y sus términos pares son los de la serie anterior, es decir, la serie

$$0 + \frac{1}{2} + 0 - \frac{1}{4} + 0 + \frac{1}{6} + 0 - \frac{1}{8} + \dots$$

Obviamente su suma es también $S/2$. La serie

$$\sum_{n=1}^{\infty} \left(\frac{(-1)^{n+1}}{n} + a_n \right)$$

converge a $3S/2$. Sus primeros términos son:

$$1 + 0 + \frac{1}{3} - \frac{1}{2} + \frac{1}{5} + 0 + \frac{1}{7} - \frac{1}{4} + \frac{1}{9} + 0 + \dots$$

Eliminando los ceros obtenemos la serie

$$1 + \frac{1}{3} - \frac{1}{2} + \frac{1}{5} + \frac{1}{7} - \frac{1}{4} + \frac{1}{9} + \frac{1}{11} - \frac{1}{6} \dots$$

Vemos que se trata de una reordenación de la serie original, pero converge a $3S/2$. \blacksquare

Capítulo III

Compacidad, conexión y completitud

Finalmente tenemos los suficientes elementos de topología como para que ésta se convierta en una herramienta eficaz. En los temas anteriores apenas hemos extraído las consecuencias más elementales de las definiciones de topología, continuidad, etc. Los resultados que veremos ahora van mucho más lejos y dan una primera muestra de las posibilidades de las técnicas topológicas.

3.1 Espacios compactos

La compacidad es en topología una propiedad similar a la “dimensión finita” en álgebra lineal. Los espacios compactos no son necesariamente finitos, pero se comportan en muchos aspectos como si lo fueran. Por ejemplo, es obvio que en un espacio finito toda sucesión ha de tomar infinitas veces un mismo valor, luego toda sucesión contiene una subsucesión constante, en particular convergente. Análogamente, toda sucesión en $\overline{\mathbb{R}}$ tiene una subsucesión convergente, pues por el teorema 1.7 toda sucesión tiene una subsucesión monótona, que por 2.81 converge a su supremo o a su ínfimo. En cambio, \mathbb{R} no tiene esta propiedad, pues la sucesión de los números naturales no contiene ninguna subsucesión convergente (ya que cualquier subsucesión converge a $+\infty$ en $\overline{\mathbb{R}}$, luego no converge en \mathbb{R}).

Veremos que esta diferencia equivale a que $\overline{\mathbb{R}}$ es un espacio compacto, mientras que \mathbb{R} no lo es. La propiedad de las subsucesiones convergentes caracteriza la compacidad en espacios métricos, pero para el caso general necesitamos otra definición más elaborada.

Definición 3.1 Sea X un espacio topológico. Un *cubrimiento abierto* de X es una familia $\{A_i\}_{i \in I}$ de abiertos de X tal que $X = \bigcup_{i \in I} A_i$.

Un *subcubrimiento* del cubrimiento dado es un cubrimiento formado por parte de los abiertos del primero.

Un espacio de Hausdorff K es *compacto* si de todo cubrimiento abierto de K se puede extraer un subcubrimiento finito.

Es obvio que si X es un espacio finito, de todo cubrimiento abierto se puede extraer un subcubrimiento finito. Basta tomar un abierto que contenga a cada uno de los puntos del espacio. Así pues, todo espacio de Hausdorff finito es compacto.

Observemos que si \mathcal{B} es una base de un espacio de Hausdorff K , se cumple que K es compacto si y sólo si todo cubrimiento de K por abiertos básicos admite un subcubrimiento finito. En efecto, si $\{A_i\}_{i \in I}$ es un cubrimiento arbitrario, el conjunto C de los elementos de \mathcal{B} contenidos en algún A_i es también un cubrimiento abierto, ya que para todo $x \in X$ existe un i tal que $x \in A_i$, y por definición de base existe un $B \in \mathcal{B}$ tal que $x \in B \subset A_i$, luego $x \in B \in C$.

Por hipótesis existe un subcubrimiento finito $\{B_1, \dots, B_n\}$, y para cada índice $j = 1, \dots, n$ existe un $i_j \in I$ tal que $B_j \subset A_{i_j}$, luego

$$K = B_1 \cup \dots \cup B_n \subset A_{i_1} \cup \dots \cup A_{i_n} \subset K.$$

Una familia de abiertos forma un cubrimiento si y sólo si la familia de sus complementarios es una familia de cerrados con intersección vacía. Por ello la compacidad puede caracterizarse así en términos de familias de cerrados:

Un espacio de Hausdorff K es compacto si y sólo si toda familia de cerrados $\{C_i\}_{i \in I}$ con la propiedad de que cualquier intersección finita de ellos no es vacía, tiene intersección total no vacía.

La propiedad de que las intersecciones finitas sean no vacías se llama *propiedad de la intersección finita*. Por lo tanto:

Teorema 3.2 *Un espacio de Hausdorff K es compacto si y sólo si toda familia de cerrados de K con la propiedad de la intersección finita tiene intersección no vacía.*

A menudo nos encontraremos con espacios que no son compactos, pero tienen subespacios compactos. Por ello resulta útil caracterizar la compacidad de un subespacio en términos de la topología de todo el espacio y no de la topología relativa. Concretamente:

Teorema 3.3 *Sea X un espacio de Hausdorff y K un subespacio de X . Entonces K es compacto si y sólo si para toda familia $\{A_i\}_{i \in I}$ de abiertos (básicos) de X tal que $K \subset \bigcup_{i \in I} A_i$ se puede extraer una subfamilia finita que cumpla lo mismo.*

DEMOSTRACIÓN: Supongamos que K es compacto. Entonces $\{A_i \cap K\}_{i \in I}$ es claramente un cubrimiento abierto de K , del que podemos extraer un subcubrimiento finito de modo que

$$K = (A_{i_1} \cap K) \cup \dots \cup (A_{i_n} \cap K),$$

luego $K \subset A_{i_1} \cup \dots \cup A_{i_n}$.

Recíprocamente, si K cumple esta propiedad y $\{A_i\}_{i \in I}$ es un cubrimiento abierto de K , entonces para cada i existe un abierto B_i de X tal que $A_i = B_i \cap K$ (concretamente, podemos elegir B_i como la unión de todos los abiertos de X con esta propiedad). Consecuentemente $K = \bigcup_{i \in I} A_i \subset \bigcup_{i \in I} B_i$, luego por hipótesis podemos tomar un número finito de conjuntos de modo que $K \subset B_{i_1} \cup \dots \cup B_{i_n}$, luego $K = (B_{i_1} \cap K) \cup \dots \cup (B_{i_n} \cap K) = A_{i_1} \cup \dots \cup A_{i_n}$. Así pues, K es compacto. ■

Si la unión de una familia de abiertos de un espacio X contiene a un subespacio K , diremos que forma un *cubrimiento abierto* de K en X . Así pues, un subespacio K de X es compacto si y sólo si de todo cubrimiento abierto de K en X puede extraerse un subcubrimiento finito (en X también).

Aquí estamos considerando la topología de X , pero deberemos tener siempre presente que la compacidad es una propiedad absoluta, y depende exclusivamente de la topología del propio espacio K .

Los teoremas siguientes muestran la anunciada similitud entre los espacios compactos y los espacios finitos. Por lo pronto, todo subespacio finito es cerrado. El análogo con compactos es el siguiente:

Teorema 3.4 *Se cumplen las propiedades siguientes:*

- a) Si X es un espacio de Hausdorff y $K \subset X$ es compacto, entonces K es cerrado en X .
- b) Si K es un compacto y $C \subset K$ es un cerrado, entonces C es compacto.
- c) Si M es un espacio métrico y $K \subset M$ es compacto, entonces K está acotado.

DEMOSTRACIÓN: a) El argumento es el mismo que emplearíamos si K fuera finito. Veamos que $X \setminus K$ es abierto. Sea $x \in X \setminus K$ y sea C la familia de todos los abiertos que son disjuntos de algún entorno de x . El hecho de que X sea un espacio de Hausdorff se traduce en que todo para todo $u \in K$ existen abiertos disjuntos A y B tales que $u \in A$ y $x \in B$, luego $A \in C$, y esto prueba que C es un cubrimiento abierto de K . Por lo tanto existe un subcubrimiento finito $K \subset A_1 \cup \dots \cup A_n$, de modo que para cada A_i existe un abierto B_i disjunto de A_i con $x \in B_i$. Finalmente, $\bigcap_{i=1}^n B_i$ es un entorno de x que no corta a K , luego $X \setminus K$ es un entorno de x .

b) Si $\{A_i\}_{i \in I}$ es un cubrimiento abierto de C , entonces $\{A_i\}_{i \in I} \cup \{K \setminus C\}$ es un cubrimiento abierto de K , luego existe un subcubrimiento finito

$$K = A_{i_1} \cup \dots \cup A_{i_n} \cup (K \setminus C).$$

Claramente entonces $C \subset A_{i_1} \cup \dots \cup A_{i_n}$, luego C es compacto.

c) Sea $x \in M$ un punto cualquiera. Para cada $u \in K$, sea $r_u = d(x, u) + 1$. Obviamente $K \subset \bigcup_{u \in K} B_{r_u}(x)$. Por compacidad podemos extraer un subcobrimiento finito de modo que $K \subset B_{r_{u_1}}(x) \cup \dots \cup B_{r_{u_n}}(x)$. Las bolas son conjuntos acotados, una unión finita de acotados es acotada y todo subconjunto de un acotado está acotado. Por tanto K está acotado. ■

Teorema 3.5 *Si K es un espacio compacto, toda sucesión en K posee un punto adherente. Por tanto si además K cumple 1AN, toda sucesión en K tiene una subsucesión convergente.*

DEMOSTRACIÓN: Sea $\{a_n\}_{n=0}^\infty$ una sucesión en K . Sea $A_n = \{a_m \mid m \geq n\}$. Obviamente

$$A_0 \supset A_1 \supset A_2 \supset A_3 \supset A_4 \supset \dots ,$$

luego también

$$\bar{A}_0 \supset \bar{A}_1 \supset \bar{A}_2 \supset \bar{A}_3 \supset \bar{A}_4 \supset \dots ,$$

y así tenemos una familia de cerrados con la propiedad de la intersección finita. Por compacidad existe un punto $x \in \bigcap_{i=0}^\infty \bar{A}_i$. Obviamente x es un punto adherente de la sucesión, pues si n es un natural y U es un entorno de x , entonces $x \in \bar{A}_n$, luego $U \cap A_n \neq \emptyset$, es decir, existe un $m \geq n$ tal que $a_m \in U$. ■

Como habíamos anunciado, esta propiedad caracteriza a los espacios métricos compactos.

Teorema 3.6 *Un espacio métrico M es compacto si y sólo si toda sucesión en M tiene una subsucesión convergente.*

DEMOSTRACIÓN: Supongamos que M no fuera compacto. Entonces existiría un cubrimiento abierto $M = \bigcup_{i \in I} A_i$ que no admite subcobrimientos finitos.

Sea $\epsilon > 0$ y $x_0 \in M$. Si $B_\epsilon(x_0) \neq M$, existe un punto $x_1 \in M$ tal que $d(x_1, x_0) \geq \epsilon$. Si $B_\epsilon(x_0) \cup B_\epsilon(x_1) \neq M$, existe un punto $x_2 \in M$ tal que $d(x_2, x_0) \geq \epsilon$, $d(x_2, x_1) \geq \epsilon$.

Si M no pudiera cubrirse por un número finito de bolas de radio ϵ , podríamos construir una sucesión $\{x_n\}_{n=0}^\infty$ con la propiedad de que $d(x_i, x_j) \geq \epsilon$ para todos los naturales i, j . Es claro que tal sucesión no puede tener subsucesiones convergentes, pues una bola de centro el límite y radio $\epsilon/2$ debería contener infinitos términos de la sucesión, que distarían entre sí menos de ϵ .

Concluimos que para todo $\epsilon > 0$ existen puntos $x_0, \dots, x_n \in M$ de modo que $M = B_\epsilon(x_0) \cup \dots \cup B_\epsilon(x_n)$.

Lo aplicamos a $\epsilon = 1$ y obtenemos tales bolas. Si todas ellas pudieran cubrirse con un número finito de abiertos A_i también M podría, luego al menos una de ellas, digamos $B_1(x_0)$ no es cubrible por un número finito de abiertos del cubrimiento.

Igualmente, con $\epsilon = 1/2$ obtenemos una bola $B_{1/2}(x_1)$ no cubrible por un número finito de abiertos del cubrimiento. En general obtenemos una sucesión de bolas $B_{1/(n+1)}(x_n)$ con esta propiedad.

Sea x un punto adherente de la sucesión $\{x_n\}_{n=0}^\infty$. Sea $i \in I$ tal que $x \in A_i$. Como A_i es un abierto existe un número natural k tal que $B_{2/(k+1)}(x) \subset A_i$. Sea $n > k$ tal que $d(x_n, x) < 1/(k+1)$. Entonces $B_{1/(n+1)}(x_n) \subset B_{2/(k+1)}(x) \subset A_i$, en contradicción con que $B_{1/(n+1)}(x_n)$ no era cubrible con un número finito de abiertos del cubrimiento. ■

Ahora tenemos probado que $\overline{\mathbb{R}}$ es compacto, por el argumento dado al principio de esta sección. De aquí obtenemos a su vez una caracterización de los subespacios compactos de \mathbb{R} :

Teorema 3.7 *Un subconjunto de \mathbb{R} es compacto si y sólo si es cerrado y acotado.*

DEMOSTRACIÓN: Por el teorema 3.4, todo compacto en \mathbb{R} ha de ser cerrado y acotado. Si C es un conjunto cerrado y acotado, toda sucesión en C tiene una subsucesión convergente en $\overline{\mathbb{R}}$. Como C es acotado su límite estará en \mathbb{R} y como es cerrado, su límite estará en C , luego toda sucesión en C tiene una subsucesión convergente en C . Por el teorema anterior C es compacto. ■

También se puede probar el teorema anterior viendo que los cerrados y acotados de \mathbb{R} son precisamente los subconjuntos cerrados de $\overline{\mathbb{R}}$ (o, si se prefiere, esto es consecuencia inmediata del teorema anterior).

Teorema 3.8 (Teorema de Tychonoff) (AE)¹ *El producto de espacios topológicos compactos es compacto.*

DEMOSTRACIÓN: Sea $K = \prod_{i \in I} K_i$ un producto de espacios compactos. Tomemos una familia \mathcal{B} de cerrados en K con la propiedad de la intersección finita. Hemos de probar que su intersección es no vacía. El conjunto de todas las familias de subconjuntos no vacíos de K (no necesariamente cerrados) que contienen a \mathcal{B} y tienen la propiedad de la intersección finita, parcialmente ordenado por la inclusión, satisface las hipótesis del lema de Zorn, lo que nos permite tomar una familia maximal \mathcal{U} . Entonces $\bigcap_{U \in \mathcal{U}} U \subset \bigcap_{B \in \mathcal{B}} B$, luego basta probar que la primera intersección es no vacía.

En primer lugar observamos que si un conjunto $A \subset K$ corta a todos los elementos de \mathcal{U} entonces está en \mathcal{U} , pues en caso contrario $\mathcal{U} \cup \{A\}$ contradiría la maximalidad de \mathcal{U} .

Sea $p_i : K \rightarrow K_i$ la proyección en el factor i -ésimo. Es fácil ver que la familia $\{p_i[U] \mid U \in \mathcal{U}\}$ tiene la propiedad de la intersección finita luego, por la compacidad de K_i , existe un punto $x_i \in K_i$ tal que $x_i \in \overline{p_i[U]}$ para todo $U \in \mathcal{U}$. Estos puntos determinan un punto $x \in K$. Basta probar que $x \in \bigcap_{U \in \mathcal{U}} U$.

Fijemos un entorno básico de x , de la forma $A = \bigcap_{i \in F} p_i^{-1}[G_i]$, donde $F \subset I$ es finito y G_i es abierto en K_i . Para cada $U \in \mathcal{U}$ tenemos que $x_i \in \overline{p_i[U]}$,

¹No se requiere AE para probar que el producto de un número finito de espacios compactos es compacto, y si el producto es numerable basta con el Principio de Elecciones Dependientes. Remitimos a [TC] (teorema 8.67 y las observaciones posteriores) para la justificación de estos hechos.

luego $G_i \cap p_i[U] \neq \emptyset$, luego $p_i^{-1}[G_i] \cap U \neq \emptyset$. Como esto es cierto para todo $U \in \mathcal{U}$, según hemos observado antes podemos concluir que $p_i^{-1}[G_i] \in \mathcal{U}$, para todo $i \in F$. Como \mathcal{U} tiene la propiedad de la intersección finita, $A \in \mathcal{U}$. De aquí se sigue que A corta a todo $U \in \mathcal{U}$ y, como A es un entorno básico de x , esto implica que $x \in \overline{U}$ para todo $U \in \mathcal{U}$. ■

Ahora podemos reconocer fácilmente los subconjuntos compactos de \mathbb{R}^n :

Teorema 3.9 *Un subconjunto de \mathbb{R}^n , es compacto si y sólo si es cerrado y acotado.*

DEMOSTRACIÓN: La acotación depende en principio de la distancia que consideremos. Hemos de entender que se trata de la inducida por cualquiera de las tres normas definidas en el capítulo II. Por el teorema 2.3, todas tienen los mismos acotados. Trabajaremos concretamente con

$$\|x\|_\infty = \text{máx}\{|x_i| \mid i = 1, \dots, n\}.$$

Ya hemos visto que un compacto ha de ser cerrado y acotado. Supongamos que K es un subconjunto de \mathbb{R}^n cerrado y acotado. Esto significa que existe un $M > 0$ tal que para todo punto $x \in C$ se cumple $\|x\| \leq M$, lo que significa que si $x \in C$, cada $x_i \in [-M, M]$ o, de otro modo, que $C \subset [-M, M]^n$.

Pero por el teorema anterior $[-M, M]^n$ es compacto y C es cerrado en él, luego C también es compacto. ■

Una de las propiedades más importantes de la compacidad es que se conserva por aplicaciones continuas (compárese con el hecho de que la imagen (continua) de un conjunto finito es finita):

Teorema 3.10 *La imagen de un espacio compacto por una aplicación continua es de nuevo un espacio compacto (supuesto que sea un espacio de Hausdorff).*

DEMOSTRACIÓN: Sea $f : K \rightarrow X$ continua y suprayectiva. Supongamos que K es compacto y que X es un espacio de Hausdorff. Si $\{A_i\}_{i \in I}$ es un cubrimiento abierto de X , entonces $\{f^{-1}[A_i]\}_{i \in I}$ es un cubrimiento abierto de K , luego admite un subcubrimiento finito $K = f^{-1}[A_{i_1}] \cup \dots \cup f^{-1}[A_{i_n}]$. Entonces $X = A_{i_1} \cup \dots \cup A_{i_n}$. ■

Este hecho tiene muchas consecuencias:

Teorema 3.11 *Si $f : K \rightarrow X$ es biyectiva y continua, K es compacto y X es un espacio de Hausdorff, entonces f es un homeomorfismo.*

DEMOSTRACIÓN: Basta probar que la inversa es continua, o sea, que transforma cerrados de K en cerrados de X , o equivalentemente, que si C es cerrado en K , entonces $f[C]$ es cerrado en X , pero es que C es compacto, luego $f[C]$ también lo es, luego es cerrado. ■

Ejemplo Una circunferencia es homeomorfa a un cuadrado.

En efecto, si C es un cuadrado de centro $(0, 0)$ en \mathbb{R}^2 , es claro que la aplicación de C en la circunferencia unidad dada por $x \mapsto x/\|x\|$ es biyectiva y continua y, como C es compacto, es un homeomorfismo. ■

Otro hecho obvio es que toda aplicación continua de un compacto a un espacio métrico está acotada. Para las funciones reales podemos decir más:

Teorema 3.12 Si $f : K \rightarrow \mathbb{R}$ es continua y K es un compacto no vacío, existen $u, v \in K$ tales que para todo $x \in K$, se cumple $f(u) \leq f(x) \leq f(v)$. Es decir, que f alcanza un valor mínimo y un valor máximo.

DEMOSTRACIÓN: Sea $C = f[K]$. Entonces C es cerrado y acotado. Sean m y M su ínfimo y su supremo, respectivamente. Así para todo $x \in K$ se cumple que $m \leq f(x) \leq M$. Sólo falta probar que m y M son imágenes de puntos de K , o sea, que $m, M \in C$. Veámoslo para M .

Si $\epsilon > 0$, entonces $M - \epsilon$ no es una cota superior de C , luego existe un punto $y \in C$ de modo que $M - \epsilon < y$, es decir, que $]M - \epsilon, M + \epsilon[\cap C \neq \emptyset$. Esto significa que todo entorno (básico) de M corta a C , o sea, $M \in \overline{C} = C$. ■

Observemos que este resultado es falso sin compacidad. Por ejemplo la función $f :]0, 1[\rightarrow \mathbb{R}$ dada por $f(x) = x$ no tiene máximo ni mínimo.

3.2 Espacios localmente compactos

Hemos visto que \mathbb{R}^n no es compacto, pero cumple una propiedad relacionada con la compacidad que conviene destacar:

Definición 3.13 Un espacio topológico de Hausdorff es *localmente compacto* si todo punto tiene un entorno compacto.

Así, todo espacio compacto es localmente compacto. El teorema siguiente implica en particular que, en un espacio localmente compacto, todo punto tiene, de hecho, una base de entornos compactos:

Teorema 3.14 Sea X un espacio localmente compacto y $K \subset V \subset X$, con K compacto y V abierto. Entonces existe un abierto W tal que $K \subset W \subset \overline{W} \subset V$ y \overline{W} es compacto.

DEMOSTRACIÓN: Lo probamos primero para el caso en que $K = \{x\}$ se reduce a un punto. Sea C un entorno compacto de p . Para cada $y \in C \setminus V$, tomamos abiertos disjuntos $x \in U_y, y \in V_y$. Como $C \setminus V$ es compacto, existen y_1, \dots, y_n tales que $C \setminus V \subset V_{y_1} \cup \dots \cup V_{y_n}$, luego $U = U_{y_1} \cap \dots \cap U_{y_n}$ cumple que $U \subset C \setminus (V_{y_1} \cup \dots \cup V_{y_n})$, luego $K \subset \overline{U} \subset C \setminus (V_{y_1} \cup \dots \cup V_{y_n}) \subset V$.

En el caso general, para cada $x \in K$ existe un abierto W_x de clausura compacta tal que $x \in W_x \subset \overline{W_x} \subset V$. Los abiertos W_x cubren a K . Tomamos un subcobrimiento finito y llamamos W a la unión de sus miembros. ■

Por ejemplo, si $x \in \mathbb{R}^n$, las bolas cerradas $B'_\epsilon(x)$ forman una base de entornos compactos de x , luego \mathbb{R}^n es localmente compacto.

Todo cuanto digamos sobre \mathbb{R}^n se aplica en particular al cuerpo \mathbb{C} de los números complejos, que no es sino $\mathbb{C} = \mathbb{R}^2$, al que consideraremos siempre como espacio topológico con la topología inducida por el valor absoluto, que no es sino la norma euclídea de \mathbb{R}^2 . En particular, tenemos que \mathbb{C} es localmente compacto.

Una de las características de los espacios localmente compactos es que se pueden “compactificar” añadiendo un punto. En efecto:

La compactificación de Alexandroff Si X es un espacio topológico localmente compacto no compacto y ∞ es cualquier conjunto que no pertenezca a X , definimos la *compactificación de Alexandroff* de X como el espacio $X^\infty = X \cup \{\infty\}$ con la topología cuyos abiertos son los de X más los de la forma $(X \setminus K) \cup \{\infty\}$, donde $K \subset X$ es compacto.

Veamos que se trata ciertamente de una topología. Claramente \emptyset y X^∞ son abiertos. Si \mathcal{F} es una familia de abiertos en X^∞ , la unión $V = \bigcup \mathcal{F}$ es abierta, pues si $\infty \notin V$ entonces es una unión de abiertos de X , y $\infty \in U \in \mathcal{F}$, entonces $U = (X \setminus K) \cup \{\infty\}$, con K compacto, y si expresamos $V = (X \setminus K') \cup \{\infty\}$, entonces $X \setminus K'$ es abierto en X , porque es unión de los abiertos $W \setminus \{\infty\}$, donde W recorre \mathcal{F} , y además $X \setminus K \subset X \setminus K'$, luego $K' \subset K$, luego K' es compacto por ser cerrado en un compacto. Así pues, V es abierto en X^∞ .

Si U_1, \dots, U_n son abiertos en X^∞ y $V = U_1 \cap \dots \cap U_n$, entonces

$$V \setminus \{\infty\} = (U_1 \setminus \{\infty\}) \cap \dots \cap (U_n \setminus \{\infty\})$$

es abierto en X , por ser intersección de abiertos. Si $\infty \notin V$, entonces V es abierto en X^∞ . En caso contrario podemos expresarlo como $V = (X \setminus K) \cup \{\infty\}$, y hay que probar que K es compacto. Pero entonces $U_i = (X \setminus K_i) \cup \{\infty\}$, donde cada K_i es compacto, y $K = K_1 \cup \dots \cup K_n$. Pero es fácil ver (usando el teorema 3.3) que una unión finita de subespacios compactos es compacta.

Además X^∞ es un espacio compacto. En efecto, como X es un espacio de Hausdorff, es claro que dos puntos distintos de X tienen entornos disjuntos en X , luego en X^∞ . Además, si $x \in X$, entonces tiene un entorno compacto K , luego podemos tomar un abierto $x \in U \subset K$, y entonces $V = (X \setminus K) \cup \{\infty\}$ es un entorno de ∞ disjunto de U . Esto prueba que X^∞ es un espacio de Hausdorff.

Además, si $\{A_i\}_{i \in I}$ es un cubrimiento abierto de X^∞ , existe un $i_0 \in I$ tal que $\infty \in A_{i_0}$. Por lo tanto, $A_{i_0} = (X \setminus K) \cup \{\infty\}$, con $K \subset X$ compacto. Entonces $\{A_i \setminus \{\infty\}\}_{i \in I}$ es un cubrimiento abierto de K , luego podemos extraer un subcubrimiento finito $\{A_{i_k} \setminus \{\infty\}\}_{k=1}^n$, y entonces es claro que

$$X^\infty = A_{i_0} \cup \dots \cup A_{i_n}.$$

Notemos que si $X = \mathbb{R}^n$, entonces el espacio X^∞ construido en la página 57 es el mismo que acabamos de construir ahora, pues allí usábamos conjuntos

cerrados y acotados donde ahora hemos usado conjuntos compactos, y ambos conceptos coinciden en \mathbb{R}^n . En particular hemos visto que $\mathbb{C}^\infty = (\mathbb{R}^2)^\infty$ es homeomorfo a una esfera y que \mathbb{R}^∞ es homeomorfo a una circunferencia. ■

Teorema 3.15 *Se cumple:*

- a) *Todo subespacio abierto o cerrado de un espacio localmente compacto es localmente compacto.*
- b) *Un espacio es localmente compacto si y sólo si es un subespacio abierto de un espacio compacto.*

DEMOSTRACIÓN: a) Sea X un espacio localmente compacto y consideremos un subespacio $Y \subset X$. Si Y es cerrado y $p \in Y$, entonces p tiene una base \mathcal{B} de entornos compactos en X , y es claro que los conjuntos $\{K \cap Y \mid K \in \mathcal{B}\}$ son una base de entornos compactos de p en Y .

Si Y es abierto, $p \in Y$ y U es un entorno abierto de p en Y , entonces U es abierto en X , luego existe un entorno compacto K de p en Y tal que $p \in K \subset U \subset Y$, pero es claro que K es también entorno de p en Y , luego hemos probado que los entornos compactos de p en Y son una base de entornos de p . Así pues, Y es localmente compacto en ambos casos.

b) Si un espacio X es abierto en un espacio compacto, entonces es localmente compacto por a), y si X es localmente compacto, entonces, o bien es compacto y es abierto en sí mismo, o bien no es compacto, y es abierto en el espacio compacto X^∞ . ■

Terminamos esta sección demostrando unos resultados que necesitaremos más adelante:

Si X es un espacio localmente compacto y $f : X \rightarrow \mathbb{R}$, llamaremos *soporte* de f a la clausura del conjunto de puntos donde f toma valores $\neq 0$. Llamaremos $C_c(X)$ al conjunto de las aplicaciones continuas $f : X \rightarrow \mathbb{R}$ con soporte compacto. Es claro que se trata de un subespacio vectorial de $C(X)$.

Usaremos las notaciones $K \prec f$ y $f \prec V$ para indicar que $f : X \rightarrow [0, 1]$, $f \in C_c(X)$, K es compacto, V es abierto, f toma el valor 1 en K y f toma el valor 0 en $X \setminus V$.

He aquí el primero de los resultados que necesitaremos en el capítulo IX: afirma que siempre podemos tomar una función continua que valga 1 en cualquier compacto prefijado, pero que descienda rápidamente hasta valer 0 fuera de cualquier entorno prefijado de dicho compacto.

Teorema 3.16 (Lema de Urysohn) *Sea X un espacio localmente compacto, sea V un abierto y $K \subset V$ compacto. Entonces existe $f \in C_c(X)$ tal que $K \prec f \prec V$.*

DEMOSTRACIÓN: Por el teorema 3.14 existe un abierto $K \subset V_0 \subset V$ de clausura compacta. A su vez, existe un abierto W_0 de clausura compacta tal que $K \subset W_0 \subset \overline{W_0} \subset V_0$. Sea $W_1 = X$. Aplicamos de nuevo el teorema a $\overline{W_0} \subset V_0$, con lo que obtenemos un abierto $W_{1/2}$ de clausura compacta tal que

$$K \subset W_0 \subset \overline{W_0} \subset W_{1/2} \subset \overline{W_{1/2}} \subset V_0 \subset W_1.$$

Dos nuevas aplicaciones del mismo teorema nos dan

$$K \subset W_0 \subset \overline{W_0} \subset W_{1/4} \subset \overline{W_{1/4}} \subset W_{1/2} \subset \overline{W_{1/2}} \subset W_{3/4} \subset \overline{W_{3/4}} \subset V_0 \subset W_1.$$

Inductivamente vamos obteniendo una familia de abiertos $\{W_r\}_{r \in R}$, donde $R = \{k/2^i \mid i \in \mathbb{N}, 0 \leq k \leq 2^i\}$, de manera que si $r < r' < 1$ son puntos de R , entonces $K \subset \overline{W_r} \subset W_{r'} \subset V_0 \subset V$.

Definimos $g : X \rightarrow [0, 1]$ mediante $g(x) = \inf\{r \in R \mid x \in W_r\}$. Así $g[K] = \{0\}$ porque $K \subset W_0$ y $g[X \setminus V_0] = \{1\}$ porque $W_r \cap (X \setminus V_0) = \emptyset$ si $r < 1$. Basta probar que g es continua, pues entonces $f = 1 - g$ cumple lo pedido (su soporte está en $\overline{V_0}$, luego es compacto).

Sea $x \in X$ y $\epsilon > 0$. Si $g(x) \neq 0$ y $g(x) \neq 1$, entonces existen $r, r' \in R$ tales que $g(x) - \epsilon < r < g(x) < r' < g(x) + \epsilon$, luego $U = W_{r'} - \overline{W_r}$ es un entorno de x que cumple

$$g[U] \subset [r, r'] \subset]g(x) - \epsilon, g(x) + \epsilon[.$$

Si $g(x) = 0$ tomamos $0 = g(x) < r < g(x) + \epsilon$ y $U = W_r$ cumple lo mismo. Si $g(x) = 1$ tomamos $g(x) - \epsilon < r < g(x)$ y el U buscado es $X \setminus \overline{W_r}$. En cualquier caso obtenemos la continuidad de g en x . ■

El resultado siguiente es una generalización del anterior que, como veremos, nos permitirá “pegar” varias funciones continuas en una sola:

Teorema 3.17 *Si X es un espacio localmente compacto, V_1, \dots, V_n son abiertos de X y $K \subset V_1 \cup \dots \cup V_n$ es compacto, entonces existen funciones $h_i \prec V_i$ tales que $h_1(x) + \dots + h_n(x) = 1$ para todo $x \in K$.*

Se dice que las funciones h_i forman una *partición de la unidad* subordinada a los abiertos dados.

DEMOSTRACIÓN: Dado $x \in K$ existe un i tal que $x \in V_i$. Existe un abierto W_x de clausura compacta tal que $x \in W_x \subset \overline{W_x} \subset V_i$. Los abiertos W_x cubren a K . Extraemos un subcubrimiento finito y llamamos H_i a la unión de todos los abiertos del subcubrimiento cuya clausura está en V_i . De este modo los H_i son abiertos de clausura compacta que cubren a K y $\overline{H_i} \subset V_i$. Por el teorema anterior existen funciones $\overline{H_i} \prec V_i$. Definimos

$$h_1 = g_1, \quad h_2 = (1 - g_1)g_2, \quad \dots \quad h_n = (1 - g_1)(1 - g_2) \cdots (1 - g_{n-1})g_n.$$

Es claro que $h_i \prec V_i$ y una simple inducción prueba que

$$h_1 + \dots + h_n = 1 - (1 - g_1) \cdots (1 - g_n).$$

Es claro entonces que la suma vale 1 sobre los puntos de K , pues una de las funciones g_i ha de tomar el valor 1. ■

3.3 Espacios conexos

Pensemos en los espacios siguientes: $[0, 1]$ y $[0, 1] \cup [2, 3]$. Hay una diferencia esencial entre ellos, y es que el primero está formado por “una sola pieza” mientras que el segundo consta de “dos piezas”. La diferencia no es conjuntista, pues también podemos dividir $[0, 1] = [0, 1/2] \cup]1/2, 1]$, pero esto no son dos piezas en el mismo sentido que en el caso de $[0, 1] \cup [2, 3]$. La diferencia es que los intervalos $[0, 1/2]$ y $]1/2, 1]$ están “pegados” mientras que los intervalos $[0, 1]$ y $[2, 3]$ están “separados”. Con más precisión, el punto $1/2$ está sólo en uno de los intervalos, el $[0, 1/2]$, pero aunque no está en el otro, está pegado a él, en el sentido de que está en su clausura.

En general, si un espacio X se expresa como $X = U \cup V$, donde U y V son disjuntos y no vacíos, podemos decir que U y V son dos “piezas” en el sentido que estamos considerando si U no contiene puntos de la clausura de V y viceversa. Ahora bien, cualquier punto de \overline{V} que no estuviera en V debería estar en U , luego la condición equivale a que $V = \overline{V}$ y $U = \overline{U}$, o sea, a que U y V sean cerrados. Por otra parte, dado que U y V son complementarios, es lo mismo decir que son cerrados o que son abiertos. Con ello llegamos a la definición de conexión:

Definición 3.18 Un espacio topológico X es *disconexo* si existen subconjuntos abiertos U y V en X tales que $X = U \cup V$, $U \cap V = \emptyset$ y $U \neq \emptyset \neq V$. En caso contrario X es *conexo*.

Según hemos dicho, es indistinto exigir que U y V sean abiertos como que sean cerrados, pues de hecho si cumplen esto son a la vez abiertos y cerrados. Por lo tanto, un espacio X es conexo si y sólo si sus únicos subconjuntos que son a la vez abiertos y cerrados son X y \emptyset .

Es obvio que $[0, 1] \cup [2, 3]$, o incluso $[0, 1/2[\cup]1/2, 1]$ son ejemplos de espacios desconexos. Notar que $[0, 1/2[$ no es cerrado en \mathbb{R} , pero sí lo es en el espacio $[0, 1/2[\cup]1/2, 1]$ (su clausura en este espacio es la intersección con él de su clausura en \mathbb{R} , que es $[0, 1/2]$, o sea, es $[0, 1/2[$).

Es importante tener claro que los intervalos $[0, 1/2[$ y $]1/2, 1]$ están separados pese a que sólo falta un punto entre ellos. La falta de ese punto es suficiente para que ambas partes no se puedan “comunicar”, en el sentido de que, por ejemplo, ninguna sucesión contenida en una de las piezas puede converger a un punto de la otra. Esto es suficiente para que ambas partes sean independientes topológicamente. Así, la función $f : [0, 1/2[\cup]1/2, 1] \rightarrow \mathbb{R}$ dada por

$$f(x) = \begin{cases} 1 & \text{si } x \in [0, 1/2[, \\ 2 & \text{si } x \in]1/2, 1] \end{cases}$$

es continua, mientras que sería imposible definir una función continua sobre $[0, 1]$ que sólo tomara los valores 1 y 2.

Si la desconexión de estos espacios es clara, no lo es tanto la conexión de espacios como $[0, 1]$.

Ejercicio: Probar que el intervalo $[0, 1] \subset \mathbb{Q}$ es desconexo.

Teorema 3.19 *Un subconjunto de $\overline{\mathbb{R}}$ es conexo si y sólo si es un intervalo.*

DEMOSTRACIÓN: Sea C un subespacio conexo de $\overline{\mathbb{R}}$. Sean a y b su ínfimo y su supremo, respectivamente. Vamos a probar que C es uno de los cuatro intervalos de extremos a y b . Por el teorema 1.34 basta ver que si $a < x < b$ entonces $x \in C$. En caso contrario los conjuntos $C \cap [-\infty, x[$ y $C \cap]x, +\infty]$ son dos abiertos disjuntos no vacíos de C cuya unión es C .

Tomemos ahora un intervalo I y veamos que es conexo. Supongamos que existen abiertos disjuntos no vacíos U y V en I de modo que $I = U \cup V$. Tomemos $x \in U$ e $y \in V$. Podemos suponer que $x < y$.

Como I es un intervalo, $[x, y] \subset I$ y $U' = U \cap [x, y]$, $V' = V \cap [x, y]$ son abiertos disjuntos no vacíos en $[x, y]$ de modo que $[x, y] = U' \cup V'$.

Sea s el supremo de U' . Entonces $s \in \overline{U'} \cap [x, y] = U'$, luego en particular $s < y$. Claramente $]s, y[\subset V'$, luego $s \in \overline{V'} \cap [x, y] = V'$, contradicción. ■

Una consecuencia de esto es que un intervalo $[a, b[$ no es homeomorfo a uno de tipo $]c, d[$. En efecto, si eliminamos un punto de un intervalo $]c, d[$ nos queda un espacio desconexo, mientras que en $[a, b[$ podemos eliminar el punto a y obtenemos un conexo. (Si fueran homeomorfos, el espacio que resultara de eliminar la imagen de a en $]c, d[$ debería ser homeomorfo a $]a, b[$).

Ejercicio: Probar que dos intervalos (acotados o no acotados) son homeomorfos si y sólo si son del mismo tipo: abierto $]a, b[$, cerrado $[a, b]$ o semiabierto $[a, b[$.

Los resultados siguientes permiten probar con facilidad la conexión de muchos espacios. El primero refleja el hecho de que las aplicaciones continuas pueden pegar pero nunca cortar.

Teorema 3.20 *Las imágenes continuas de los espacios conexos son conexas.*

DEMOSTRACIÓN: Si $f : X \rightarrow Y$ es una aplicación continua y suprayectiva pero Y no es conexo, entonces X tampoco puede serlo, pues si A es un abierto cerrado no vacío en Y y distinto de Y , entonces $f^{-1}[A]$ cumple lo mismo en X . ■

Teorema 3.21 *Sea $\{A_i\}_{i \in I}$ una familia de subespacios conexos de un espacio X tal que $\bigcap_{i \in I} A_i \neq \emptyset$. Entonces $\bigcup_{i \in I} A_i$ es conexo.*

DEMOSTRACIÓN: Supongamos que $\bigcup_{i \in I} A_i = U \cup V$, donde U y V son abiertos disjuntos. Entonces para un i cualquiera se tendrá que $A_i = (U \cap A_i) \cup (V \cap A_i)$, pero $U \cap A_i$, $V \cap A_i$ son abiertos disjuntos en A_i , luego uno de ellos es vacío, y así $A_i \subset U$ o bien $A_i \subset V$.

Pero si $A_i \subset U$, entonces U contiene a $\bigcap_{i \in I} A_i$, luego U corta a todos los A_i y por conexión los contiene a todos. Así $\bigcup_{i \in I} A_i = U$, y $V = \emptyset$. Igualmente, si $A_i \subset V$ se deduce que U es vacío. ■

Ejemplo Las circunferencias son conexas.

Sea $f : [-1, 1] \rightarrow \{(x, y) \in \mathbb{R}^2 \mid x^2 + y^2 = 1, y \geq 0\}$ la aplicación dada por $f(x) = (x, \sqrt{1 - x^2})$.

Claramente f es continua y suprayectiva, lo que prueba que la semicircunferencia es conexa. Igualmente se prueba que la semicircunferencia opuesta es conexa, y como ambas se cortan en los puntos $(\pm 1, 0)$, su unión, es decir, la circunferencia, es conexa. ■

Teorema 3.22 Si A es un subespacio conexo de un espacio X , entonces \bar{A} es conexo.

DEMOSTRACIÓN: Supongamos que $\bar{A} = U \cup V$, donde U y V son abiertos disjuntos en \bar{A} . Entonces $A = (U \cap A) \cup (V \cap A)$, y $U \cap A$, $V \cap A$ son abiertos disjuntos en A . Por conexión uno es vacío, luego $A \subset U$ o bien $A \subset V$. Digamos $A \subset U \subset \bar{A}$.

Pero U es cerrado en \bar{A} , luego $\bar{A} \subset \bar{U} = U$, es decir, $U = \bar{A}$ y $V = \emptyset$. Esto prueba que \bar{A} es conexo. ■

Hemos dicho que un espacio desconexo es un espacio formado por varias “piezas” ahora podemos dar una definición rigurosa de lo que entendemos por una “pieza”.

Definición 3.23 Sea X un espacio topológico y $x \in X$. Llamaremos *componente conexa* de x a la unión $C(x)$ de todos los subconjuntos conexos de X que contienen a x . Por el teorema 3.21, $C(x)$ es un conexo, el mayor subespacio conexo de X que contiene a x .

Es obvio que si $x, y \in X$, entonces $C(x)$ y $C(y)$ son iguales o disjuntas. En efecto, si tienen puntos en común, por el teorema 3.21 resulta que $C(x) \cup C(y)$ es un conexo, luego $C(x) \cup C(y) \subset C(x)$ y $C(x) \cup C(y) \subset C(y)$, con lo que $C(x) = C(x) \cup C(y) = C(y)$.

En resumen, todo espacio X está dividido en componentes conexas disjuntas. Las componentes conexas son cerradas por el teorema 3.22. En efecto, $\overline{C(x)}$ es un conexo que contiene a x , luego $\overline{C(x)} \subset C(x)$.

Sin embargo las componentes conexas no siempre son abiertas. Si un espacio tiene un número finito de componentes conexas, éstas serán abiertas y cerradas a la vez, evidentemente, pero si hay infinitas componentes ya no es necesario. Por ejemplo, ningún subconjunto de \mathbb{Q} con más de un punto es conexo, porque no es un intervalo de \mathbb{R} , luego las componentes conexas de \mathbb{Q} son los puntos, que no son abiertos.

Espacios arcoconexos Para probar que un espacio es conexo, resulta útil el concepto de arco.

Un *arco* en un espacio X es una aplicación continua $a : [0, 1] \rightarrow X$. El espacio X es *arco-conexo* si para todo par de puntos $x, y \in X$ existe un arco $a : [0, 1] \rightarrow X$ tal que $a(0) = x$, $a(1) = y$.



Todo espacio arcoconexo es conexo, pues si x e y son dos puntos cualesquiera y a es un arco que los une, entonces x e y están en la imagen del arco a , que es un conexo, luego x e y están en la misma componente conexa de X , luego sólo hay una. El recíproco no es cierto, pero no vamos a dar un ejemplo.

Espacios localmente convexos Dados dos puntos $x, y \in \mathbb{R}^n$, el segmento que los une está formado por los puntos de la forma $y + \lambda(x - y)$, con $\lambda \in [0, 1]$. Esto se puede definir en cualquier \mathbb{R} -espacio vectorial. Si V es un espacio vectorial topológico y $x, y \in V$, entonces la aplicación $a : [0, 1] \rightarrow V$ dada por $a(\lambda) = (1 - \lambda)x + \lambda y$ es un arco, el *segmento* que une x con y .

Un subconjunto A de un \mathbb{R} -espacio vectorial V es *convexo* si para todos los puntos $x, y \in A$ y todo $\lambda \in [0, 1]$ se cumple $(1 - \lambda)x + \lambda y \in A$, es decir, si cuando A contiene a dos puntos, también contiene al segmento que los une.

Uniendo todo esto, resulta que en un espacio vectorial topológico sobre \mathbb{R} todo convexo es arco-conexo, luego conexo. En particular todo espacio vectorial topológico sobre \mathbb{R} es conexo. En particular \mathbb{R}^n es conexo.

Ejercicio: Probar que toda esfera de centro O en \mathbb{R}^n es imagen continua de $\mathbb{R}^n \setminus \{0\}$. Probar que $\mathbb{R}^n \setminus \{0\}$ es conexo y deducir de aquí la conexión de la esfera.

Ejercicio: Probar que \mathbb{R}^2 no es homeomorfo a \mathbb{R} .

Los conjuntos convexos tienen una propiedad que en general no cumplen los conexos, y es que, claramente, la intersección de convexos es convexa.

Ejemplo Las bolas en los espacios normados (sobre \mathbb{R}) son convexas.

En efecto, si $x, y \in B_\epsilon(z)$, entonces $\|x - z\| < \epsilon$, $\|y - z\| < \epsilon$, luego para todo $\lambda \in [0, 1]$ se cumple

$$\begin{aligned} \|(1 - \lambda)x + \lambda y - z\| &= \|(1 - \lambda)x + \lambda y - (1 - \lambda)z - \lambda z\| \leq \\ &(1 - \lambda)\|x - z\| + \lambda\|y - z\| < (1 - \lambda)\epsilon + \lambda\epsilon = \epsilon. \end{aligned}$$

Por lo tanto $(1 - \lambda)x + \lambda y \in B_\epsilon(z)$.

(Cambiando las desigualdades estrictas por desigualdades no estrictas se prueba que las bolas cerradas son convexas.) ■

Espacios localmente convexos, conexos o arco-conexos Un espacio vectorial topológico (sobre \mathbb{R}) es *localmente convexo* si tiene una base formada por conjuntos convexos. El ejemplo anterior prueba que todos los espacios normados son localmente convexos.

Igualmente, un espacio topológico es *localmente conexo* o *localmente arco-conexo* si tiene una base formada por abiertos conexos o arco-conexos, respectivamente.

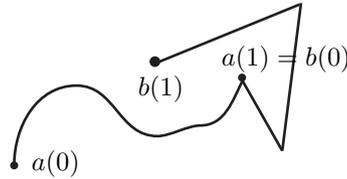
Así, los espacios localmente convexos son localmente arco-conexos y los espacios localmente arco-conexos son localmente conexos.

Un hecho importante es que si A es un abierto en un espacio localmente conexo X , entonces las componentes conexas de A son abiertas y cerradas.

En efecto, si C es una componente conexa de A y $x \in C$, entonces existe un abierto (básico) U de X tal que U es conexo y $x \in U \subset A$. Como C es la componente conexa de x , ha de ser $U \subset C$, luego C es un entorno de x , o sea, C es entorno de todos sus puntos, luego es un abierto.

Ejercicio: Probar que todo abierto no vacío en \mathbb{R} es la unión disjunta de una cantidad numerable de intervalos abiertos.

Sean a y b dos arcos en un espacio X de modo que $a(1) = b(0)$. Entonces la aplicación $b_1 : [1, 2] \rightarrow X$ dada por $b_1(t) = b(t-1)$ es continua y cumple que $b_1(1) = b(0)$, $b_1(2) = b(1)$ (se trata de la composición con b del homeomorfismo $i : [1, 2] \rightarrow [0, 1]$ dado por $i(t) = t-1$).



Ahora, la unión $c = a \cup b_1 : [0, 2] \rightarrow X$, esto es, la aplicación que restringida a $[0, 1]$ es a y restringida a $[1, 2]$ es b_1 , es continua (porque restringida a los dos cerrados $[0, 1]$ y $[1, 2]$ lo es), y su imagen es la unión de las imágenes de a y b .

Finalmente, llamamos $a \cup b : [0, 1] \rightarrow X$ a la función $(a \cup b)(t) = c(2t)$, es decir, la composición de c con el homeomorfismo $j : [0, 1] \rightarrow [0, 2]$ definido mediante $j(t) = 2t$. Claramente $a \cup b$ es un arco cuya imagen es la unión de las imágenes de a y b . En particular $(a \cup b)(0) = a(0)$ y $(a \cup b)(1) = b(1)$.

Esto significa que si podemos unir un punto x con un punto y a través de un arco a , y podemos unir un punto y con un punto z a través de un arco b , entonces podemos unir x con z mediante el arco $a \cup b$.

Por otra parte, si a es un arco en un espacio X , la aplicación $-a : [0, 1] \rightarrow X$ dada por $(-a)(t) = a(1-t)$ es un arco con la misma imagen pero de modo que $(-a)(0) = a(1)$ y $(-a)(1) = a(0)$. También es obvio que un arco constante une un punto consigo mismo.

De todo esto se sigue que la relación “ x e y se pueden unir mediante un arco” es una relación de equivalencia en todo espacio X . Las clases de equivalencia reciben el nombre de *componentes arco-conexas* del espacio X y son una partición de X en subespacios arco-conexos.

Las componentes arco-conexas no tienen por qué ser abiertas ni cerradas, pero si X es localmente arco-conexo, es claro que las componentes arco-conexas son abiertas, luego también son cerradas, puesto que el complementario de una

es la unión de las demás, y esto implica que las componentes arcoconexas de X coinciden con sus componentes conexas. A su vez, esto implica:

Teorema 3.24 *Sea X un espacio localmente arco-conexo. Entonces un abierto de X es conexo si y sólo si es arco-conexo.*

DEMOSTRACIÓN: Obviamente los abiertos arco-conexos son conexos. Si A es un abierto conexo no vacío, entonces A es también localmente arco-conexo, luego tiene una única componente arco-conexa abierta y cerrada, luego, al ser A conexo, tiene que ser todo A , luego A es arco-conexo. ■

Una *poligonal* es una unión de un número finito de segmentos. Una pequeña modificación del teorema anterior permite probar que en un espacio localmente conexo, un abierto es conexo si y sólo si es conexo por poligonales, es decir, si todo par de puntos se puede unir por una poligonal.

Aplicaciones de la conexión Una de las consecuencias más relevantes de la conexión de un espacio topológico es el siguiente hecho obvio:

Teorema 3.25 (Teorema de los valores intermedios) *Si X es un espacio conexo, $f : X \rightarrow \mathbb{R}$ es una aplicación continua, x, y son puntos de X y $f(x) < \alpha < f(y)$, entonces existe un punto $z \in X$ tal que $f(z) = \alpha$.*

DEMOSTRACIÓN: Se cumple que $f[X]$ es un conexo, luego un intervalo. Como $f(x)$ y $f(y)$ están en $f[X]$, también $\alpha \in f[X]$. ■

A pesar de su simplicidad, las consecuencias de este teorema son importantes. Por ejemplo, no hay polinomios irreducibles de grado impar sobre \mathbb{R} , salvo los de grado 1:

Teorema 3.26 *Todo polinomio $p(x) \in \mathbb{R}[x]$ de grado impar tiene al menos una raíz en \mathbb{R} .*

DEMOSTRACIÓN: Como $p(x)$ tiene una raíz si y sólo si la tiene $-p(x)$, podemos suponer que su coeficiente director es positivo.

Entonces $\lim_{x \rightarrow +\infty} p(x) = +\infty$, mientras que $\lim_{x \rightarrow -\infty} p(x) = -\infty$. En particular existe un $u \in \mathbb{R}$ tal que $p(u) < 0$ y existe un $v \in \mathbb{R}$ tal que $p(v) > 0$. Por el teorema de los valores intermedios también existe un $a \in \mathbb{R}$ tal que $p(a) = 0$. ■

Por supuesto el teorema es falso para polinomios de grado par. Basta pensar en $x^2 + 1$.

Si $a > 0$, el teorema de los valores intermedios aplicado al polinomio $x^n - a$ nos permite concluir la existencia de un $b > 0$ tal que $b^n = a$. Es claramente único, pues si $b^n = c^n$, entonces $(b/c)^n = 1$, de donde $b/c = \pm 1$, luego si ambos son positivos $b = c$.

Definición 3.27 Para cada natural $n > 0$ y cada número real $a > 0$ definimos la raíz n -ésima de a como el único número $b > 0$ tal que $b^n = a$. Lo representaremos $b = \sqrt[n]{a}$

Unas comprobaciones rutinarias muestran que si m, n son números enteros $n > 0$ y $a > 0$ entonces el número $a^{m/n} = (\sqrt[n]{a})^m$ depende sólo de la fracción m/n , con lo que tenemos definida la exponencial a^r para todo número real positivo a y todo número racional r y extiende a la exponencial entera. También se comprueba que $a^{r+s} = a^r a^s$, $(a^r)^s = a^{rs}$.

3.4 Espacios completos

Retomamos ahora el concepto de espacio métrico completo, que introdujimos en el capítulo I. Veamos algunas propiedades elementales de la completitud:

Teorema 3.28 *Se cumple:*

- a) *Todo espacio métrico compacto es completo.*
- b) *Todo cerrado en un espacio métrico completo es completo.*
- c) *Todo subespacio completo de un espacio métrico es cerrado.*

DEMOSTRACIÓN: a) Toda sucesión tiene una subsucesión convergente, luego si es de Cauchy es convergente.

b) Si M es un espacio métrico completo y C es un cerrado en M , entonces toda sucesión de Cauchy en C converge en M , y como C es cerrado su límite estará en C , luego la sucesión converge en C .

c) Si C es un subespacio completo de un espacio métrico M , dado un punto x en la clausura de C , existe una sucesión en C que converge a x , luego la sucesión es de Cauchy, luego converge en C , luego x está en C , luego C es cerrado. ■

Por supuesto no todo espacio completo es compacto. A veces es útil conocer lo que separa a un espacio completo de la compacidad:

Definición 3.29 Un espacio métrico M es *precompacto* si para cada $\epsilon > 0$ existen puntos x_1, \dots, x_n en M tales que $M = B_\epsilon(x_1) \cup \dots \cup B_\epsilon(x_n)$.

Obviamente todo espacio compacto es precompacto (basta extraer un subcubrimiento finito del cubrimiento formado por todas las bolas de radio ϵ).

El teorema 3.6 contiene la demostración de que un espacio precompacto y completo es compacto. En efecto, partiendo de la hipótesis sobre subsucesiones convergentes, en primer lugar se prueba que el espacio es precompacto. Con ayuda de la precompacidad se construye una sucesión de bolas $B_{1/(n+1)}(x_n)$ en las que podemos exigir que cada una de ellas corte a la anterior. Esto garantiza que la sucesión de los centros es de Cauchy, con lo que podemos garantizar su convergencia por la completitud y probamos la compacidad del espacio. Así pues:

Teorema 3.30 *Un espacio métrico es compacto si y sólo si es precompacto y completo.*

Estudiamos ahora los espacios normados completos:

Definición 3.31 Un *espacio de Banach* es un espacio normado completo.

Ejercicio: Probar que si $E \neq 0$ es un espacio de Banach sobre un cuerpo métrico \mathbb{K} , entonces \mathbb{K} es completo.

Vamos a probar que \mathbb{R}^n es un espacio de Banach con cualquiera de las normas consideradas en el teorema 2.2. De hecho, vamos a probar un resultado más general. Para ello conviene introducir la definición siguiente:

Sea E un espacio vectorial sobre un cuerpo métrico \mathbb{K} . Se dice que dos normas $\|\cdot\|_1$ y $\|\cdot\|_2$ en V son *equivalentes* si existen números reales $M_1, M_2 > 0$ tales que

$$\|x\|_1 \leq M_2 \|x\|_2 \quad \text{y} \quad \|x\|_2 \leq M_1 \|x\|_1$$

para todo $x \in E$.

Dos normas equivalentes inducen la misma topología en E , pues claramente $B_{\epsilon/M_1}^1(a) \subset B_\epsilon^2(a)$ y $B_{\epsilon/M_2}^2(a) \subset B_\epsilon^1(a)$, donde el superíndice en las bolas indica la norma respecto a la que se calculan. Estas inclusiones hacen que todo abierto para una norma lo sea también para la otra.

También es claro que si dos normas son equivalentes entonces E es completo con una si y sólo si lo es con la otra, así como que ambas determinan los mismos subconjuntos acotados.

Teorema 3.32 *Sea E un espacio vectorial de dimensión finita sobre un cuerpo métrico completo \mathbb{K} . Entonces todas las normas en E son equivalentes, luego todas determinan los mismos conjuntos acotados y además E es completo con cualquiera de ellas.*

DEMOSTRACIÓN: Supongamos primero que $E = \mathbb{K}^n$ y consideremos la norma $\|\cdot\|_\infty$ definida en 2.2. Si $\{x_m\}_{m=0}^\infty$ es una sucesión de Cauchy en E , es inmediato comprobar que las sucesiones $\{x_{mi}\}_{m=0}^\infty$ son de Cauchy en \mathbb{K} , para $i = 1, \dots, n$, luego convergen a ciertos valores $x_i \in \mathbb{K}$, que determinan un $x \in \mathbb{K}^n$, y es claro que la sucesión original converge a x (por ejemplo, porque sabemos que la topología de \mathbb{K}^n es la topología producto).

Basta probar que cualquier otra norma $\|\cdot\|$ en \mathbb{K}^n es equivalente a $\|\cdot\|_\infty$.

Sea e_1, \dots, e_n la base canónica de \mathbb{K}^n . Entonces para todo $x \in \mathbb{K}^n$ se cumple

$$\|x\| = \|x_1 e_1 + \dots + x_n e_n\| \leq |x_1| \|e_1\| + \dots + |x_n| \|e_n\| \leq M \|x\|_\infty,$$

donde $M = \|e_1\| + \dots + \|e_n\|$.

Ahora hemos de probar la relación opuesta. Basta ver que existen números reales $N_i > 0$ de modo que si $x \in \mathbb{K}^n$, entonces $|x_i| \leq N_i \|x\|$, pues en tal caso $N = \max N_i$ cumple $\|x\|_\infty \leq N \|x\|$. Lo probaremos por inducción sobre n .

Si $n = 1$ basta tomar $N_1 = 1/\|1\|$. Supuesto cierto para $n - 1$, identificamos \mathbb{K}^{n-1} con los elementos de \mathbb{K}^n cuya última coordenada es nula. Las restricciones a \mathbb{K}^{n-1} de las dos normas consideradas son normas en \mathbb{K}^{n-1} . Por hipótesis de inducción son equivalentes y \mathbb{K}^{n-1} es completo para la restricción de la norma $\|\cdot\|$, luego es cerrado en \mathbb{K}^n para la topología inducida por esta norma.

Supongamos, por reducción al absurdo, que para todo natural m existe un $w_m \in \mathbb{K}^n$ de manera que $|w_{mn}| > m\|w_m\|$. Podemos suponer que $w_{mn} = 1$ (basta dividir w_m entre w_{mn} si es preciso), y entonces la desigualdad se reduce a $\|w_m\| < 1/m$. Por otra parte esta condición adicional implica también que $w_m - e_n \in \mathbb{K}^{n-1}$.

De este modo tenemos que $\{w_m\}_{m=0}^\infty$ tiende a 0 y que $w_m - e_n$ tiende a $-e_n$ pero, como \mathbb{K}^{n-1} es cerrado, esto implica que $e_n \in \mathbb{K}^{n-1}$, lo cual es absurdo. Por lo tanto existe un m tal que $|w_n| \leq m\|w\|$ para todo $w \in \mathbb{K}^n$. El mismo razonamiento se aplica a cualquier otro índice.

Esto termina la prueba para $E = \mathbb{K}^n$. Si E es ahora un espacio vectorial cualquiera de dimensión n sobre \mathbb{K} , cada norma en E induce una en \mathbb{K}^n a través de un isomorfismo de espacios vectoriales. Del hecho de que las normas inducidas sean equivalentes se sigue obviamente que las normas de partida también lo sean. Igualmente se concluye que E es completo con cualquiera de ellas. ■

En particular $\mathbb{C} = \mathbb{R}^2$ es un cuerpo métrico completo. Luego todos los resultados que probamos para cuerpos métricos completos \mathbb{K} se aplican tanto a $\mathbb{K} = \mathbb{R}$ como a $\mathbb{K} = \mathbb{C}$.

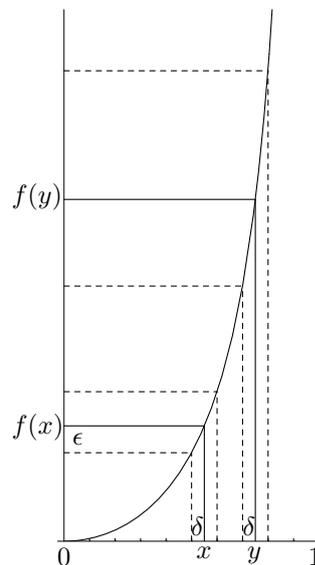
Al contrario que la compacidad o la conexión, la completitud no es una propiedad topológica, sino métrica, pues, por ejemplo, sabemos que \mathbb{R} es un espacio métrico completo, mientras que $] -1, 1[$ no lo es.

Si analizamos lo que falla, vemos que en $] -1, 1[$ hay sucesiones de Cauchy no convergentes, por ejemplo las que convergen a 1 en \mathbb{R} , pero cuando las transformamos por un homeomorfismo entre $] -1, 1[$ y \mathbb{R} se convierten en sucesiones que tienden a $+\infty$, que ya no son de Cauchy, luego no violan la completitud de \mathbb{R} . El problema es que mientras una aplicación continua transforma sucesiones convergentes en sucesiones convergentes, no transforma necesariamente sucesiones de Cauchy en sucesiones de Cauchy. Y a su vez esto se debe a que si se estira infinitamente una sucesión de Cauchy, ésta deja de serlo. Esto nos lleva a conjeturar que la completitud se conservará por aplicaciones continuas que no produzcan estiramientos infinitos. Vamos a definir este tipo de aplicaciones.

Definición 3.33 Una aplicación $f : M \rightarrow N$ entre espacios métricos es *uniformemente continua* si para todo $\epsilon > 0$ existe un $\delta > 0$ de modo que si $x, y \in M$ cumplen $d(x, y) < \delta$, entonces $d(f(x), f(y)) < \epsilon$.

Desde un punto de vista lógico, la diferencia entre “aplicación continua” y “uniformemente continua” es sutil. Conviene confrontarlas. Recordemos que f es continua si para todo punto x y todo $\epsilon > 0$ existe un $\delta > 0$ tal que si $y \in M$ cumple $d(x, y) < \delta$, entonces $d(f(x), f(y)) < \epsilon$.

La diferencia es, pues, que cuando f es uniformemente continua el mismo δ verifica la definición de continuidad para un ϵ dado simultáneamente en todos los puntos x , mientras que la definición de continuidad permite tomar un δ distinto en función del punto x . Por ejemplo, el homeomorfismo f entre $] -1, 1[$ y \mathbb{R} estira más los puntos cuanto más próximos están de ± 1 . Si tomamos un punto x y queremos garantizar que $f(z)$ diste de $f(x)$ menos que ϵ , tendremos que exigir que z diste de x menos que un cierto δ , pero si consideramos los puntos que distan menos que δ de un punto y más cercano a 1, vemos que muchos de ellos se transforman en puntos que distan de $f(y)$ mucho más que ϵ , y que si queremos que no sobrepasen esta cota hemos de tomar un δ mucho menor.



La imagen de una sucesión de Cauchy por una aplicación uniformemente continua es una sucesión de Cauchy. En efecto, si $\{a_n\}_{n=0}^{\infty}$ es de Cauchy y f es uniformemente continua, dado $\epsilon > 0$ existe un $\delta > 0$ de manera que si $d(x, y) < \delta$, entonces $d(f(x), f(y)) < \epsilon$. Existe un natural n_0 tal que si $n, m \geq n_0$ entonces $d(a_m, a_n) < \delta$, luego $d(f(a_m), f(a_n)) < \epsilon$. Esto prueba que la sucesión $\{f(a_n)\}_{n=0}^{\infty}$ es de Cauchy.

Como consecuencia, si dos espacios métricos X e Y son *uniformemente homeomorfos*, esto es, si existe una biyección uniformemente continua con inversa uniformemente continua entre ellos, uno es completo si y sólo si lo es el otro.

No obstante, una imagen uniformemente continua de un espacio completo no tiene por qué ser completa.

Teorema 3.34 Sea $f : K \rightarrow M$ una aplicación continua entre espacios métricos y supongamos que K es compacto. Entonces f es uniformemente continua.

DEMOSTRACIÓN: Sea $\epsilon > 0$. Para cada punto $x \in K$, como f es continua en x , existe un $\delta(x) > 0$ tal que si $d(y, x) < \delta(x)$, entonces $d(f(y), f(x)) < \epsilon/2$. Cubramos el espacio K por todas las bolas abiertas de centro cada punto x y de radio $\delta(x)/2$ y tomemos un subcubrimiento finito. Digamos que las bolas que forman este subcubrimiento tienen centros en los puntos x_1, \dots, x_n . Sea $\delta > 0$ el mínimo del conjunto $\{\delta(x_1)/2, \dots, \delta(x_n)/2\}$.

Si x, y son dos puntos cualesquiera de K tales que $d(x, y) < \delta$, entonces x estará en una de las bolas $B_{\delta(x_i)/2}(x_i)$. Entonces

$$d(y, x_i) \leq d(y, x) + d(x, x_i) < \frac{\delta(x_i)}{2} + \frac{\delta(x_i)}{2} = \delta(x_i).$$

Por lo tanto $x, y \in B_{\delta(x_i)}(x_i)$, luego $d(f(x), f(x_i)) < \epsilon/2$, $d(f(y), f(x_i)) < \epsilon/2$. Consecuentemente, $d(f(x), f(y)) < \epsilon$. ■

Otro caso en el que la continuidad equivale a la continuidad uniforme es el de las aplicaciones lineales entre espacios normados:

Teorema 3.35 *Sea $f : E \rightarrow F$ una aplicación lineal entre espacios normados. Las siguientes condiciones son equivalentes:*

- a) f es continua en E .
- b) f es continua en 0 .
- c) f está acotada en $\overline{B}_1(0)$.
- d) Existe un $M \geq 0$ tal que para todo $x \in E$ se cumple $\|f(x)\| \leq M\|x\|$.
- e) f es uniformemente continua en E .

DEMOSTRACIÓN: a) \Rightarrow b) es obvio.

b) \Rightarrow c), pues existe un $\delta > 0$ tal que $\|x-0\| \leq \delta$, entonces $\|f(x)-f(0)\| \leq 1$.

Por lo tanto si $\|x\| \leq 1$, se cumple $\|\delta x\| \leq \delta$, $\|f(\delta x)\| \leq 1$, $\|f(x)\| \leq 1/\delta$, o sea, que $1/\delta$ es una cota de f en $\overline{B}_1(0)$.

c) \Rightarrow d), pues si M es una cota de f en $\overline{B}_1(0)$, dado cualquier $x \neq 0$ se cumple que $x/\|x\| \in \overline{B}_1(0)$, luego $\|f(x/\|x\|)\| \leq M$, de donde $\|f(x)\| \leq M\|x\|$, y esto también es cierto si $x = 0$.

d) \Rightarrow e), pues para todos los x, y en E :

$$\|f(x) - f(y)\| = \|f(x - y)\| \leq M\|x - y\|,$$

luego dado $\epsilon > 0$, si $\|x - y\| < \epsilon/M$, se cumple $\|f(x) - f(y)\| < \epsilon$.

e) \Rightarrow a) es evidente. ■

En particular, si \mathbb{K} es un cuerpo métrico completo, todo isomorfismo entre dos \mathbb{K} -espacios vectoriales de dimensión finita es un homeomorfismo uniforme para cualquier par de normas.

3.5 Espacios de Hilbert

En esta sección \mathbb{K} representará al cuerpo \mathbb{R} de los números reales o bien al cuerpo \mathbb{C} de los números complejos. Si $\alpha \in \mathbb{K}$, la notación $\bar{\alpha}$ representará al conjugado de α si $\mathbb{K} = \mathbb{C}$ o simplemente $\bar{\alpha} = \alpha$ si $\mathbb{K} = \mathbb{R}$.

Definición 3.36 Si H es un \mathbb{K} -espacio vectorial, un *producto escalar* en H es una aplicación $\cdot : H \times H \rightarrow \mathbb{K}$ que cumple las propiedades siguientes (donde $x, y, z \in H$ y $\alpha \in \mathbb{K}$):

- a) $x \cdot y = \overline{y \cdot x}$,
- b) $(x + y) \cdot z = x \cdot z + y \cdot z$,
- c) $(\alpha x) \cdot y = \alpha(x \cdot y)$,
- d) $x \cdot x \geq 0$ y $x \cdot x = 0$ si y sólo si $x = 0$.

Notemos que a) y b) implican también la propiedad distributiva por la derecha: $x \cdot (y + z) = x \cdot y + x \cdot z$. De las propiedades a) y c) se sigue que

$$x \cdot (\alpha y) = \overline{(\alpha y)} \cdot x = \overline{\alpha} \overline{y} \cdot x = \overline{\alpha} x \cdot y.$$

Un *espacio prehilbertiano* es un par (H, \cdot) , donde H es un \mathbb{K} -espacio vectorial y \cdot es un producto escalar en H . En la práctica escribiremos simplemente H en lugar de (H, \cdot) .

Si H es un espacio prehilbertiano, definimos su *norma* asociada como la aplicación $\| \cdot \| : H \rightarrow \mathbb{R}$ dada por $\|x\| = \sqrt{x \cdot x}$.

Vamos a probar que ciertamente se trata de una norma en H , pero antes conviene probar este hecho:

Teorema 3.37 (Desigualdad de Schwarz) *Si x, y son elementos de un espacio prehilbertiano, entonces $|x \cdot y| \leq \|x\| \|y\|$.*

DEMOSTRACIÓN: Sean $A = \|x\|^2$, $B = |x \cdot y|$ y $C = \|y\|^2$. Existe un número complejo α tal que $|\alpha| = 1$ y $\alpha(y \cdot x) = B$. Para todo número real r se cumple

$$0 \leq (x - r\alpha y) \cdot (x - r\alpha y) = x \cdot x - r\alpha(y \cdot x) - r\bar{\alpha}(x \cdot y) + r^2 y \cdot y.$$

Notemos que $\bar{\alpha}(x \cdot y) = \bar{B} = B$, luego $A - 2Br + Cr^2 \geq 0$. Si $C = 0$ ha de ser $B = 0$, o de lo contrario la desigualdad sería falsa para r grande. Si $C > 0$ tomamos $r = B/C$ y obtenemos $B^2 \leq AC$, y basta tomar la raíz cuadrada de ambos miembros. ■

Ahora ya es inmediato que la norma determinada por un producto escalar es realmente una norma: De la propiedad d) se sigue que $\|x\| = 0$ si y sólo si $x = 0$, y por otra parte

$$\|\alpha x\| = \sqrt{(\alpha x) \cdot (\alpha x)} = \sqrt{\alpha \bar{\alpha} x \cdot x} = \sqrt{|\alpha|^2} \sqrt{x \cdot x} = |\alpha| \|x\|.$$

Por último, usando la desigualdad de Schwarz vemos que

$$\begin{aligned} \|x + y\|^2 &= (x + y) \cdot (x + y) = x \cdot x + x \cdot y + y \cdot x + y \cdot y \\ &\leq \|x\|^2 + 2\|x\| \|y\| + \|y\|^2 = (\|x\| + \|y\|)^2, \end{aligned}$$

donde hemos usado que $x \cdot y + y \cdot x = x \cdot y + \overline{x \cdot y}$ es un número real, luego

$$x \cdot y + y \cdot x \leq |x \cdot y + y \cdot x| \leq |x \cdot y| + |y \cdot x| \leq 2\|x\| \|y\|.$$

Calculando la raíz cuadrada llegamos a que $\|x + y\| \leq \|x\| + \|y\|$.

Así pues, todo espacio prehilbertiano es un espacio normado con la norma derivada de su producto escalar. Un *espacio de Hilbert* es un espacio prehilbertiano H que sea completo como espacio normado, es decir, que sea un espacio de Banach.

Ejemplo Un producto escalar en el espacio \mathbb{K}^n viene dado por

$$x \cdot y = x_1 \bar{y}_1 + \cdots + x_n \bar{y}_n.$$

De este modo, $\|x\| = \sqrt{|x_1|^2 + \cdots + |x_n|^2}$. El teorema 3.32 implica que \mathbb{K}^n es un espacio de Hilbert con este producto escalar. ■

Teorema 3.38 Si H es un espacio prehilbertiano, entonces el producto escalar en H es una función continua.

DEMOSTRACIÓN: Para todos los $x, x', y, y' \in H$ se cumple

$$\begin{aligned} |x \cdot y - x' \cdot y'| &\leq |(x - x') \cdot y| + |x' \cdot (y - y')| \leq \|x - x'\| \|y\| + \|x'\| \|y - y'\| \\ &\leq \|x - x'\| \|y\| + \|x' - x\| \|y - y'\| + \|x\| \|y - y'\|. \end{aligned}$$

Así, dado un par $(x, y) \in H \times H$ y un $\epsilon > 0$, todo par $(x', y') \in H \times H$ que cumpla $\|x' - x\|, \|y' - y\| < \epsilon/3M$, donde $M > \|x\|, \|y\|$, cumple también que

$$|x \cdot y - x' \cdot y'| < \frac{\epsilon}{3M} M + \frac{\epsilon^2}{9M^2} + \frac{\epsilon}{3M} M < \epsilon.$$

■

Definición 3.39 Si H es un espacio prehilbertiano, diremos que $x, y \in H$ son *ortogonales*, y lo representaremos por $x \perp y$, si $x \cdot y = 0$.

Sabemos que la ortogonalidad en \mathbb{R}^n coincide con el concepto geométrico de perpendicularidad. Es fácil generalizar el teorema de Pitágoras:

$$\text{Si } x \perp y, \text{ entonces } \|x + y\|^2 = \|x\|^2 + \|y\|^2.$$

Así mismo, las propiedades del producto escalar dan inmediatamente la fórmula conocida como *identidad del paralelogramo*:

$$\|x + y\|^2 + \|x - y\|^2 = 2\|x\|^2 + 2\|y\|^2.$$

Para cada $A \subset H$, definimos

$$A^\perp = \{x \in H \mid x \perp a \text{ para todo } a \in A\}.$$

Es claro que A^\perp es un subespacio vectorial de H . Más aún, puesto que $\{a\}^\perp$ es la antiimagen de 0 por la aplicación continua $x \mapsto a \cdot x$, se cumple que $\{a\}^\perp$ es cerrado, y como

$$A^\perp = \bigcap_{a \in A} \{a\}^\perp,$$

vemos que A^\perp es un subespacio cerrado de H .

Si V es un subespacio vectorial de H es claro que $V \cap V^\perp = 0$. Vamos a probar que si V es cerrado entonces $H = V \oplus V^\perp$. Para ello necesitamos un resultado previo:

Teorema 3.40 *Sea M un subconjunto no vacío, cerrado y convexo de un espacio de Hilbert H . Entonces M contiene un único elemento de norma mínima.*

DEMOSTRACIÓN: Sea δ el ínfimo de las normas de los elementos de M . Aplicando la identidad del paralelogramo a $\frac{1}{2}x, \frac{1}{2}y$ tenemos

$$\frac{1}{4}\|x - y\|^2 = \frac{1}{2}\|x\|^2 + \frac{1}{2}\|y\|^2 - \left\| \frac{x + y}{2} \right\|^2.$$

Si $x, y \in M$, por convexidad $(x + y)/2 \in M$, luego

$$\|x - y\|^2 \leq 2\|x\|^2 + 2\|y\|^2 - 4\delta^2. \quad (3.1)$$

Si $\|x\| = \|y\| = \delta$ esto implica $x = y$, lo que nos da la unicidad. Es fácil construir una sucesión $\{x_n\}_{n=0}^\infty \subset M$ tal que $\lim_n \|x_n\| = \delta$. Aplicando 3.1 a x_m y x_n concluimos fácilmente que la sucesión es de Cauchy, luego converge a un punto x que, por continuidad de la norma, cumplirá $\|x\| = \delta$. Además, como M es cerrado, ha de ser $x \in M$ y claramente su norma es la mínima en M . ■

Teorema 3.41 *Sea V un subespacio cerrado de un espacio de Hilbert H . Entonces*

- a) *Todo $x \in H$ se descompone de forma única como $x = Px + Qx$, donde $Px \in V, Qx \in V^\perp$.*
- b) *Px y Qx son los puntos de V y V^\perp más próximos a x .*
- c) *Las aplicaciones $P : H \rightarrow V$ y $Q : H \rightarrow V^\perp$ son lineales y continuas.*
- d) $\|x\|^2 = \|Px\|^2 + \|Qx\|^2$.

DEMOSTRACIÓN: El conjunto $x + V$ es cerrado y convexo, luego podemos definir Qx como el elemento de norma mínima en $x + V$. Definimos $Px = x - Qx$. Obviamente $Px \in V$. Veamos que $Qx \in V^\perp$. Para ello probaremos que $(Qx) \cdot y = 0$ para todo $y \in V$. No perdemos generalidad si suponemos $\|y\| = 1$. Por definición de Qx tenemos que

$$(Qx) \cdot (Qx) = \|Qx\|^2 \leq \|Qx - \alpha y\|^2 = (Qx - \alpha y) \cdot (Qx - \alpha y),$$

para todo $\alpha \in \mathbb{K}$. Simplificando queda

$$0 \leq -\alpha(y \cdot Qx) - \bar{\alpha}(Qx \cdot y) + \alpha\bar{\alpha},$$

y si hacemos $\alpha = (Qx) \cdot y$ queda $0 \leq -|(Qx) \cdot y|^2$, luego $(Qx) \cdot y = 0$.

Esto prueba la existencia de la descomposición de a). La unicidad se debe a que $V \cap V^\perp = 0$. En definitiva, $H = V \oplus V^\perp$.

Ahora observamos que si $y \in V$ entonces

$$\|x - y\|^2 = \|Qx + (Px - y)\|^2 = \|Qx\|^2 + \|Px - y\|^2,$$

luego la mínima distancia entre x y un punto $y \in V$ se alcanza cuando $y = Px$. Esto prueba b). El apartado d) es el teorema de Pitágoras. La linealidad de P y Q es obvia. La continuidad se debe a que por d) tenemos $\|Px\| \leq \|x\|$, $\|Qx\| \leq \|x\|$, y basta aplicar el teorema 3.35. ■

El teorema siguiente caracteriza las aplicaciones lineales continuas de un espacio de Hilbert en \mathbb{K} .

Teorema 3.42 *Sea H un espacio de Hilbert y $f : H \rightarrow \mathbb{K}$ una aplicación lineal continua. Entonces existe un único $y \in H$ tal que $f(x) = x \cdot y$ para todo $x \in H$.*

DEMOSTRACIÓN: Si f es la aplicación nula tomamos $y = 0$. En otro caso sea V el núcleo de f , que será un subespacio cerrado propio de H . Por el teorema anterior $V^\perp \neq 0$, luego podemos tomar $z \in V^\perp$ con $\|z\| = 1$. Sea $u = f(x)z - f(z)x$. Como $f(u) = f(x)f(z) - f(z)f(x) = 0$, tenemos que $u \in V$, luego $u \cdot z = 0$. La definición de u implica

$$f(x) = f(x)(z \cdot z) = f(z)(x \cdot z).$$

Tomando $y = \overline{f(z)}z$ resulta $f(x) = x \cdot y$.

La unicidad es obvia, pues si dos puntos y, y' cumplen el teorema, entonces $y - y'$ es ortogonal a todo $x \in H$. ■

Hasta ahora, los únicos espacios de Hilbert que conocemos son los espacios de dimensión finita \mathbb{K}^n . Es fácil probar que todo espacio prehilbertiano H de dimensión finita es isométrico a \mathbb{K}^n , es decir, existe un isomorfismo $H \rightarrow \mathbb{K}^n$ que conserva el producto escalar (porque el teorema de ortonormalización de Gram-Schmidt² asegura la existencia de una base ortonormal, y un isomorfismo que transforme una base ortonormal de H en una de \mathbb{K}^n es una isometría). Veamos ahora un ejemplo de espacio de Hilbert de dimensión infinita:

Definición 3.43 Llamaremos ℓ^2 al conjunto de todas las sucesiones $\{x_n\}_{n=0}^\infty$ en \mathbb{K} tales que $\sum_{n=0}^\infty |x_n|^2$ es convergente.

Teorema 3.44 ℓ^2 es un espacio de Hilbert con el producto escalar dado por

$$x \cdot y = \sum_{n=0}^\infty x_n \bar{y}_n.$$

DEMOSTRACIÓN: En primer lugar demostramos que la serie que define el producto escalar es convergente. Para ello basta probar que es absolutamente convergente. Ahora bien, sean $A = \sum_{n=0}^\infty |x_n|^2$ y $B = \sum_{n=0}^\infty |y_n|^2$ y supongamos que $A \neq 0 \neq B$. Entonces

$$\left(\frac{|x_n|}{A} - \frac{|y_n|}{B} \right)^2 = \frac{|x_n|^2}{A^2} + \frac{|y_n|^2}{B^2} - 2 \frac{|x_n||y_n|}{AB} \geq 0,$$

²En [G 4.19] está probado para espacios sobre $\mathbb{K} = \mathbb{R}$, pero el mismo argumento vale si $\mathbb{K} = \mathbb{C}$.

luego

$$\frac{|x_n||y_n|}{AB} \leq \frac{1}{2A}|x_n|^2 + \frac{1}{2B}|y_n|^2,$$

luego

$$\sum_{n=0}^{\infty} \frac{|x_n||y_n|}{AB} \leq 1,$$

luego $\sum_{n=0}^{\infty} |x_n \bar{y}_n| \leq AB$. Observemos que si $A = 0$ o $B = 0$ una de las sucesiones es idénticamente nula y la conclusión es trivial.

Ahora ya es fácil probar que ℓ^2 es un espacio vectorial. Lo único que no es trivial es que la suma de dos elementos de ℓ^2 está en ℓ^2 , pero

$$\begin{aligned} \sum_{n=0}^{\infty} |x_n + y_n|^2 &\leq \sum_{n=0}^{\infty} (|x_n| + |y_n|)^2 = \sum_{n=0}^{\infty} (|x_n|^2 + |y_n|^2 + 2|x_n||y_n|) \\ &= \sum_{n=0}^{\infty} |x_n|^2 + \sum_{n=0}^{\infty} |y_n|^2 + 2 \sum_{n=0}^{\infty} |x_n \bar{y}_n| \leq \sum_{n=0}^{\infty} |x_n|^2 + \sum_{n=0}^{\infty} |y_n|^2 + \sum_{n=0}^{\infty} |x_n|^2 \sum_{n=0}^{\infty} |y_n|^2, \end{aligned}$$

luego la serie asociada a la suma es convergente.

Una vez probada la convergencia de la serie que define al producto se concluye sin dificultad que dicho producto es ciertamente un producto escalar en ℓ^2 , que es, por consiguiente, un espacio prehilbertiano. La norma asociada es

$$\|x\|_2 = \sqrt{\sum_{n=0}^{\infty} |x_n|^2}.$$

Para probar que ℓ^2 es completo tomamos una sucesión de Cauchy $\{x^n\}_{n=0}^{\infty}$ en ℓ^2 . Así, dado $\epsilon > 0$, existe un n_0 tal que si $m, n \geq n_0$ entonces

$$|x_i^n - x_i^m| \leq \|x^n - x^m\|_2 < \epsilon.$$

Esto implica que las sucesiones $\{x_i^n\}_{i=0}^{\infty}$ son de Cauchy en \mathbb{K} , luego son convergentes, luego podemos tomar $x_i = \lim_i x_i^n \in \mathbb{K}$. Ahora usamos que

$$\sum_{i=0}^N |x_i^n - x_i^m|^2 \leq \sum_{i=0}^{\infty} |x_i^n - x_i^m|^2 = \|x^n - x^m\|_2^2 < \epsilon,$$

siempre que $m, n \geq n_0$. Haciendo tender m a infinito queda

$$\sum_{i=0}^N |x_i^n - x_i|^2 \leq \epsilon^2. \quad (3.2)$$

Por la desigualdad triangular en \mathbb{K}^N tenemos que

$$\sqrt{\sum_{i=0}^N |x_i|^2} \leq \sqrt{\sum_{i=0}^N |x_i^n - x_i|^2} + \sqrt{\sum_{i=0}^N |x_i^n|^2} \leq \epsilon + \|x^n\|_2.$$

Así pues, las sumas parciales de $\sum_{i=0}^{\infty} |x_i|^2$ están acotadas, luego la serie converge y $x \in \ell^2$. Más aún, haciendo tender N a infinito en (3.2) queda $\|x^n - x\|_2^2 \leq \epsilon^2$, para todo $n \geq n_0$, luego $\lim_n x_n = x$. ■

Es inmediato que las sucesiones $e_n \in \ell^2$ dadas por

$$e_{ni} = \begin{cases} 1 & \text{si } n = i, \\ 0 & \text{si } n \neq i \end{cases}$$

son linealmente independientes, por lo que ℓ^2 tiene dimensión infinita. Observemos que no constituyen una base de ℓ^2 en el sentido algebraico pues, por ejemplo, la sucesión $\{\frac{1}{n+1}\}_{n=0}^{\infty} \in \ell^2$ no puede expresarse como combinación lineal de las sucesiones e_n . Sin embargo $\{e_n\}_{n=0}^{\infty}$ es una base de ℓ^2 en un sentido analítico que vamos a introducir en breve. Primero estudiamos las sucesiones ortonormales en un espacio de Hilbert:

Definición 3.45 Sea H un espacio de Hilbert. Una sucesión $\{x_n\}_{n=0}^{\infty}$ en H es *ortogonal* si $\langle x_m, x_n \rangle = 0$ para todo $m \neq n$. Si además $\langle x_n, x_n \rangle = 1$ para todo n diremos que la sucesión es *ortonormal*.

Es claro que la sucesión $\{e_n\}_{n=0}^{\infty}$ que acabamos de introducir es una sucesión ortonormal en ℓ^2 . Además tiene la propiedad de que, para todo $x \in \ell^2$, se cumple que $e_n \cdot x = x_n$.

Obviamente, toda sucesión ortogonal en un espacio de Hilbert que no contenga al 0 da lugar a la sucesión ortonormal $\{x_n/\|x_n\|\}_{n=0}^{\infty}$. Además es linealmente independiente, pues si tenemos una combinación lineal nula

$$\sum_{n=0}^k \lambda_n x_n = 0,$$

multiplicando por x_m obtenemos que $\lambda_m(x_m \cdot x_m) = 0$, luego $\lambda_m = 0$. En particular esto vale para todo sistema ortonormal.

Más en general, si $\{x_n\}_{n=0}^{\infty}$ es una sucesión ortonormal en X y

$$x = \sum_{n=0}^k \lambda_n x_n,$$

al multiplicar por x_n vemos que necesariamente $\lambda_n = x \cdot x_n$, pero todavía podemos enunciar algo más general en esta línea:

Teorema 3.46 Sean x_0, \dots, x_k elementos ortonormales en un espacio de Hilbert H y sea $x \in H$ un elemento arbitrario. Entonces, los escalares $\lambda_n \in \mathbb{K}$ que hacen mínima la expresión

$$\left\| x - \sum_{n=0}^k \lambda_n x_n \right\|$$

son $\lambda_n = x \cdot x_n$. Equivalentemente,

$$\sum_{n=0}^k (x \cdot x_n) x_n$$

es la combinación lineal de x_0, \dots, x_n que mejor aproxima a x .

DEMOSTRACIÓN: Basta desarrollar

$$\begin{aligned} \left\| x - \sum_{n=0}^k \lambda_n x_n \right\|^2 &= \|x\|^2 + \sum_{n=0}^k |\lambda_n|^2 - \sum_{n=0}^k (\bar{\lambda}_n (x \cdot x_n) + \lambda_n \overline{x \cdot x_n}) \\ &= \|x\|^2 - \sum_{n=0}^k |x \cdot x_n|^2 + \sum_{n=0}^k |\lambda_n - x \cdot x_n|^2. \end{aligned}$$

La conclusión es inmediata. ■

Observemos finalmente que si un elemento $x \in H$ admite un desarrollo en serie infinita

$$x = \sum_{n=0}^{\infty} \lambda_n x_n,$$

multiplicando por x_n y usando que el producto escalar es continuo concluimos igualmente que $\lambda_n = \langle x, x_n \rangle$. Ahora bien, nada nos asegura en principio que la serie

$$x = \sum_{n=0}^{\infty} (x \cdot x_n) x_n,$$

vaya a converger ni, en caso de que lo haga, que su suma sea precisamente x . El teorema siguiente garantiza la convergencia de la serie:

Teorema 3.47 (Desigualdad de Bessel) *Sea $\{x_n\}_{n=0}^{\infty}$ una sucesión ortonormal en un espacio de Hilbert H . Entonces, para cada $x \in H$, se cumple que*

$$\sum_{n=0}^{\infty} |x \cdot x_n|^2 \leq \|x\|^2.$$

Por consiguiente, la serie

$$\sum_{n=0}^{\infty} (x \cdot x_n) x_n$$

converge en H .

DEMOSTRACIÓN: En la demostración de 3.46 hemos visto que

$$\left\| x - \sum_{n=0}^k (x \cdot x_n) x_n \right\|^2 = \|x\|^2 - \sum_{n=0}^k |x \cdot x_n|^2,$$

luego

$$\sum_{n=0}^k |(x \cdot x_n)|^2 \leq \|x\|^2.$$

Esto implica que la serie converge y además

$$\sum_{n=0}^{\infty} |x \cdot x_n|^2 \leq \|x\|^2.$$

Llamemos ahora

$$S_k = \sum_{n=0}^k (x \cdot x_n) x_n,$$

de modo que

$$\|S_k - S_m\|^2 = \left\| \sum_{n=m+1}^k (x \cdot x_n) x_n \right\|^2 = \sum_{n=m+1}^k |x \cdot x_n|^2.$$

Esto implica que $\|S_k - S_m\|$ puede hacerse arbitrariamente pequeño, luego $\{S_k\}_k$ es una sucesión de Cauchy, luego converge en H . ■

No obstante, una mínima reflexión nos convence de que, aunque la serie converja, no tiene por qué converger a x . Por ejemplo, si la serie converge a x y eliminamos del sistema ortonormal un elemento tal que $x \cdot x_n \neq 0$, la serie asociada al nuevo sistema “incompleto” seguirá convergiendo, pero ya no a x , es decir, la serie definida a partir de un elemento x por un sistema ortonormal puede no converger a x porque al sistema “le faltan” elementos. Esto nos lleva a la definición siguiente:

Definición 3.48 Una sucesión ortonormal $\{x_n\}_{n=0}^{\infty}$ en un espacio de Hilbert H es *completa* (también se dice que es una *base ortonormal* de H) si todo $x \in H$ se puede escribir de la forma

$$x = \sum_{n=0}^{\infty} \lambda_n x_n,$$

para ciertos $\lambda_n \in \mathbb{K}$, necesariamente únicos, pues han de ser $\lambda_n = x_n \cdot x$ (por la continuidad del producto escalar). Más en general, diremos que una sucesión ortogonal $\{x_n\}_{n=0}^{\infty}$ es *completa*, (o que es una *base ortogonal* de H) si sus elementos son no nulos y la sucesión $\{x_n/\|x_n\|\}_{n=0}^{\infty}$ es una base ortonormal de H .

Así, todo elemento de H se expresa de forma única como combinación lineal infinita de los elementos de una base ortonormal. Para que una base ortonormal en este sentido sea una base en el sentido algebraico es necesario (y suficiente) que sea finita.

Teniendo en cuenta que $e_n \cdot x = x_n$, es inmediato que la sucesión $\{e_n\}_{n=0}^{\infty}$ es una base ortonormal de ℓ^2 . La igualdad

$$\left\| x - \sum_{n=0}^k (x \cdot x_n) x_n \right\|^2 = \|x\|^2 - \sum_{n=0}^k |x \cdot x_n|^2.$$

implica de forma inmediata la caracterización siguiente:

Teorema 3.49 (Identidad de Parseval) Una sucesión ortonormal $\{x_n\}_{n=0}^{\infty}$ en un espacio de Hilbert es completa si y sólo si

$$\|x\|^2 = \sum_{n=0}^{\infty} |x_n \cdot x|^2$$

para todo $x \in H$.

A menudo resulta más útil la caracterización siguiente:

Teorema 3.50 Una sucesión ortogonal (sin términos nulos) en un espacio de Hilbert H es completa si y sólo si todo elemento de H se puede aproximar arbitrariamente por una combinación lineal finita de términos de la sucesión.

DEMOSTRACIÓN: Con más precisión, si $\{x_n\}_{n=0}^{\infty}$ es la sucesión del enunciado (que podemos suponer ortonormal sin pérdida de generalidad), la hipótesis es que, dado $x \in H$ y dado $\epsilon > 0$, existen $\lambda_0, \dots, \lambda_N \in \mathbb{K}$ tales que

$$\left\| x - \sum_{n=0}^N \lambda_n x_n \right\| < \epsilon.$$

El teorema 3.46 implica que, para todo $m > N$,

$$\left\| x - \sum_{n=0}^m (x \cdot x_n) x_n \right\| < \epsilon,$$

y esto implica que

$$x = \sum_{n=0}^{\infty} (x \cdot x_n) x_n. \quad \blacksquare$$

La identidad de Parseval implica también el resultado siguiente:

Teorema 3.51 Si H es un espacio de Hilbert con una base ortonormal (infinita) entonces H es isométrico a ℓ^2 , es decir, existe un isomorfismo $f: H \rightarrow \ell^2$ que conserva el producto escalar.

DEMOSTRACIÓN: Si $\{x_n\}_{n=0}^{\infty}$ es una base ortonormal en H , basta definir $f(x) = \{x \cdot x_n\}_{n=0}^{\infty}$. La identidad de Parseval implica que $f(x) \in \ell^2$. El mismo argumento empleado en la prueba de la desigualdad de Bessel prueba que f es suprayectiva. La continuidad del producto escalar nos da que

$$\begin{aligned} x \cdot y &= \left(\sum_{n=0}^{\infty} (x \cdot x_n) x_n \right) \cdot \left(\sum_{n=0}^{\infty} (y \cdot x_n) x_n \right) = \\ \lim_N \left(\sum_{n=0}^N (x \cdot x_n) x_n \right) \cdot \left(\sum_{n=0}^N (y \cdot x_n) x_n \right) &= \lim_N \sum_{n=0}^N (x \cdot x_n) (\overline{y \cdot x_n}) = \\ \sum_{n=0}^{\infty} (x \cdot x_n) (\overline{y \cdot x_n}) &= f(x) \cdot f(y). \quad \blacksquare \end{aligned}$$

Así pues, todo espacio de Hilbert con una base ortonormal (infinita) es isométrico a ℓ^2 .

Ejercicio: Probar que un espacio de Hilbert tiene un subconjunto denso numerable si y sólo si es isométrico a un espacio \mathbb{K}^n o a ℓ^2 . AYUDA: Si $\{d_n\}_n$ es un conjunto denso, construir una sucesión ortonormal $\{x_n\}_n$ tal que $\langle d_0, \dots, d_n \rangle \subset \langle x_0, \dots, x_n \rangle$. Si la sucesión es infinita aplicar 3.50.

3.6 Espacios de funciones

Uno de los éxitos de la topología consiste en que sus técnicas, desarrolladas en principio para estudiar espacios “geométricos” como \mathbb{R}^n , se aplican igualmente a objetos más abstractos, como son los conjuntos de funciones entre espacios topológicos. Las definiciones siguientes no corresponden en realidad a conceptos nuevos desde un punto de vista topológico:

Definición 3.52 Sea Y un espacio topológico y X un conjunto cualquiera. Una *sucesión funcional* de X en Y es una sucesión $\{f_n\}_{n=0}^\infty$ en el espacio Y^X de todas las aplicaciones de X en Y , es decir, para cada n se cumple $f_n : X \rightarrow Y$.

Si Y es un espacio vectorial topológico (en especial si Y es un cuerpo métrico \mathbb{K}) cada sucesión funcional define la correspondiente *serie funcional* $\sum_{n=0}^\infty f_n$, es decir, la sucesión cuyos términos son las funciones $S_n = \sum_{k=0}^n f_k : X \rightarrow Y$.

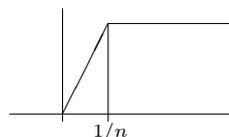
Diremos que una sucesión funcional $\{f_n\}_{n=0}^\infty$ *converge puntualmente* a una función $f \in Y^X$ si para todo $x \in X$ se cumple $\lim_n f_n(x) = f(x)$. En tal caso escribiremos $\lim_n f_n = f$. Para series de funciones podemos definir de manera obvia la convergencia puntual absoluta y la convergencia puntual condicional.

En realidad no estamos introduciendo un nuevo concepto de convergencia. Notemos que Y^X es el producto cartesiano del espacio Y por sí mismo tantas veces como elementos tiene X , luego podemos considerarlo como espacio topológico con la topología producto. Las sucesiones convergen en esta topología si y sólo si convergen coordenada a coordenada, o sea, si y sólo si convergen puntualmente. Por ello a la topología producto en Y^X se la llama también *topología de la convergencia puntual*.

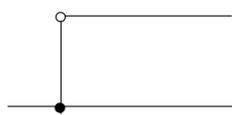
Sin embargo, la convergencia puntual no es la convergencia más natural que puede definirse sobre las sucesiones funcionales. De hecho, presenta grandes inconvenientes.

Ejemplo Para cada $n \geq 1$ sea $f_n : \mathbb{R} \rightarrow \mathbb{R}$ la función dada por

$$f_n(x) = \begin{cases} 0 & \text{si } x \leq 0 \\ nx & \text{si } 0 \leq x \leq 1/n \\ 1 & \text{si } 1/n \leq x \end{cases}$$



Si $x \leq 0$ entonces $f_n(x)$ es constante igual a 0 y si $x > 1/n$ entonces $f_n(x)$ es finalmente constante igual a 1, luego esta sucesión funcional converge puntualmente a la función f que muestra la figura de la izquierda. Tenemos, pues, una sucesión de funciones continuas cuyo límite puntual no es continuo. ■



En general el hecho de que una función sea límite puntual de una sucesión de funciones aporta muy poca información. La razón es que los entornos de la topología puntual son muy grandes. En efecto, en el caso de $\mathbb{R}^{\mathbb{R}}$ no es difícil ver que un entorno básico de una función f es un conjunto de la forma

$$\{g \in \mathbb{R}^{\mathbb{R}} \mid |g(x_i) - f(x_i)| < \epsilon, i = 1, \dots, n\},$$

donde $\epsilon > 0$ y $x_1, \dots, x_n \in \mathbb{R}$, es decir, si una sucesión funcional tiende a f , lo máximo que podemos garantizar tomando un índice grande es que los términos de la sucesión se parecerán a f en un número finito de puntos, pero dos funciones pueden parecerse en un número finito de puntos y ser muy diferentes.

Es mucho más natural considerar que dos funciones están próximas cuando distan menos de un ϵ en todos los puntos a la vez. Por ello, si Y es un espacio métrico, definimos la *topología de la convergencia uniforme* en Y^X como la que tiene por base de entornos abiertos de una función f a los conjuntos de la forma

$$B(f, \epsilon) = \{g \in Y^X \mid d(f(x), g(x)) < \epsilon \text{ para todo } x \in X\}.$$

De este modo, cuando una función g está en un entorno de f suficientemente pequeño, ambas funciones se parecen realmente. Es fácil comprobar que los conjuntos $B(f, \epsilon)$ cumplen las condiciones del teorema 2.12 y por tanto definen, según hemos dicho, una topología en Y^X .

Es inmediato comprobar que una sucesión funcional $\{f_n\}_{n=0}^{\infty}$ converge uniformemente (es decir, en la topología de la convergencia uniforme) a una función f si y sólo si para todo $\epsilon > 0$ existe un n_0 tal que si $n \geq n_0$, entonces $d(f_n(x), f(x)) < \epsilon$ para todo $x \in X$.

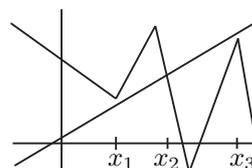
La diferencia, pues, entre la convergencia uniforme y la convergencia puntual es que cuando la convergencia es uniforme hay un n_0 a partir del cual todos los $f_n(x)$ distan de su límite menos de un ϵ dado, mientras que si la convergencia es puntual cada punto x puede requerir un n_0 mayor para acercarse a su límite en menos de ϵ , de manera que ningún n_0 sirva simultáneamente para todos los puntos.

Es obvio que si una sucesión funcional converge uniformemente a una función, también converge puntualmente a dicha función. Tanto la topología de la convergencia uniforme como la topología de la convergencia puntual son de Hausdorff, luego los límites son únicos.

Ejemplo Tomemos $X = [0, 1]$ y consideremos la sucesión de polinomios de $\mathbb{R}[t]$ que cumplen la definición recurrente:

$$P_0(t) = 0, \quad P_{n+1}(t) = P_n(t) + \frac{1}{2}(t - P_n(t)^2).$$

Vamos a probar que converge uniformemente a la función \sqrt{t} en X .



Para ello observamos en primer lugar que $P_n(t) \leq \sqrt{t}$, para todo $t \in X$. Esto es trivialmente cierto para $n = 0$ y, si vale para n , entonces

$$\begin{aligned} \sqrt{t} - P_{n+1}(t) &= \sqrt{t} - P_n(t) - \frac{1}{2}(\sqrt{t} - P_n(t))(\sqrt{t} + P_n(t)) \\ &= (\sqrt{t} - P_n(t))\left(1 - \frac{1}{2}(\sqrt{t} + P_n(t))\right) \geq (\sqrt{t} - P_n(t))\left(1 - \frac{1}{2}2\sqrt{t}\right) \geq 0. \end{aligned}$$

Teniendo esto en cuenta, la definición recurrente implica de forma inmediata que la sucesión $\{P_n(t)\}_{n=0}^\infty$ es monótona creciente (para cada $t \in X$). Como además está acotada superiormente por \sqrt{t} , tenemos que es convergente. Llamemos $f(t) = \lim_n P_n(t)$.

Tomando límites en la definición recurrente de la sucesión obtenemos que $f(t) = f(t) + \frac{1}{2}(t - f(t)^2)$, de donde concluimos que $f(t) = \sqrt{t}$.

Con esto hemos probado que la sucesión de funciones $\{P_n\}_{n=0}^\infty$ converge puntualmente a \sqrt{t} . Ahora vamos a probar que la convergencia es uniforme. Para ello, dado $\epsilon > 0$, observamos que los conjuntos

$$F_n = \{t \in X \mid f(t) - P_n(t) \geq \epsilon\}$$

son cerrados y, como la sucesión de polinomios es creciente, los conjuntos forman una sucesión decreciente respecto de la inclusión. La convergencia puntual implica que $\bigcap_n F_n = \emptyset$. Como X es compacto, la familia de los cerrados F_n no puede tener la propiedad de la intersección finita, lo cual significa que existe un n_0 tal que, para todo $n \geq n_0$, $F_n = \emptyset$. A su vez, esto significa que todo $t \in X$ cumple $|f(t) - P_n(t)| < \epsilon$, es decir, $d(f(t), P_n(t)) < \epsilon$. ■

Si en el espacio Y tomamos la distancia $d'(x, y) = \min\{1, d(x, y)\}$, los conjuntos $B(f, \epsilon)$ son los mismos cuando $\epsilon < 1$, luego d' induce la misma topología de convergencia uniforme en Y^X .

La ventaja de esta métrica es que no toma valores mayores que 1, por lo que podemos definir

$$d(f, g) = \sup\{d'(f(x), g(x)) \mid x \in X\},$$

y es claro que esta d es una distancia en Y^X para la cual $B(f, \epsilon)$ es simplemente la bola abierta de centro f y radio ϵ . Esto convierte a Y^X en un espacio métrico cuya topología es precisamente la de la convergencia uniforme.

Teorema 3.53 *Sea X un espacio topológico e Y un espacio métrico. Entonces el conjunto $C(X, Y)$ de las aplicaciones continuas de X en Y es cerrado en Y^X , cuando en éste consideramos la topología de la convergencia uniforme.*

DEMOSTRACIÓN: Sea $\{f_n\}_{n=0}^\infty$ una sucesión de aplicaciones continuas que converja uniformemente a una función f . Basta ver que f es continua. Sea $\epsilon > 0$. Sea n_0 tal que si $n \geq n_0$ y $x \in X$, entonces $d(f_n(x), f(x)) < \epsilon/3$.

Sea $x_0 \in X$ (vamos a probar que f es continua en x_0). Como f_{n_0} es continua existe un entorno U de x_0 tal que si $x \in U$, entonces $d(f_{n_0}(x), f_{n_0}(x_0)) < \epsilon/3$. Por lo tanto, si $x \in U$, se cumple que

$$d(f(x), f(x_0)) \leq d(f(x), f_{n_0}(x)) + d(f_{n_0}(x), f_{n_0}(x_0)) + d(f_{n_0}(x_0), f(x_0)) < \epsilon. \quad \blacksquare$$

Éste es un primer ejemplo del buen comportamiento de la convergencia uniforme. Veamos otros:

Teorema 3.54 *Sea Y un espacio métrico completo y X un conjunto. Entonces Y^X es completo con la métrica inducida a partir de la métrica de Y . Por lo tanto si X es un espacio topológico, $C(X, Y)$ también es completo.*

DEMOSTRACIÓN: Ante todo notemos que si Y es completo con su métrica d , también lo es con la métrica d' que resulta de tomar el mínimo con 1, luego podemos suponer que la métrica en Y está acotada.

Sea $\{f_n\}_{n=0}^\infty$ una sucesión de Cauchy en Y^X . Esto significa que para todo $\epsilon > 0$ existe un n_0 tal que si $m, n \geq n_0$ y $x \in X$, entonces $d(f_m(x), f_n(x)) < \epsilon$.

En particular esto implica que cada sucesión $\{f_n(x)\}_{n=0}^\infty$ es de Cauchy en Y , luego converge a un cierto punto $f(x)$. Con esto tenemos una función $f \in Y^X$ a la cual $\{f_n\}_{n=0}^\infty$ converge puntualmente. Basta ver que también converge uniformemente.

Sea $\epsilon > 0$ y tomemos un natural n_0 como antes. Así, si $n_0 \leq n \leq m$, se cumple $d(f_m(x), f_n(x)) < \epsilon$, luego $f_m(x)$ está en la bola cerrada de centro $f_n(x)$ y radio ϵ , luego el límite $f(x)$ estará en esta misma bola, o sea, se cumplirá $d(f(x), f_n(x)) \leq \epsilon$, y esto para todo $n \geq n_0$ y todo $x \in X$. Esto significa que la sucesión converge uniformemente a f . La completitud de $C(X, Y)$ se sigue de que es un cerrado. \blacksquare

En general no podemos convertir a Y^X en un espacio normado aunque Y lo sea. El problema es que no podemos transformar la norma en una norma acotada. Lo único que podemos hacer es definir $(Y^X)^*$ como el espacio de las funciones acotadas de X en Y , es decir, las funciones f tales que $f[X]$ está acotado. En este conjunto podemos definir la *norma supremo* dada por $\|f\|_\infty = \sup\{\|f(x)\| \mid x \in X\}$, que obviamente genera las bolas que definen la topología de la convergencia uniforme (restringida a $(Y^X)^*$). Así pues, si Y es un espacio normado, $(Y^X)^*$ también lo es, y como es fácil ver que el límite uniforme de funciones acotadas está acotado, resulta que $(Y^X)^*$ es cerrado en Y^X , luego si Y es un espacio de Banach, $(Y^X)^*$ también lo es.

Si X es un espacio topológico e Y es un espacio normado, definimos $C^*(X, Y)$ como el espacio de las funciones continuas y acotadas de X en Y . Obviamente se trata de la intersección de dos cerrados, luego es cerrado, es un espacio normado y si Y es un espacio de Banach, $C^*(X, Y)$ también lo es.

Por el teorema 3.4, tenemos que si X es compacto entonces $C^*(X, Y) = C(X, Y)$, luego en este caso sí que tenemos una estructura de espacio normado en $C(X, Y)$, con la topología de la convergencia uniforme.

En particular, si \mathbb{K} es un cuerpo métrico completo, $C^*(X, \mathbb{K})$ es un espacio de Banach. En general no podemos dotar a $C(X, \mathbb{K})$ de estructura de espacio normado. De hecho, es fácil ver que $C^*(X, \mathbb{K})$ es abierto y cerrado en $C(X, \mathbb{K})$, luego, salvo que $C(X, \mathbb{K}) = C^*(X, \mathbb{K})$, tenemos que $C(X, \mathbb{K})$ no es conexo, luego ni siquiera es un espacio vectorial topológico.

Sea $L(E, F)$ el conjunto de las aplicaciones lineales continuas entre dos espacios normados. Teniendo en cuenta 3.35, la aplicación $L(E, F) \xrightarrow{\text{restricción}} C^*(\overline{B_1(0)}, F)$ definida por restricción es claramente lineal e inyectiva (pues $B_1(0)$ contiene una base de E). Si transportamos la norma de este segundo espacio al primero obtenemos que, para cada aplicación lineal y continua $f : E \rightarrow F$, su norma es el supremo de f en $\overline{B_1(0)}$, y por lo tanto cumple $\|f(v)\| \leq \|f\| \|v\|$ para todo $v \in E$.

Ejercicio: Probar que $L(E, F)$ es un subespacio cerrado de $C^*(\overline{B_1(0)}, F)$. Por consiguiente, si F es un espacio de Banach, $L(E, F)$ también lo es.

Sea $K \subset \mathbb{R}^n$ un subconjunto compacto infinito. Entonces podemos considerar al anillo de polinomios $\mathbb{R}[x_1, \dots, x_n]$ como subespacio del espacio $C(K)$ de todas las funciones reales continuas en K . Vamos a probar que es denso en $C(K)$, es decir, que toda función continua en K puede aproximarse uniformemente por un polinomio. En realidad probamos un resultado más general:

Teorema 3.55 (Stone-Weierstrass) *Sea X un espacio compacto y $A \subset C(X)$ una subálgebra que contenga a la función 1 y que separe puntos de X (es decir, que para cada par de puntos distintos $x, y \in X$ exista $f \in A$ tal que $f(x) \neq f(y)$). Entonces A es densa en $C(X)$ respecto a la topología de la convergencia uniforme.*

DEMOSTRACIÓN: Si \overline{A} es la clausura de A respecto de la topología de la convergencia uniforme, es fácil ver que \overline{A} sigue cumpliendo las hipótesis del teorema. Por ejemplo, dadas $f, g \in \overline{A}$, existen sucesiones $\{f_n\}_n, \{g_n\}_n$ en A que convergen a f y g , respectivamente, y es fácil ver que $\{f_n + g_n\}_n$ y $\{f_n g_n\}_n$ convergen a $f + g$ y fg , respectivamente, luego $f + g, fg \in \overline{A}$. Esto implica que \overline{A} es también una subálgebra de $C(X)$.

Así pues, cambiando A por \overline{A} podemos suponer que A es cerrada en $C(X)$, y ahora tenemos que probar que $A = C(X)$. Observemos que el hecho de que A sea una subálgebra y contenga a la función constante 1 implica que, de hecho, contiene a todas las constantes.

En primer lugar demostramos que si $f \in A$, entonces $|f| \in A$. Sea $c \in \mathbb{R}^+$ tal que $|f(x)| \leq c$ para todo $x \in X$. Entonces $f/c \in A$ y basta probar que $|f/c| \in A$. Equivalentemente, podemos suponer que $|f(x)| \leq 1$ para todo $x \in X$.

Sea $\{P_n\}_{n=0}^\infty$ la sucesión de polinomios considerada en el ejemplo de la página 128 y sea $f_n(x) = P_n(f(x)^2)$. El hecho de que A sea un álgebra que contiene a las constantes implica que $f_n(x) \in A$, y el hecho de que la sucesión de polinomios converja uniformemente a \sqrt{t} en $[0, 1]$ implica claramente que $\{f_n\}_{n=0}^\infty$ converge uniformemente a $\sqrt{f^2} = |f|$ en X . Como A es cerrado, $|f| \in A$.

A continuación vemos que si $f, g \in A$, entonces $\max(f, g)$ y $\min(f, g) \in A$.

Basta observar las expresiones:

$$\min(f, g) = \frac{1}{2}(f + g) - \frac{1}{2}|f - g|, \quad \max(f, g) = \frac{1}{2}(f + g) + \frac{1}{2}|f - g|.$$

Fijemos $f \in C(X)$ y tomemos dos puntos distintos $x, y \in X$. Por hipótesis existe $h \in A$ tal que $h(x) \neq h(y)$. Podemos encontrar números reales α, β tales que la función $f_{xy} = \alpha h + \beta 1$ cumpla $f_{xy}(x) = f(x)$, $f_{xy}(y) = f(y)$. Claramente $f_{xy} \in A$.

Sea B_y un entorno abierto de y tal que $f_{xy}(z) < f(z) + \epsilon$ para todo $z \in B_y$. Como X es compacto podemos cubrirlo por un número finito de estos abiertos, digamos B_{y_1}, \dots, B_{y_n} . Sea $f_x = \min(f_{y_1}, \dots, f_{y_n}) \in A$. Se cumple entonces que $f_x(x) = f(x)$ y $f_x(z) < f(z) + \epsilon$ para todo $z \in X$. Sea B_x un entorno abierto de x tal que $f_x(z) > f(z) - \epsilon$ para todo $z \in B_x$. Tomamos un subcubrimiento finito, digamos B_{x_1}, \dots, B_{x_m} y llamamos $g = \max(f_{x_1}, \dots, f_{x_m}) \in A$, de modo que $g(z) > f(z) - \epsilon$ para todo $z \in X$. También cumple que $g(z) < f(z) + \epsilon$, luego $\|f - g\| \leq \epsilon$. Como ϵ es arbitrario, concluimos que $f \in A$. ■

Por ejemplo, si $K \subset \mathbb{R}^n$ es un compacto infinito, es inmediato que el álgebra de los polinomios $\mathbb{R}[x_1, \dots, x_n]$, vista como subálgebra de $C(K)$ cumple las hipótesis del teorema anterior, luego si $f \in C(K)$, dado $\epsilon > 0$ existe un polinomio p tal que $|f(x) - p(x)| < \epsilon$, para todo $x \in K$.

Veamos ahora un resultado importante sobre convergencia de series funcionales. Previamente notemos lo siguiente: si una serie $\sum_{n=0}^{\infty} f_n$ con valores en un cuerpo métrico completo converge absoluta y uniformemente en un conjunto X , es decir, si la serie $\sum_{n=0}^{\infty} |f_n|$ converge uniformemente, entonces $\sum_{n=0}^{\infty} f_n$ converge uniformemente. La prueba es la misma que la de 2.86.

Teorema 3.56 (Criterio de Mayoración de Weierstrass) *Sea \mathbb{K} un cuerpo métrico completo y sea $\sum_{n=0}^{\infty} f_n$ una serie funcional en un espacio X , de modo que cada $f_n : X \rightarrow \mathbb{K}$. Sea $\{M_n\}_{n=0}^{\infty}$ una sucesión en el intervalo $[0, +\infty[$ tal que para todo natural n y todo $x \in X$ se cumpla $|f_n(x)| \leq M_n$. Si la serie $\sum_{n=0}^{\infty} M_n$ es convergente, entonces la serie $\sum_{n=0}^{\infty} f_n$ es absoluta y uniformemente convergente en X .*

DEMOSTRACIÓN: La serie M_n es de Cauchy, luego dado $\epsilon > 0$ existe un n_0 tal que si $n_0 \leq m \leq p$, entonces $\sum_{n=m}^p M_n < \epsilon$. Así

$$\left| \sum_{n=m}^p |f_n(x)| \right| = \sum_{n=m}^p |f_n(x)| \leq \sum_{n=m}^p M_n < \epsilon$$

para todo $x \in X$. Esto significa que la serie $\sum_{n=0}^{\infty} |f_n|$ es de Cauchy en $C(X, \mathbb{K})$, luego (uniformemente) convergente, luego $\sum_{n=0}^{\infty} f_n$ es absoluta y uniformemente convergente. ■

Para terminar caracterizamos los subespacios compactos de un espacio de funciones continuas $C(X, Y)$, con X compacto:

Teorema 3.57 (Ascoli-Arzelà) *Sea X un espacio métrico compacto, sea Y un espacio métrico completo y $A \subset C(X, Y)$. Entonces A es relativamente compacto (es decir, tiene clausura compacta) respecto de la topología de la convergencia uniforme si y sólo si cumple las dos condiciones siguientes:*

- A es equicontinuo, es decir, para todo $\epsilon > 0$ existe un $\delta > 0$ tal que si $f \in A$ y $x, y \in X$ cumplen $d(x, y) < \delta$, entonces $d(f(x), f(y)) < \epsilon$.*
- Para cada $x \in X$, el conjunto $\{f(x) \mid f \in A\}$ es relativamente compacto en Y .*

DEMOSTRACIÓN: Para cada número natural $n \geq 1$, podemos tomar un conjunto $D_n = \{x_1, \dots, x_k\} \subset X$ tal que $X = B_{1/n}(x_1) \cup \dots \cup B_{1/n}(x_k)$. Llamamos $D = \{d_n \mid n \in \mathbb{N}\}$ a la unión de todos los conjuntos D_n .

Supuesto que A sea equicontinuo y uniformemente acotado, basta probar que toda sucesión $\{f_n\}_{n=0}^{\infty}$ en A tiene una subsucesión (uniformemente) convergente, pues en tal caso lo mismo vale para \bar{A} . En efecto, si $\{g_n\}_{n=0}^{\infty}$ es una sucesión en \bar{A} , para cada n podemos tomar $f_n \in A$ tal que $d(f_n, g_n) < 1/(n+1)$, entonces $\{f_n\}_{n=0}^{\infty}$ tiene una subsucesión tal que existe $\lim_n f_n = f \in C(X, Y)$, y es claro entonces que también $\lim_n g_n = f$.

Dada $\{f_n\}_{n=0}^{\infty}$ en A , veamos que tiene una subsucesión que converge puntualmente en cada punto de D . Para ello observamos que la sucesión $\{f_n(d_0)\}_{n=0}^{\infty}$ está contenida en un subconjunto compacto de Y , luego tiene una subsucesión convergente, digamos $\{f_{n_0, k}(d_0)\}_{k=0}^{\infty}$. Similarmente, de $\{f_{n_0, k}\}_{k=0}^{\infty}$ podemos extraer una subsucesión $\{f_{n_1, k}\}_{k=0}^{\infty}$ tal que $\{f_{n_1, k}(d_1)\}_{k=0}^{\infty}$ sea convergente.

Procediendo de este modo obtenemos sucesiones $\{f_{n_m, k}\}_{k=0}^{\infty}$, y podemos considerar entonces la subsucesión $\{f_{n_k, k}\}_{k=0}^{\infty}$, que es, de hecho, una subsucesión de todas las anteriores, luego converge en todos los puntos de D . A partir de aquí la representaremos por $\{f_{n_k}\}_{k=0}^{\infty}$.

Dado $\epsilon > 0$, por la equicontinuidad existe un $\delta > 0$ tal que si $d(x, y) < \delta$, entonces $d(f_{n_k}(x), f_{n_k}(y)) < \epsilon/3$, para todo k . Tomemos un número natural $r > 1/\delta$. Como la subsucesión converge en todos los puntos de D_r , existe un $k_0 > 0$ tal que si $k, l \geq k_0$ y $d \in D_r$ entonces $d(f_{n_k}(d), f_{n_l}(d)) < \epsilon/3$.

Ahora, dado $x \in X$, tenemos que existe un $d \in D_r$ tal que $x \in B_{1/r}(d)$, luego $d(x, d) < \delta$, luego, para $k, l \geq k_0$,

$$d(f_{n_k}(x), f_{n_l}(x)) \leq d(f_{n_k}(x), f_{n_k}(d)) + d(f_{n_k}(d), f_{n_l}(d)) + d(f_{n_l}(d), f_{n_l}(x)) < \epsilon.$$

Esto significa que la sucesión $\{f_{n_k}\}_{k=0}^\infty$ es de Cauchy en $C(X, Y)$, que es completo, luego es convergente.

Veamos el recíproco: Si A es relativamente compacto, dado $\epsilon > 0$ para cada $g \in \bar{A}$ existe un $f \in A$ tal que $g \in B_{\epsilon/3}(f)$, luego estas bolas cubren el compacto \bar{A} , luego existen $f_1, \dots, f_k \in A$ tales que $A \subset B_{\epsilon/3}(f_1) \cup \dots \cup B_{\epsilon/3}(f_k)$.

Por lo tanto, si $x \in X$, tenemos que

$$\{f(x) \mid f \in A\} \subset B_{\epsilon/3}(f_1(x)) \cup \dots \cup B_{\epsilon/3}(f_k(x)),$$

luego

$$\overline{\{f(x) \mid f \in A\}} \subset B_\epsilon(f_1(x)) \cup \dots \cup B_\epsilon(f_k(x)).$$

Esto prueba que la clausura es precompacta y, como es cerrada en Y , que es un espacio métrico completo, también es completa, luego es compacta, por 3.30.

Por otra parte, como X es compacto, cada f_i es uniformemente continua en X , luego existe un $\delta > 0$ tal que si $x, y \in X$ cumplen $d(x, y) < \delta$, entonces $d(f_i(x), f_i(y)) < \epsilon/3$. Así, si $f \in A$, existe un i tal que $d(f, f_i) < \epsilon/3$, luego si $d(x, y) < \delta$, tenemos que

$$d(f(x), f(y)) \leq d(f(x), f_i(x)) + d(f_i(x), f_i(y)) + d(f_i(y), f(y)) < \epsilon.$$

Esto prueba que A es equicontinuo. ■

Nota Si $Y = \mathbb{R}^n$, sus subconjuntos relativamente compactos son simplemente sus subconjuntos acotados, por lo que la condición b) del teorema anterior equivale a que A sea *puntualmente acotado*, es decir, a que todos los conjuntos $\{f(x) \mid f \in A\}$ estén acotados. Sin embargo, si A es compacto, se cumple algo más fuerte, y es que A está acotado respecto de la norma de $C(X, \mathbb{R}^n)$, lo que implica trivialmente la acotación puntual. Por lo tanto, si X es un espacio topológico compacto, un subconjunto $A \subset C(X, \mathbb{R}^n)$ es compacto con la topología de la convergencia uniforme si y sólo si es cerrado, acotado y equicontinuo. ■

3.7 El teorema de Baire

Terminamos el capítulo con un resultado que no nos va a hacer falta más adelante, pero que es útil en algunos contextos más avanzados. Necesitamos un resultado previo:

Teorema 3.58 *Sea M un espacio métrico completo. Toda familia decreciente $\{C_n\}_n$ de cerrados en M no vacíos cuyos diámetros cumplan $\lim_n d(C_n) = 0$ tiene intersección no vacía.*

DEMOSTRACIÓN: Para cada n tomamos $x_n \in C_n$. Como $\lim_n d(C_n) = 0$, es claro que la sucesión $(x_n)_n$ es de Cauchy. Su límite x está en cada C_n por ser éste cerrado y $\{C_n\}_n$ decreciente. ■

Teorema 3.59 (Teorema de Baire) *En un espacio métrico completo, la intersección de una familia numerable de abiertos densos es un conjunto denso.*

DEMOSTRACIÓN: Sea M un espacio métrico completo y $G = \bigcap_{n=1}^{\infty} G_n$ una intersección numerable de abiertos G_n densos en M . Basta probar que G corta a toda bola abierta $B_r(x)$. Como G_1 es denso en M existe $x_1 \in G_1 \cap B_r(x)$. Como $G_1 \cap B_r(x)$ es abierto, existe un $r_1 > 0$, que podemos tomar menor que $r/2$, tal que $\overline{B_{r_1}(x_1)} \subset G_1 \cap B_r(x)$.

Inductivamente podemos construir una sucesión $\{x_n\}_n$ de puntos de M y una sucesión $\{r_n\}_n$ de números reales positivos de modo que

$$\overline{B_{r_n}(x_n)} \subset G_n \cap B_{r_{n-1}}(x_{n-1}) \quad (3.3)$$

y $r_n < r/n$ para todo n .

Por el teorema anterior, $\bigcap_{n=1}^{\infty} \overline{B_{r_n}(x_n)} \neq \emptyset$, luego (3.3) implica que

$$G \cap B_r(x) \supset \bigcap_{n=1}^{\infty} (G_n \cap B_{r_{n-1}}(x_{n-1})) \neq \emptyset. \quad \blacksquare$$

Sin más que tener en cuenta que el complementario de un conjunto denso es un conjunto con interior vacío tenemos una forma equivalente del teorema de Baire:

Teorema 3.60 (Teorema de Baire) *En un espacio métrico completo, toda unión numerable de cerrados de interior vacío tiene interior vacío.*

Conviene observar que la tesis del teorema de Baire (en cualquiera de sus dos formas equivalentes) se cumple también sobre espacios topológicos localmente compactos, no necesariamente metrizable. La prueba es, de hecho, más sencilla, y se obtiene sustituyendo el teorema 3.58 por el hecho de que la intersección de una familia decreciente de compactos no vacíos es no vacía.

Aunque, según hemos indicado, no necesitaremos el teorema de Baire, vamos a tratar de explicar su interés. Para ello necesitamos algunas definiciones:

Definición 3.61 Sea X un espacio topológico. Un subconjunto A de X es *diseminado* si $X \setminus A$ contiene un abierto denso o, equivalentemente, si A está contenido en un cerrado de interior vacío o, también, si $\text{int } \overline{A} = \emptyset$.

Informalmente, la idea es que un abierto denso (y cualquier conjunto que lo contenga) es un conjunto “muy grande” desde el punto de vista topológico, pues todo abierto contiene un abierto contenido en tal conjunto; los conjuntos diseminados son topológicamente “muy pequeños”. Como todo conjunto que contenga a un conjunto que contenga a un abierto denso contiene un abierto denso, tomando complementarios obtenemos que los subconjuntos de los conjuntos diseminados son diseminados. Esta noción de conjunto diseminado resulta ser muy restrictiva, esencialmente a causa de que no se conserva por uniones numerables, por ello se definen los conjuntos de primera categoría:

Definición 3.62 Un subconjunto A de un espacio topológico X es de *primera categoría* si es unión numerable de conjuntos diseminados. A es de *segunda categoría* si no es de primera categoría.

Es evidente que toda unión numerable de conjuntos de primera categoría es de primera categoría. Así, los conjuntos de primera categoría son conjuntos topológicamente “pequeños”, aunque no necesariamente “muy pequeños”, mientras que los conjuntos de segunda categoría son los topológicamente “grandes”.

No obstante, estas nociones no sirven de nada sin el teorema de Baire, que puede enunciarse en una tercera forma equivalente:

Teorema 3.63 (Teorema de Baire) *En un espacio métrico completo, los conjuntos de primera categoría tienen interior vacío.*

DEMOSTRACIÓN: Consideremos un conjunto C de primera categoría. Entonces

$$C = \bigcup_n A_n \subset \bigcup_n \bar{A}_n = C',$$

donde los conjuntos A_n son diseminados, luego sus clausuras son cerrados de interior vacío, luego C' tiene interior vacío (por la versión que ya hemos probado del teorema de Baire) y C también. ■

Así pues, si probamos que un conjunto C es “pequeño”, en el sentido de que es de primera categoría, el teorema de Baire nos da que todo abierto va a contener puntos que no están en C . Éste es esencialmente el interés del teorema de Baire.

Terminaremos con una aplicación del teorema de Baire, que no es de las más típicas, pero tal vez la más sencilla:

Ejemplo \mathbb{Q} no puede expresarse como una intersección numerable de abiertos de \mathbb{R} .

En efecto, en tal caso sería una intersección numerable de abiertos densos, luego $\mathbb{R} \setminus \mathbb{Q}$ sería una unión numerable de cerrados de interior vacío, al igual que $\mathbb{R} = \mathbb{Q} \cup (\mathbb{R} \setminus \mathbb{Q})$. El teorema de Baire nos daría entonces que \mathbb{R} tiene interior vacío (en sí mismo), lo cual es absurdo. ■

Los conjuntos que pueden expresarse como intersección numerable de abiertos se llaman conjuntos G_δ . El ejemplo anterior, junto con el teorema siguiente, muestra que no puede existir una función $f : \mathbb{R} \rightarrow \mathbb{R}$ continua en los puntos de \mathbb{Q} y discontinua en los de $\mathbb{R} \setminus \mathbb{Q}$. (Si el lector cree que esto es evidente, debería pensar en el ejercicio que sigue al teorema.)

Teorema 3.64 *Sea M un espacio métrico completo. El conjunto de puntos de continuidad de toda función $f : M \rightarrow \mathbb{R}$ es un G_δ .*

DEMOSTRACIÓN: Para cada natural no nulo n definimos

$$G_n = \{x \in M \mid \text{existe } \delta > 0 \text{ tal que } \sup_{y \in B_\delta(x)} f(y) - \inf_{y \in B_\delta(x)} f(y) < 1/n\}.$$

Claramente los conjuntos G_n son abiertos. Basta probar que el conjunto de puntos de continuidad de f es $G = \bigcap_{n=1}^{\infty} G_n$.

Si f es continua en x , dado un $n > 0$ existe un $\delta > 0$ tal que si $y \in B_\delta(x)$ entonces $|f(x) - f(y)| < 1/(4n)$. Por lo tanto, si $y, y' \in B_\delta(x)$ se cumple que $|f(y) - f(y')| < 1/(2n)$ y, tomando el supremo en y y el ínfimo en y' , concluimos que $x \in G_n$.

Recíprocamente, supongamos que $x \in G$. Dado $\epsilon > 0$ tomamos n tal que $1/n < \epsilon$. Como $x \in G_n$, existe $\delta > 0$ tal que

$$\sup_{y \in B_\delta(x)} f(y) - \inf_{y \in B_\delta(x)} f(y) < 1/n.$$

Entonces si $y \in B_\delta(x)$ se tiene que

$$|f(x) - f(y)| \leq \sup_{y \in B_\delta(x)} f(y) - \inf_{y \in B_\delta(x)} f(y) < 1/n < \epsilon,$$

luego f es continua en x . ■

Ejercicio: Consideremos la función $f: \mathbb{R} \rightarrow \mathbb{R}$ dada por

$$f(x) = \begin{cases} 1/q & \text{si } x = p/q \text{ con } p, q \in \mathbb{Z}, (p, q) = 1, q > 0, \\ 0 & \text{en caso contrario.} \end{cases}$$

Demostrar que $\lim_{x \rightarrow x_0} f(x) = 0$ para todo $x_0 \in \mathbb{R}$. Deducir que f es continua en $\mathbb{R} \setminus \mathbb{Q}$ y discontinua en \mathbb{Q} .

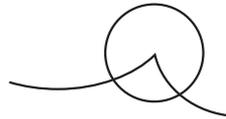
Capítulo IV

Cálculo diferencial de una variable

En este capítulo estudiaremos una de las ideas más importantes y fructíferas que posee la matemática actual. Su núcleo está en la observación de que muchas curvas se parecen localmente a rectas. Por ejemplo, visto suficientemente de cerca, un arco de circunferencia es indistinguible de un segmento de recta. Una prueba de ello está en el horizonte que separa el cielo del mar en un día despejado. Se trata de un arco de circunferencia, pero ¿se ve que es así?

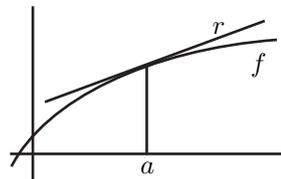
La razón por la que el horizonte parece recto no es que la Tierra sea muy grande, sino que la vemos muy de cerca. Vista desde la Luna es claramente esférica y cualquier circunferencia al microscopio parece una recta. Lo mismo sucede con muchas curvas, que si las vemos muy de cerca parecen rectas. Pero esto no es cierto para todas. Pensemos en la curva de la figura.

Vista de cerca podría pasar por una recta en un entorno de cualquiera de sus puntos excepto el señalado con el círculo, donde tiene un “pico”. Por más que nos acerquemos nunca dejaremos de ver ese pico que delatará que no se trata de una recta. Vamos a estudiar las curvas que localmente se parecen a rectas. Pero ¿qué es exactamente parecerse a una recta?



4.1 Derivación

Consideremos una función f definida en un entorno de un punto $a \in \mathbb{R}$ y con imagen en \mathbb{R} . Supongamos que su gráfica se parece mucho a una recta en un entorno del punto. La pregunta es ¿a qué recta se parece? Por lo pronto a una que pasa por el punto $(a, f(a))$. Las rectas que pasan por dicho punto (a excepción de la vertical, que no nos va a interesar) son de la forma $r(x) = m(x - a) + f(a)$, donde $m \in \mathbb{R}$.



La interpretación geométrica de m es sencilla. En general, si tenemos una recta $r(x) = mx + n$, en un punto a tomará el valor $r(a) = ma + n$. Si nos trasladamos a un punto $a + h$ con $h \neq 0$ obtendremos $r(a + h) = m(a + h) + n$. El incremento que ha experimentado la función es $r(a + h) - r(a) = mh$, y si lo dividimos por el desplazamiento h obtenemos, independientemente de cuál sea h , el valor m . Es decir,

$$m = \frac{r(a + h) - r(a)}{h}.$$

Geométricamente esto no es sino el teorema de Tales. En definitiva, m expresa lo que aumenta la función por unidad de avance: si nos desplazamos $h = 1$ unidad, la recta aumenta en m unidades, si avanzamos $h = 2$ unidades, la recta aumenta $2m$, etc. Por lo tanto, si el valor de m es grande la recta subirá muy rápidamente, será una recta muy empinada. Si $m = 0$ la recta no sube, es horizontal. Si m es negativo la recta baja, más rápidamente cuanto mayor sea m en módulo. El número m se llama pendiente de la recta. Una recta viene determinada por dos de sus puntos o bien por uno de sus puntos y su pendiente (pues conocido un punto (a, b) y la pendiente m conocemos más puntos: $(a + 1, b + m)$, por ejemplo).

Volviendo a nuestro problema, tenemos la recta $r(x) = m(x - a) + f(a)$, cuya pendiente es m . Nos falta determinar m para que sea la recta que se parece a f . Consideremos la expresión

$$m(h) = \frac{f(a + h) - f(a)}{h}. \quad (4.1)$$

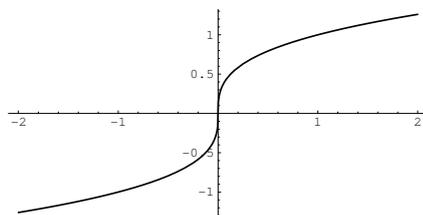
Esto no es una constante (salvo que f sea una recta), pero si ciertamente f se parece a una recta r , esta expresión debería parecerse a la pendiente de r . Si f se parece más a r cuanto más de cerca la miramos, esto es, cuando consideramos puntos más cercanos al punto a , el valor $m(h)$ debería parecerse más a la pendiente de r cuanto menor es h . Por ello definimos:

Definición 4.1 Sea $f : A \rightarrow \mathbb{R}$ y a un punto interior de A . Diremos que f es *derivable* en a si existe (en \mathbb{R})

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a + h) - f(a)}{h}.$$

Cuando esto sucede, a la recta $r(x) = f'(a)(x - a) + f(a)$ se le llama *recta tangente* a f en el punto $(a, f(a))$ (o para abreviar, en el punto a). El número $f'(a)$ es la *derivada* de f en el punto a .

Según hemos dicho, una función es derivable en un punto cuando su gráfica se confunde en un entorno de dicho punto con la de una recta, la recta tangente a la función en el punto. Esto no es exacto, pues en realidad hay funciones que se parecen a rectas en los alrededores de un punto y pese a ello no son derivables. Esto ocurre cuando la recta tangente es vertical, con lo que su pendiente es infinita y no existe (en \mathbb{R}) el límite que define la derivada. Un ejemplo lo proporciona la función $\sqrt[3]{x}$ en $x = 0$.



Observamos que el eje vertical es tangente a la función en 0, pese a lo cual, según veremos, la función no es derivable en 0.

Diremos que una función es *derivable* en un abierto A si es derivable en todos los puntos de A . Una función es *derivable* si su dominio es un abierto y es derivable en todos sus puntos.

Si $f : A \rightarrow \mathbb{R}$ es derivable, tenemos definida otra función $f' : A \rightarrow \mathbb{R}$ que a cada punto $a \in A$ le asigna su derivada $f'(a)$. A esta función la llamamos (función) *derivada* de f en A .

Teniendo en cuenta la motivación que hemos dado para el concepto de derivada, es claro que toda recta no vertical, $f(x) = mx + n$ es derivable en \mathbb{R} y su derivada es su pendiente, o sea, m . La razón es que, según hemos visto, el cociente (4.1) es en este caso constante igual a m , luego el límite cuando h tiende a 0 es igualmente m . En particular, la derivada de una función constante, $f(x) = a$, es $f'(x) = 0$.

Ejemplo Calculemos la derivada de la función $f(x) = x^2$.

$$f'(x) = \lim_{h \rightarrow 0} \frac{(x+h)^2 - x^2}{h} = \lim_{h \rightarrow 0} \frac{x^2 + 2xh + h^2 - x^2}{h} = \lim_{h \rightarrow 0} 2x + h = 2x. \quad \blacksquare$$

Enseguida veremos que es muy fácil reconocer las funciones derivables, así como calcular sus derivadas. Primero demostremos un hecho básico. Obviamente, lo primero que ha de hacer una función para parecerse a una recta es ser continua.

Teorema 4.2 *Si una función es derivable en un punto a , entonces es continua en a .*

DEMOSTRACIÓN: Sea $f : A \rightarrow \mathbb{R}$ derivable en a . Entonces existe

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h}.$$

Como $\lim_{h \rightarrow 0} h = 0$, multiplicando obtenemos que

$$\lim_{h \rightarrow 0} (f(a+h) - f(a)) = f'(a) \cdot 0 = 0.$$

Por lo tanto $\lim_{h \rightarrow 0} f(a+h) = f(a)$. Teniendo en cuenta la definición de límite, es fácil ver que esto equivale a que $\lim_{x \rightarrow a} f(x) = f(a)$. Esto significa que f es continua en a . ■

En particular, las funciones derivables son continuas, pero no toda función continua es derivable.

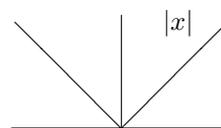
Ejemplos Ya hemos dicho que $\sqrt[3]{x}$ no es derivable en 0. Es fácil probarlo. La derivada en 0 sería

$$\lim_{h \rightarrow 0} \frac{\sqrt[3]{h} - 0}{h} = \lim_{h \rightarrow 0} \frac{1}{\sqrt[3]{h^2}} = +\infty.$$

Aquí la razón es que la pendiente de la función se vuelve infinita en 0.

Otra causa de no derivabilidad (a pesar de la continuidad) es que la función forme un “pico”. Por ejemplo $f(x) = |x|$ en $x = 0$. La derivada sería

$$\lim_{h \rightarrow 0} \frac{|h| - 0}{h} = \lim_{h \rightarrow 0} \operatorname{sig} h,$$



pero es claro que el límite por la izquierda es -1 y el límite por la derecha es $+1$, luego no existe tal límite. ■

4.2 Cálculo de derivadas

El teorema siguiente recoge las propiedades básicas que nos permiten derivar las funciones más simples:

Teorema 4.3 Sean $f, g : A \rightarrow \mathbb{R}$ funciones derivables en un punto $a \in A$ y sea $\alpha \in \mathbb{R}$.

- a) $f + g$ es derivable en a y $(f + g)'(a) = f'(a) + g'(a)$.
- b) αf es derivable en a y $(\alpha f)'(a) = \alpha f'(a)$.
- c) fg es derivable en a y $(fg)'(a) = f'(a)g(a) + f(a)g'(a)$.
- d) Si $g(a) \neq 0$, f/g es derivable en a y

$$(f/g)'(a) = \frac{f'(a)g(a) - f(a)g'(a)}{g^2(a)}.$$

En particular, las funciones derivables en A forman una subálgebra de $C(A)$.

DEMOSTRACIÓN: Las propiedades a) y b) son muy sencillas. Veamos c).

$$\begin{aligned}
 (fg)'(a) &= \lim_{h \rightarrow 0} \frac{g(a+h)g(a+h) - f(a)g(a)}{h} \\
 &= \lim_{h \rightarrow 0} \frac{f(a+h)g(a+h) - f(a)g(a+h) + f(a)g(a+h) - f(a)g(a)}{h} \\
 &= \lim_{h \rightarrow 0} \left(\frac{f(a+h) - f(a)}{h} g(a+h) + f(a) \frac{g(a+h) - g(a)}{h} \right) \\
 &= f'(a)g(a) + f(a)g'(a),
 \end{aligned}$$

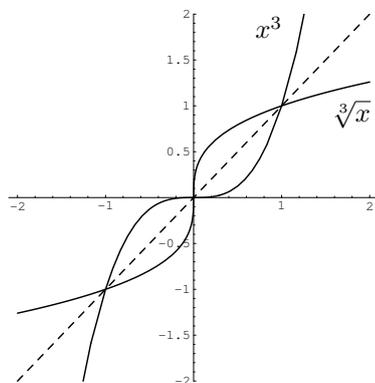
donde hemos usado la continuidad de g en a al afirmar que $\lim_{h \rightarrow 0} g(a+h) = g(a)$.

La prueba de d) es similar. ■

Aplicando inductivamente la propiedad c) se obtiene que $(x^n)' = nx^{n-1}$, para todo natural $n \geq 1$. Aplicando ahora d) resulta que esto es cierto para todo entero $n \neq 0$. En particular todos los polinomios y fracciones algebraicas son derivables en sus dominios.

Por ejemplo, la derivada de $3x^4 - 2x^3 + x^2 + 5x - 3$ es igual a $12x^3 - 6x^2 + 2x + 5$. Observemos que la derivada de un polinomio coincide con su derivada formal en el sentido algebraico.

La derivación de las raíces se seguirá un resultado general sobre funciones inversas. Por ejemplo, la función $\sqrt[3]{x}$ es la inversa de la función x^3 . Esto significa que un punto (x, y) está en la gráfica de $\sqrt[3]{x}$ si y sólo si el punto (y, x) está en la de x^3 . La gráfica de una se obtiene de la de la otra cambiando x por y . Geométricamente esto equivale a reflejar la gráfica respecto a la diagonal:



Es claro geoméricamente que si una función tiene tangente en el punto (x, y) , su inversa tendrá tangente en el punto (y, x) , y que la recta tangente a la inversa resultará de girar respecto a la diagonal la recta tangente a la función original. Ahora bien, ¿qué relación hay entre la pendiente de una recta y la pendiente de la recta que resulta de girarla respecto a la diagonal?. Si una recta es $y = mx + n$ (con pendiente m), girar respecto a la diagonal es cambiar y por x , o sea, pasar a $x = my + n$. Despejando y obtenemos $y = (1/m)x - n/m$.

Por lo tanto la pendiente es $1/m$. Si la recta original tiene pendiente 0 (es horizontal), la recta girada es vertical, tiene pendiente infinita.

Por lo tanto es de esperar que si una función biyectiva f cumple $f(a) = b$ y tiene derivada $m \neq 0$ en a , entonces su inversa tiene derivada $1/m$ en b . Antes de probarlo analíticamente damos un sencillo resultado técnico.

Definición 4.4 Sea A un intervalo y $f : A \rightarrow \mathbb{R}$, diremos que f es *creciente* en A si cuando $x < y$ son dos puntos de A , se cumple $f(x) \leq f(y)$. Si de hecho se cumple $f(x) < f(y)$ diremos que f es *estrictamente creciente* en A . Se dice que f es *decreciente* en A si cuando $x < y$ son puntos de A , se cumple $f(y) \leq f(x)$. Si se cumple $f(y) < f(x)$ se dice que es *estrictamente decreciente* en A . La función f es (estrictamente) monótona en A si es (estrictamente) creciente o decreciente en A .

Por ejemplo, la función x^3 es estrictamente creciente en \mathbb{R} .

Teorema 4.5 Sea A un intervalo y $f : A \rightarrow \mathbb{R}$ una función inyectiva y continua. Entonces f es estrictamente monótona en A .

DEMOSTRACIÓN: Sean $a < b$ dos puntos cualesquiera de A . Supongamos que $f(a) < f(b)$. Entonces todo $a < x < b$ ha de cumplir $f(a) < f(x) < f(b)$, pues si, por ejemplo, $f(x) < f(a) < f(b)$, por el teorema de los valores intermedios, en el intervalo $]x, b[$ habría un punto cuya imagen sería $f(a)$, y f no sería inyectiva.

De aquí se sigue que f es creciente en $[a, b]$, pues si $a \leq x < y \leq b$, hemos visto que $f(a) < f(x) < f(b)$, y aplicando lo mismo a los puntos x, b , resulta que $f(x) < f(y) < f(b)$. Igualmente, de $f(b) < f(a)$ llegaríamos a que f es decreciente en $[a, b]$. Por lo tanto f es monótona en cualquier intervalo $[a, b]$ contenido en A .

Pero si f no fuera monótona en A existirían puntos $u < v$ y $r < s$ tales que $f(u) < f(v)$ y $f(s) < f(r)$. Tomando el mínimo y el máximo de estos cuatro puntos obtendríamos los extremos de un intervalo en el que f no sería monótona. ■

Teorema 4.6 (Teorema de la función inversa) Sea A un intervalo abierto y $f : A \rightarrow \mathbb{R}$ una función inyectiva y derivable en A tal que f' no se anule en ningún punto de A . Entonces

- a) $B = f[A]$ es un intervalo abierto.
- b) La función inversa $g = f^{-1} : B \rightarrow A$ es derivable en B .
- c) Para todo $a \in A$, si $f(a) = b$, se cumple que $g'(b) = 1/f'(a)$.

DEMOSTRACIÓN: Por el teorema anterior sabemos que f es estrictamente monótona. Digamos que es monótona creciente. Si $a < b$ son dos puntos de A , por conexión $f[]a, b[$ ha de ser un intervalo, y de la monotonía se sigue fácilmente que $f[]a, b[=]f(a), f(b)[$.

Dado $a \in A$, podemos tomar un $\epsilon > 0$ tal que $[a - \epsilon, a + \epsilon] \subset A$, con lo que

$$f(a) \in]f(a - \epsilon), f(a + \epsilon)[= f[]a - \epsilon, a + \epsilon[\subset B.$$

Así pues B es un entorno de $f(a)$ para todo $a \in A$, o sea, para todos los puntos de B . Por lo tanto B es abierto. Por conexión es un intervalo. Además hemos visto que f envía abiertos básicos $]a, b[$ a abiertos básicos, y esto significa que g es continua.

Sea ahora $f(a) = b$. Por la monotonía, si $h \neq 0$, entonces $g(b + h) \neq g(b)$. Sea $k = g(b + h) - g(b) \neq 0$. Así $g(b + h) = k + a$, luego $b + h = f(k + a)$, y $h = f(k + a) - f(a)$. Por lo tanto

$$\frac{g(b + h) - g(b)}{h} = \frac{1}{\frac{f(a+k) - f(a)}{k}}$$

Ahora, k es una función de h y, como g es continua, $\lim_{h \rightarrow 0} k(h) = 0$. La derivabilidad de f en el punto a nos da que

$$g'(b) = \lim_{h \rightarrow 0} \frac{g(b + h) - g(b)}{h} = \frac{1}{f'(a)}.$$

■

Ejemplo Sea n un número natural no nulo y consideremos $f(x) = x^n$ definida en $]0, +\infty[$. Sabemos que es inyectiva y derivable. Su derivada es nx^{n-1} , que no se anula en $]0, +\infty[$. Por lo tanto su inversa, que es $g(x) = \sqrt[n]{x}$, es derivable en su dominio y si $y^n = x$ (con lo que $y = \sqrt[n]{x}$), entonces $g'(x) = 1/f'(y)$, o sea,

$$(\sqrt[n]{x})' = \frac{1}{ny^{n-1}} = \frac{1}{n\sqrt[n]{x}^{n-1}} = \frac{1}{n}x^{-1+1/n}.$$

Así pues, tenemos probado que la regla de derivación $x^r \mapsto rx^{r-1}$ es válida cuando r es entero o de la forma $1/n$, donde n es un número natural no nulo. Aplicando la regla del producto se concluye por inducción que vale de hecho para todo número racional r . ■

Con todo lo anterior todavía no sabemos derivar funciones como $\sqrt{x^2 + 1}$. El teorema siguiente nos permite de derivar cualquiera de las funciones que podemos construir a partir de polinomios y raíces.

Teorema 4.7 (regla de la cadena) Sean $f : A \rightarrow \mathbb{R}$ y $g : B \rightarrow \mathbb{R}$. Sea un punto $a \in A$ tal que f sea derivable en a y g sea derivable en $f(a)$. Entonces la función compuesta $f \circ g$ es derivable en a y $(f \circ g)'(a) = g'(f(a))f'(a)$.

DEMOSTRACIÓN: Notemos que B es un entorno de $f(a)$ y f es continua en a , luego $f^{-1}[B]$ es un entorno de a sobre el que está definida $f \circ g$. Sea $b = f(a)$. Para $k \neq 0$, llamemos

$$G(k) = \frac{g(b + k) - g(b)}{k} - g'(b).$$

La función G está definida para los puntos k tales que $b+k \in B$. Como B es abierto, G está definida al menos para h en un intervalo $]-\epsilon, \epsilon[\setminus \{0\}$. Como g es derivable en b , existe $\lim_{k \rightarrow 0} G(k) = 0$, luego si definimos $G(0) = 0$ tenemos que G es continua en un entorno de 0. Claramente además

$$g(b+k) - g(b) = (g'(b) + G(k))k.$$

Ahora tomamos $h \neq 0$ tal que $a+h \in A$ y $k = f(a+h) - f(a)$, con lo que se cumple $f(a+h) = b+k$, luego $b+k \in B$ y está definido $G(k)$. Entonces

$$\begin{aligned} g(f(a+h)) - g(f(a)) &= (g'(f(a)) + G(k))k \\ &= (g'(f(a)) + G(k))(f(a+h) - f(a)). \end{aligned}$$

En consecuencia

$$\frac{(f \circ g)(a+h) - (f \circ g)(a)}{h} = (g'(f(a)) + G(f(a+h) - f(a))) \frac{f(a+h) - f(a)}{h}.$$

Usando la continuidad de f en a y la de G en 0, tomamos el límite cuando h tiende a 0 y queda que existe $(f \circ g)'(a) = g'(f(a))f'(a)$. ■

Ejemplo La función $h(x) = \sqrt{x^2 + 1}$ es derivable en \mathbb{R} , pues es la composición del polinomio $f(x) = x^2 + 1$ con la función $g(x) = \sqrt{x}$, y ambas funciones son derivables en sus dominios. Sabemos que $f'(x) = 2x$ y $g'(x) = 1/(2\sqrt{x})$. La regla de la cadena nos da que

$$h'(x) = g'(x^2 + 1)f'(x) = \frac{2x}{2\sqrt{x^2 + 1}} = \frac{x}{\sqrt{x^2 + 1}}. \quad \blacksquare$$

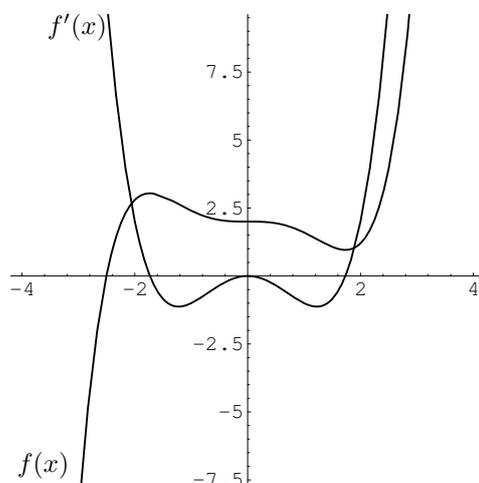
El lector que no esté familiarizado con el cálculo de derivadas debería practicar hasta que la derivación le resultara un acto mecánico. Hay muchos libros adecuados para ello, por lo que en adelante dejaremos de justificar los cálculos de derivadas.

4.3 Propiedades de las funciones derivables

La derivada de una función contiene mucha información sobre ésta. Consideremos por ejemplo

$$f(x) = \frac{x^5}{10} - \frac{x^3}{2} + 2, \quad f'(x) = \frac{x^4}{2} - \frac{3}{2}x^2.$$

La figura de la página siguiente muestra las gráficas. Fijémonos en el signo de la derivada. Es fácil ver que f' tiene una raíz doble en 0 y dos raíces simples en $\pm\sqrt{3}$. La restricción de f' a cada uno de los intervalos $]-\infty, -\sqrt{3}[$, $]-\sqrt{3}, 0[$, $]0, \sqrt{3}[$ y $]\sqrt{3}, +\infty[$ es una función continua que no se anula, luego por el teorema de los valores intermedios f' tiene signo constante en cada uno de ellos. Es fácil ver entonces que el signo de f' varía como indica la gráfica, es decir, f' es positiva en los dos intervalos no acotados y negativa en los acotados.



En los puntos donde f' es positiva la tangente a f tiene pendiente positiva, y vemos en la gráfica que esto se traduce en que f es creciente. Por el contrario, en los intervalos donde f' es negativa la función f es decreciente.

En los puntos donde f' se anula la tangente a f es horizontal. En $-\sqrt{3}$, la derivada pasa de ser positiva a ser negativa, luego f pasa de ser creciente a ser decreciente, y por ello el punto es un máximo relativo, en el sentido de que f toma en $-\sqrt{3}$ un valor mayor que en los puntos de alrededor. En cambio, en $\sqrt{3}$ la derivada pasa de negativa a positiva, f pasa de decreciente a creciente y el punto es un mínimo relativo. El caso del 0 es distinto, pues f' es negativa a la izquierda, toca el 0 y vuelve a bajar, con lo que sigue siendo negativa. Por ello f es creciente en 0 y no tiene ni un máximo ni un mínimo en 0.

Vemos así que conociendo la derivada podemos formarnos una idea de la función: dónde crece, dónde decrece, dónde tiene máximos y mínimos, y más cosas de las que no hemos hablado. Vamos a desarrollar estas ideas.

Definición 4.8 Sea $f : A \subset \mathbb{R} \rightarrow \mathbb{R}$. Diremos que f tiene un *máximo relativo* en un punto $a \in A$ si existe un entorno V de a contenido en A de modo que para todo $x \in V$ se cumple $f(x) \leq f(a)$. Diremos que f tiene un *mínimo relativo* en a si existe un entorno V de a contenido en A tal que para todo $x \in V$ se cumple $f(a) \leq f(x)$. La función f tiene un *extremo relativo* en a si tiene un máximo o un mínimo relativo en a .

Teorema 4.9 Si $f : A \rightarrow \mathbb{R}$ es una función derivable en un punto $a \in A$ y f tiene un extremo relativo en a , entonces $f'(a) = 0$.

DEMOSTRACIÓN: Supongamos, por reducción al absurdo, que $f'(a) > 0$. (El caso $f'(a) < 0$ se razona análogamente.) Entonces $]0, +\infty[$ es un entorno de $f'(a)$, luego por definición de límite y de derivada existe un $\epsilon > 0$ de manera que $]a - \epsilon, a + \epsilon[\subset A$ y si $0 < |h| < \epsilon$ entonces

$$\frac{f(a+h) - f(a)}{h} > 0.$$

Esto se traduce en que $f(a+h) > f(a)$ si $h > 0$ y $f(a+h) < f(a)$ si $h < 0$, lo que contradice que f tenga un extremo relativo en a . ■

La función del ejemplo anterior muestra que $f'(a) = 0$ no implica que a sea un extremo relativo. Más adelante volveremos sobre este punto. Ahora probemos una consecuencia sencilla de este teorema:

Teorema 4.10 (Teorema de Rolle) *Sea $f : [a, b] \rightarrow \mathbb{R}$ una función continua en $[a, b]$ y derivable en $]a, b[$. Si $f(a) = f(b)$, entonces existe un $c \in]a, b[$ tal que $f'(c) = 0$.*

DEMOSTRACIÓN: Como $[a, b]$ es compacto, la función f alcanza un valor mínimo m y un valor máximo M . Si se cumpliera que $m = M = f(a) = f(b)$, entonces f sería constante y su derivada sería nula, luego cualquier $c \in]a, b[$ cumpliría el teorema.

Supongamos que $m < M$. Entonces, o bien $m \neq f(a)$ o bien $M \neq f(a)$. Digamos por ejemplo $M \neq f(a)$. Sea $c \in [a, b]$ el punto donde $f(c) = M$. Como $M \neq f(a) = f(b)$, ha de ser $a < b < c$, y como f toma en c su valor máximo, en particular c es un máximo relativo de f , luego por el teorema anterior $f'(c) = 0$. ■

El interés principal del teorema de Rolle es que permite probar el teorema del valor medio, como veremos a continuación, pero de hecho conviene probar un resultado técnico algo más general:

Teorema 4.11 (Teorema de Cauchy) *Sean $f, g : [a, b] \rightarrow \mathbb{R}$ funciones continuas en $[a, b]$ y derivables en $]a, b[$. Entonces existe un $c \in]a, b[$ tal que*

$$g'(c)(f(b) - f(a)) = f'(c)(g(b) - g(a)).$$

DEMOSTRACIÓN: Consideremos la función dada por

$$h(x) = f(x)(g(b) - g(a)) - g(x)(f(b) - f(a)).$$

Se cumple que $h(a) = h(b) = f(a)g(b) - g(a)f(b)$. Además h es continua en $[a, b]$ y derivable en $]a, b[$. Por el teorema de Rolle existe un punto $c \in]a, b[$ tal que $h'(c) = 0$, pero $h'(x) = f'(x)(g(b) - g(a)) - g'(x)(f(b) - f(a))$, luego

$$f'(c)(g(b) - g(a)) - g'(c)(f(b) - f(a)) = 0 \quad \blacksquare$$

Más adelante tendremos ocasión de usar este resultado en toda su generalidad, pero de momento nos basta con el caso particular que resulta de tomar como función g la dada por $g(x) = x$. Entonces tenemos:

Teorema 4.12 (Teorema del valor medio) *Sea $f : [a, b] \rightarrow \mathbb{R}$ una función continua en $[a, b]$ y derivable en $]a, b[$. Entonces existe un $c \in]a, b[$ tal que*

$$f(b) - f(a) = f'(c)(b - a).$$

Notar que el teorema de Rolle es un caso particular del teorema del valor medio. Este teorema tiene una interpretación geométrica. La expresión

$$\frac{f(b) - f(a)}{b - a}$$

puede interpretarse como la “pendiente media” de f en $[a, b]$, es decir, es el cociente de lo que aumenta f cuando la variable x pasa de a a b dividido entre lo que ha aumentado la variable. Lo que dice el teorema del valor medio es que hay un punto en el intervalo donde la función toma el valor medio de su pendiente.

La importancia de este teorema es que nos relaciona una magnitud global, la pendiente media, con una magnitud local, la derivada en un punto. Las consecuencias son muchas. Por ejemplo, ahora podemos probar algo que ya habíamos observado al principio de esta sección:

Teorema 4.13 *Si $f : [a, b] \rightarrow \mathbb{R}$ es continua en $[a, b]$, derivable en $]a, b[$ y su derivada es ≥ 0 , (resp. > 0 , ≤ 0 , < 0) en $]a, b[$, entonces f es creciente (resp. estrictamente creciente, decreciente, estrictamente decreciente) en $[a, b]$.*

DEMOSTRACIÓN: Tomemos $a \leq x < y \leq b$. Por el teorema del valor medio, existe un $x < c < y$ de modo que

$$f(y) - f(x) = f'(c)(y - x),$$

y por la hipótesis sobre la derivada concluimos que $f(y) - f(x) \geq 0$ (resp. > 0 , ≤ 0 , < 0), lo que nos da la conclusión. ■

Sabemos que las funciones constantes tienen derivada nula. El teorema del valor medio nos da el recíproco:

Teorema 4.14 *Si una función tiene derivada nula en todos los puntos de un intervalo abierto, entonces es constante.*

DEMOSTRACIÓN: Sea f una función derivable en un intervalo A con derivada nula. Sean $a < b$ dos puntos cualesquiera de A . Entonces

$$f(b) - f(a) = f'(c)(b - a) = 0,$$

donde c es un punto de $]a, b[$. Por lo tanto f es constante. ■

Una consecuencia inmediata es el teorema siguiente, que afirma que una función derivable está unívocamente determinada por su derivada y su valor en un punto cualquiera.

Teorema 4.15 *Si f y g son funciones derivables en un intervalo abierto y $f' = g'$, entonces existe un $k \in \mathbb{R}$ tal que $f = g + k$.*

DEMOSTRACIÓN: La función $f - g$ tiene derivada nula, luego $f - g = k$. ■

La derivada de una función también está relacionada con su concavidad y convexidad en el sentido siguiente:

Definición 4.16 Diremos que una función $f :]a, b[\rightarrow \mathbb{R}$ es convexa (resp. cóncava) si para todo $a < x < y < b$ y todo $0 < \lambda < 1$ se cumple que

$$f((1 - \lambda)x + \lambda y) \leq (1 - \lambda)f(x) + \lambda f(y)$$

o

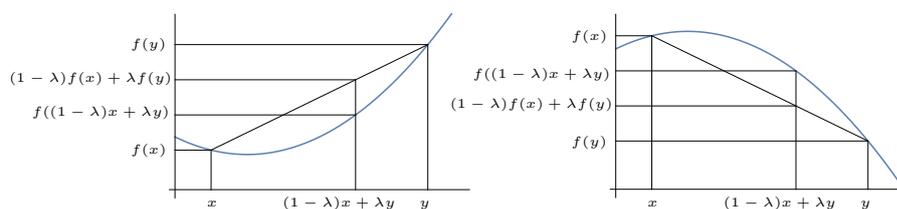
$$f((1 - \lambda)x + \lambda y) \geq (1 - \lambda)f(x) + \lambda f(y),$$

respectivamente.

Estas definiciones tienen una interpretación geométrica muy simple. Basta tener en cuenta que, cuando λ varía entre 0 y 1, la expresión $z = (1 - \lambda)x + \lambda y$ recorre todos los puntos intermedios entre x e y , e igualmente, $(1 - \lambda)f(x) + \lambda f(y)$ recorre todos los puntos intermedios entre $f(x)$ y $f(y)$. Más precisamente, la expresión

$$((1 - \lambda)x + \lambda y, (1 - \lambda)f(x) + \lambda f(y)) = (1 - \lambda)(x, f(x)) + \lambda(y, f(y))$$

recorre el segmento de extremos $(x, f(x))$ e $(y, f(y))$. Así, la figura de la izquierda ilustra la definición de función convexa, que expresa que el segmento de extremos $(x, f(x))$ e $(y, f(y))$ está por encima de la gráfica de f , mientras que en el caso de una función cóncava (la figura de la derecha) el segmento está por debajo de la gráfica.



Es inmediato comprobar que una función f es cóncava si y sólo si $-f$ es convexa, y esto permite traducir inmediatamente a funciones cóncavas todos los resultados que vamos a demostrar para funciones convexas.

Si en $z = (1 - \lambda)x + \lambda y$ despejamos λ e incluimos la expresión en la definición de convexidad, obtenemos la versión equivalente: si $a < x < z < y < b$, entonces

$$f(z) \leq \frac{y - z}{y - x}f(x) + \frac{z - x}{y - x}f(y).$$

Operando:

$$yf(x) - zf(x) + zf(y) - xf(y) - yf(z) + xf(z) \geq 0.$$

Una comprobación rutinaria muestra ahora que esto equivale a cualquiera de las tres desigualdades siguientes:

$$\frac{f(z) - f(x)}{z - x} \leq \frac{f(y) - f(x)}{y - x} \leq \frac{f(y) - f(z)}{y - z}$$

(sólo hay que desarrollarlas y comprobar que llevan a la misma desigualdad precedente). Si llamamos

$$p(x, y) = \frac{f(y) - f(x)}{y - x},$$

podemos expresar así las desigualdades precedentes:

$$p(x, z) \leq p(x, y) \leq p(z, y),$$

es decir, que la pendiente del segmento que une dos puntos de la gráfica de f aumenta cuando se desplaza hacia la derecha cualquiera de ellos.

Veamos la relación entre la convexidad y las derivadas (el lector puede constatarla en el ejemplo dado al principio de esta sección):

Teorema 4.17 *Si $f :]a, b[\rightarrow \mathbb{R}$ es una función derivable y su derivada f' es creciente (decreciente), entonces f es convexa (cóncava).*

DEMOSTRACIÓN: Si $a < x < z < y < b$, por el teorema del valor medio existen $x < c < z < c' < y$ tales que

$$\frac{f(z) - f(x)}{z - x} = f'(c) \leq f'(c') = \frac{f(y) - f(z)}{y - z},$$

que es una de las tres caracterizaciones de la convexidad en términos de pendientes. ■

En particular, si f' es también derivable y su derivada f'' es ≥ 0 (resp. ≤ 0) en $]a, b[$, entonces f es convexa (resp. cóncava).

Otro hecho elemental es que una función convexa (cóncava) está siempre por encima (por debajo) de sus rectas tangentes:

Teorema 4.18 *Si $f :]a, b[\rightarrow \mathbb{R}$ es una función convexa derivable en un punto x , entonces, para todo $y \in]a, b[$ se cumple que*

$$f(x) + f'(x)(y - x) \leq f(y).$$

DEMOSTRACIÓN: Supongamos que $x < y$ y sea $x < z < y$. Entonces sabemos que

$$\frac{f(z) - f(x)}{z - x} \leq \frac{f(y) - f(x)}{y - x}.$$

Tomando el límite cuando $z \rightarrow x$ resulta que

$$f'(x) \leq \frac{f(y) - f(x)}{y - x}.$$

Si es $y < x$ usamos que

$$\frac{f(y) - f(x)}{y - x} = \frac{f(x) - f(y)}{x - y} \leq \frac{f(x) - f(z)}{x - z} = \frac{f(z) - f(x)}{z - x} \rightarrow f'(x)$$

y que $y - x < 0$, luego al despejarlo se invierte la desigualdad. ■

Una función cóncava o convexa no es necesariamente derivable (por ejemplo, la función $|x|$ es convexa), pero conviene observar que sí que es necesariamente continua:

Teorema 4.19 *Toda función cóncava o convexa $f :]a, b[\rightarrow \mathbb{R}$ es continua.*

DEMOSTRACIÓN: Veamos que si $a < a' < b' < b$, para cada par de números $a' \leq x < y \leq b'$ existe un $C > 0$ tal que $|p(x, y)| \leq C$. En efecto, tomamos $a < a'' < a' < b' < b'' < b$ y observamos que

$$p(a', a'') \leq p(a', y) \leq p(x, y) \leq p(x, b') \leq p(b'', b'),$$

y basta tomar $C = \max\{|p(a', a'')|, |p(b'', b')|\}$, pues entonces $-C \leq p(x, y) \leq C$. Así pues, $|f(y) - f(x)| \leq C|y - x|$, lo cual significa que f tiene la propiedad de Lipschitz en $[a', b']$, luego es continua en dicho intervalo, luego también en $]a, b[$. ■

Las derivadas proporcionan un teorema muy útil para el cálculo de límites. De momento no podemos estimar su valor porque los límites de las funciones que conocemos (polinomios, fracciones algebraicas. etc.) son fáciles de calcular directamente, pero más adelante tendremos ocasión de aprovecharlo.

Teorema 4.20 (Regla de L'Hôpital) *Sean $f, g :]a, b[\rightarrow \mathbb{R}$ funciones derivables tales que $\lim_{x \rightarrow a} f(x) = \lim_{x \rightarrow a} g(x) = 0$ y de modo que g y g' no se anulen en $]a, b[$. Si existe*

$$\lim_{x \rightarrow a} \frac{f'(x)}{g'(x)} = L,$$

entonces también existe

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = L.$$

DEMOSTRACIÓN: Extendamos f y g al intervalo $[a, b[$ estableciendo que $f(a) = g(a) = 0$. Así siguen siendo continuas.

Si $a < x < b$, por el teorema de Cauchy existe un punto $c \in]a, x[$ tal que

$$(f(x) - f(a))g'(c) = (g(x) - g(a))f'(c),$$

o sea, $f(x)g'(c) = g(x)f'(c)$, y como $g(x) \neq 0 \neq g'(c)$, podemos escribir

$$\frac{f(x)}{g(x)} = \frac{f'(c)}{g'(c)}. \quad (4.2)$$

Por definición de límite, si $\epsilon > 0$, existe un $\delta > 0$ tal que si $0 < c - a < \delta$, entonces

$$\left| \frac{f'(c)}{g'(c)} - L \right| < \epsilon. \quad (4.3)$$

Así tenemos que si $0 < x - a < \delta$, existe un $c \in]a, x[$ que cumple (4.2) y (4.3). Por consiguiente, para todo $x \in]a, a + \delta[$ se cumple

$$\left| \frac{f(x)}{g(x)} - L \right| < \epsilon.$$

Esto significa que

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = L. \quad \blacksquare$$

Obviamente la regla de L'Hôpital también es válida cuando en las hipótesis cambiamos a por b . Combinando las dos versiones obtenemos la regla de L'Hôpital para funciones definidas en intervalos $]a - \epsilon, a + \epsilon[\setminus \{a\}$ y tomando límites en a (si existe el límite del cociente de derivadas, existen los límites por la derecha y por la izquierda y coinciden, por los casos correspondientes de la regla, existen los límites de los cocientes de las funciones por ambos lados y coinciden, luego existe el límite y coincide con el de las derivadas).

Los teoremas siguientes demuestran otras variantes de la regla de L'Hôpital de no menor interés:

Teorema 4.21 (Regla de L'Hôpital) Sean $f, g :]a, +\infty[\rightarrow \mathbb{R}$ dos funciones derivables tales que $\lim_{x \rightarrow +\infty} f(x) = \lim_{x \rightarrow +\infty} g(x) = 0$ y de modo que g y g' no se anulan en $]a, +\infty[$. Si existe

$$\lim_{x \rightarrow +\infty} \frac{f'(x)}{g'(x)} = L,$$

entonces también existe

$$\lim_{x \rightarrow +\infty} \frac{f(x)}{g(x)} = L.$$

DEMOSTRACIÓN: Consideremos $F(x) = f(1/x)$ y $G(x) = g(1/x)$, definidas en $]0, 1/a[$. Claramente F y G son continuas, y $\lim_{x \rightarrow 0} F(x) = \lim_{x \rightarrow 0} G(x) = 0$. Además por la regla de la cadena son funciones derivables y sus derivadas son

$$F'(x) = -\frac{f'(1/x)}{x^2}, \quad G'(x) = -\frac{g'(1/x)}{x^2}.$$

También es claro que ni G ni G' se anulan en su dominio y

$$\frac{F'(x)}{G'(x)} = \frac{f'(1/x)}{g'(1/x)},$$

luego existe

$$\lim_{x \rightarrow 0} \frac{F'(x)}{G'(x)} = L.$$

El caso ya probado de la regla de L'Hôpital nos da ahora que también existe

$$\lim_{x \rightarrow 0} \frac{F(x)}{G(x)} = L.$$

Obviamente entonces

$$\lim_{x \rightarrow +\infty} \frac{f(x)}{g(x)} = L. \quad \blacksquare$$

Igualmente se prueba la regla de L'Hôpital para funciones definidas en intervalos $]-\infty, b[$ y cuando x tiende a $-\infty$.

Así pues, si tenemos una indeterminación de tipo $0/0$ y al derivar numerador y denominador podemos calcular el límite, la función original tiene ese mismo límite. Ahora veremos que la regla de L'Hôpital es aplicable también a indeterminaciones del tipo ∞/∞ .

Teorema 4.22 (Regla de L'Hôpital) Sean $f, g :]a, +\infty[\rightarrow \mathbb{R}$ dos funciones derivables tales que $\lim_{x \rightarrow +\infty} f(x) = \lim_{x \rightarrow +\infty} g(x) = \infty$ y de modo que g y g' no se anulan en $]a, +\infty[$. Si existe

$$\lim_{x \rightarrow +\infty} \frac{f'(x)}{g'(x)} = L,$$

entonces también existe

$$\lim_{x \rightarrow +\infty} \frac{f(x)}{g(x)} = L.$$

DEMOSTRACIÓN: Por definición de límite, dado $\epsilon > 0$, existe un $M > a$ tal que si $x > M$ entonces

$$\left| \frac{f'(x)}{g'(x)} - L \right| < \epsilon.$$

Por el teorema de Cauchy, si $x > M$, existe un $y \in]M, x[$ de modo que

$$\frac{f(x) - f(M)}{g(x) - g(M)} = \frac{f'(y)}{g'(y)},$$

luego

$$\left| \frac{f(x) - f(M)}{g(x) - g(M)} - L \right| < \epsilon.$$

(Notar que, como g' no se anula, la función g es monótona, luego el denominador es no nulo).

Puesto que $\lim_{x \rightarrow +\infty} f(x) = \infty$, existe un $N > M$ tal que si $x > N$ entonces $|f(x)| > |f(M)|$, y en particular $f(x) - f(M) \neq 0$. Por ello, para $x > N$ podemos escribir

$$\frac{f(x)}{g(x)} = \frac{f(x) - f(M)}{g(x) - g(M)} \frac{f(x)}{f(x) - f(M)} \frac{g(x) - g(M)}{g(x)}.$$

Si en los dos últimos factores dividimos numerador y denominador entre $f(x)$ y $g(x)$ respectivamente, queda claro que tienden a 1 cuando $x \rightarrow +\infty$, luego tomando N suficientemente grande podemos suponer que si $x > N$ entonces el producto de ambos dista de 1 menos de ϵ . De este modo, para x suficientemente grande, el cociente $f(x)/g(x)$ se puede expresar como producto de dos números reales, uno arbitrariamente próximo a L y otro arbitrariamente próximo a 1. De la continuidad del producto se sigue que

$$\lim_{x \rightarrow +\infty} \frac{f(x)}{g(x)} = L.$$

■

Naturalmente la regla de L'Hôpital también es válida en el caso ∞/∞ cuando $x \rightarrow -\infty$. El mismo argumento que nos ha permitido pasar del caso finito al caso infinito en la indeterminación $0/0$ nos permite pasar ahora al caso finito. Es fácil probar:

Teorema 4.23 (Regla de L'Hôpital) Sean $f, g :]a, b[\rightarrow \mathbb{R}$ derivables tales que $\lim_{x \rightarrow a} f(x) = \lim_{x \rightarrow a} g(x) = \infty$ y de modo que g y g' no se anulan en $]a, b[$. Si existe

$$\lim_{x \rightarrow a} \frac{f'(x)}{g'(x)} = L,$$

entonces también existe

$$\lim_{x \rightarrow a} \frac{f(x)}{g(x)} = L.$$

También se cumple la versión correspondiente cuando x tiende a b y cuando x tiende a un punto por la izquierda y la derecha a la vez.

4.4 La diferencial de una función

Consideremos una función $f : A \rightarrow \mathbb{R}$ derivable en un punto $a \in A$. Sea Δx un número próximo¹ a 0 de modo que $a + \Delta x \in A$. Al calcular f en el punto $a + \Delta x$ obtenemos una variación o incremento de f dado por $\Delta_a(f) = f(a + \Delta x) - f(a)$. La expresión $\Delta_a(f)$ representa a una función de la variable Δx , definida en un entorno de 0. La derivada de f en a es, por definición,

$$f'(a) = \lim_{\Delta x \rightarrow 0} \frac{\Delta_a(f)}{\Delta x}.$$

Llamemos

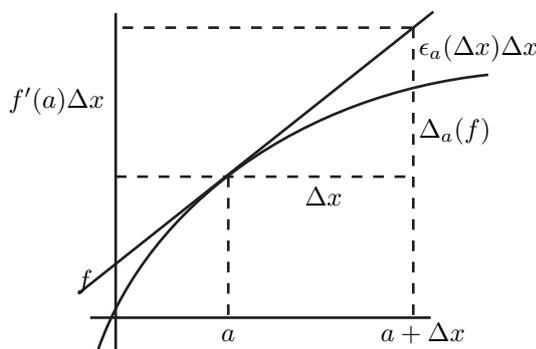
$$\epsilon_a(\Delta x) = \frac{\Delta_a(f)}{\Delta x} - f'(a).$$

Así, $\lim_{\Delta x \rightarrow 0} \epsilon_a(\Delta x) = 0$. Desarrollando las definiciones tenemos que

$$\Delta_a(f) = f(a + \Delta x) - f(a) = f'(a)\Delta x + \epsilon(\Delta x)\Delta x.$$

La figura muestra la situación:

¹El término Δx se lee "incremento de x ", porque representa un pequeño aumento de la variable x en el punto a .



El hecho de que $\epsilon_a(\Delta x)$ tienda a 0 expresa simplemente el hecho de que, para valores pequeños de Δx , se cumple $f(a + \Delta x) - f(a) \approx f'(a)\Delta x$, donde el signo \approx significa “aproximadamente igual”, es decir, que el valor de la función $f(a + \Delta x)$ es similar al de la recta tangente $f(a) + f'(a)\Delta x$. En la figura ambos valores son muy diferentes porque hemos tomado un Δx grande para mayor claridad.

Llamaremos *diferencial* de f en el punto a a la aplicación $df(a) : \mathbb{R} \rightarrow \mathbb{R}$ dada por $df(a)(\Delta x) = f'(a)\Delta x$. Así $df(a)$ es una aplicación lineal en \mathbb{R} de modo que $f(a + \Delta x) - f(a) \approx df(a)(\Delta x)$, o sea, la función $df(a)$ aproxima las diferencias entre las imágenes de f en puntos cercanos al punto a y la imagen de a . De aquí su nombre.

Si consideramos la función polinómica x , su derivada es 1, luego la diferencial de x es simplemente $dx(a)(\Delta x) = \Delta x$. Por ello podemos escribir $df(a)(\Delta x) = f'(a) dx(a)(\Delta x)$, luego tenemos la igualdad funcional

$$df(a) = f'(a) dx(a).$$

Si la función f es derivable en todo punto de A , la igualdad anterior se cumple en todo punto, luego si consideramos a df y dx como aplicaciones de A en el espacio de aplicaciones lineales de \mathbb{R} en \mathbb{R} , tenemos la igualdad funcional

$$df = f' dx,$$

donde dx es la función constante que a cada $a \in A$ le asigna la función identidad en \mathbb{R} .

Del mismo modo que la estructura topológica permite hablar de los puntos de alrededor de un punto dado, pese a que ningún punto en particular está alrededor de otro, así mismo la diferencial de una función recoge el concepto de “incremento infinitesimal” de una función, pese a que ningún incremento en particular es infinitamente pequeño. Por ejemplo, la igualdad $dx^2 = 2x dx$ expresa que cuando la variable x experimenta un incremento infinitesimal dx , la función x^2 experimenta un incremento infinitesimal de $2x dx$. Con rigor, dx^2 no es un incremento infinitesimal, sino la función que a cada incremento Δx le asigna una aproximación al incremento correspondiente de x^2 , de modo que

lo que propiamente tenemos es la aproximación que resulta de evaluar dx en incrementos concretos, es decir, $\Delta_x(x^2) \approx 2x\Delta x$. El error de esta aproximación se puede hacer arbitrariamente pequeño tomando Δx suficientemente pequeño.

Por ejemplo, $(1.1)^2 \approx 1^2 + dx^2(1)(0.1) = 1 + 2 \cdot 0.1 = 1.2$. En realidad $(1.1)^2 = 1.21$, luego el error cometido es de una centésima.

Dada la igualdad $df = f' dx$, representaremos también la derivada de f mediante la notación

$$f'(x) = \frac{df}{dx},$$

que expresa que $f'(x)$ es la proporción entre las funciones df y dx , o también que $f'(x)$ es la razón entre un incremento infinitesimal de f respecto al incremento infinitesimal de x que lo ocasiona.

Es costumbre, especialmente en física, nombrar las funciones, no por la expresión que las determina, sino por la magnitud que determinan. Por ejemplo, supongamos que la posición e de un objeto depende del tiempo viene dada por la relación $e(t) = t^2$. Entonces la velocidad del móvil es

$$v(t) = \frac{de}{dt} = 2t,$$

lo que nos permite expresar t en función de v , mediante $t(v) = v/2$. A su vez, esto nos permite calcular la posición en función de la velocidad, mediante la función $e(v) = v^2/4$.

De este modo, llamamos e tanto a la función $e(t)$ como a la función $e(v)$, que son funciones distintas. La letra v representa a una función en $v = 2t$ y a una variable en $t = v/2$. Estos convenios no provocan ninguna ambigüedad, al contrario, en muchos casos resultan más claros y permiten expresar los resultados de forma más elegante. Por ejemplo, si tenemos dos funciones $y = y(x)$ y $z = z(y)$, entonces la función compuesta se expresa, en estos términos, como $z = z(x)$. Si las funciones son derivables en sus dominios, la regla de la cadena se convierte en

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}.$$

La primera derivada es la de la función compuesta $z(x)$, mientras que la segunda es la de $z(y)$. No es necesario indicar que dicha derivada ha de calcularse en $y(x)$, pues esto ya está implícito en el hecho de que se trata de una función de y (no de x). Por ejemplo,

$$\frac{de}{dt} = \frac{de}{dv} \frac{dv}{dt} = \frac{v}{2} \cdot 2 = v = 2t,$$

como cabía esperar.

Similarmente, si $y(x)$ es una función inyectiva y derivable con derivada no nula, su inversa se representa por $x(y)$, y el teorema de la función inversa se expresa así:

$$\frac{dy}{dx} = \frac{1}{\frac{dx}{dy}}.$$

Por ejemplo,

$$\frac{dv}{dt} = \frac{1}{\frac{dt}{dv}} = \frac{1}{1/2} = 2.$$

Vemos, pues, que en estos términos las propiedades de las derivadas son formalmente análogas a las de las fracciones.

4.5 El teorema de Taylor

Sea $f : A \rightarrow \mathbb{R}$ una función derivable en el abierto A y tal que $f' : A \rightarrow \mathbb{R}$ también sea derivable en A . Entonces a la derivada de f' se la denomina *derivada segunda* de f en A , y se representa por f'' .

A su vez la derivada segunda puede ser derivable, y entonces está definida la derivada tercera, y así sucesivamente. Si una función admite n derivadas en A , a la derivada n -sima se la representa por $f^{(n)} : A \rightarrow \mathbb{R}$. Conviene usar la notación $f^{(0)}$ para referirse a la propia función f .

La derivada n -sima de un producto se puede calcular con una fórmula análoga a la del binomio de Newton:

Teorema 4.24 (Leibniz) *Si f y g son funciones n veces derivables en un punto x , entonces*

$$(fg)^{(n)}(x) = \sum_{i=0}^n \binom{n}{i} f^{(i)}(x) g^{(n-i)}(x).$$

DEMOSTRACIÓN: Por inducción sobre n . Para $n = 0$ es inmediato. Si vale para n , entonces

$$\begin{aligned} (fg)^{(n+1)}(x) &= \sum_{i=0}^n \binom{n}{i} (f^{(i+1)} g^{(n-i)}(x) + f^{(i)} g^{(n+1-i)}(x)) = \\ &f^{(n+1)}(x)g(x) + f(x)g^{(n+1)}(x) + \sum_{i=1}^n \binom{n}{i-1} f^{(i)}(x)g^{(n+1-i)}(x) + \sum_{i=1}^n \binom{n}{i} f^{(i)}(x)g^{(n+1-i)}(x) \\ &= \sum_{i=0}^{n+1} \binom{n+1}{i} f^{(i)}(x)g^{(n+1-i)}(x). \quad \blacksquare \end{aligned}$$

Llamaremos $C^n(A)$ al conjunto de las funciones definidas en A que admiten n derivadas y todas ellas son continuas en A . Si llamamos $C^0(A) = C(A)$, es decir, al conjunto de las funciones continuas en A , entonces tenemos

$$C^0(A) \supset C^1(A) \supset C^2(A) \supset C^3(A) \supset C^4(A) \supset \dots$$

Llamaremos $C^\infty(A)$ al conjunto de las funciones infinitamente derivables en A . Por ejemplo, los polinomios y las fracciones algebraicas son de clase C^∞ en su dominio. Es inmediato que todos estos conjuntos son subálgebras de $C(A)$.

Las inclusiones son todas estrictas. Por ejemplo, si $a \in A$ es fácil ver que la función dada por

$$f(x) = \begin{cases} (x-a)^{n+1} & \text{si } x \geq a \\ -(x-a)^{n+1} & \text{si } x \leq a \end{cases}$$

es una función de clase $C^n(A)$ pero no de clase $C^{n+1}(A)$.

Según sabemos, si una función f es derivable en un punto a , entonces alrededor de a la función f puede ser aproximada por su recta tangente, esto es, por el polinomio $f(a) + f'(a)(x - a)$. La recta tangente es el único polinomio $P(x)$ de grado 1 que cumple $P(a) = f(a)$ y $P'(a) = f'(a)$.

Cabe suponer que si una función f admite dos derivadas y tomamos un polinomio P de grado 2 tal que $P(a) = f(a)$, $P'(a) = f'(a)$ y $P''(a) = f''(a)$, el polinomio P nos dará una aproximación mejor de la función f que la recta tangente. Esto no siempre es así, pero hay bastante de verdad en ello. Vamos a investigarlo.

Ante todo, si K es un cuerpo y $a \in K$, la aplicación $u : K[x] \rightarrow K[x]$ dada por $u(p) = p(x - a)$ es un isomorfismo de K -espacios vectoriales. Como los polinomios $1, x, x^2, x^3, \dots$ son una K -base de $K[x]$, resulta que los polinomios

$$1, (x - a), (x - a)^2, (x - a)^3, (x - a)^4, \dots$$

son también una K -base, luego todo polinomio de $K[x]$ se expresa de forma única como

$$P(x) = c_0 + c_1(x - a) + c_2(x - a)^2 + \dots + c_n(x - a)^n, \quad (4.4)$$

para cierto natural n y ciertos coeficientes $c_0, \dots, c_n \in K$.

Si una función f admite n derivadas en un punto a , las ecuaciones

$$f(a) = P(a), \quad f'(a) = P'(a), \quad \dots \quad f^{(n)}(a) = P^{(n)}(a)$$

son satisfechas por un único polinomio de grado $\leq n$. En efecto, si $P(x)$ viene dado por (4.4), entonces $P(a) = c_0$, luego ha de ser $c_0 = f(a)$. Derivando obtenemos

$$P'(x) = c_1 + 2c_2(x - a) + \dots + nc_n(x - a)^{n-1},$$

de donde $P'(a) = c_1$, y ha de ser $c_1 = f'(a)$. Similarmente, $P''(a) = 2c_2$, luego $c_2 = f''(a)/2$. Igualmente se obtiene $c_3 = f'''(a)/6$ y, en general, $c_k = f^{(k)}(a)/k!$. En resumen:

$$P(x) = \sum_{k=0}^n \frac{f^{(k)}(a)}{k!} (x - a)^k.$$

Recíprocamente, es fácil ver que el polinomio $P(x)$ así definido cumple que $P^{(k)}(a) = f^{(k)}(a)$ para $k = 0, \dots, n$.

Definición 4.25 Sea f una función derivable n veces en un punto a . Llamaremos *polinomio de Taylor* de grado n de f en a al polinomio

$$P_n(f)(x) = \sum_{k=0}^n \frac{f^{(k)}(a)}{k!} (x - a)^k \in \mathbb{R}[x].$$

El polinomio de Taylor es el único polinomio P de grado menor o igual que n que cumple $P^{(k)}(a) = f^{(k)}(a)$ para $k = 0, \dots, n$. En particular si f es un polinomio de grado menor o igual que n se ha de cumplir $P_n(f) = f$.

Notemos también que $P_0(f) = f(a)$, y que $P_1(f) = f(a) + f'(a)(x - a)$ es la recta tangente a f en a . Nuestra conjetura es que $P_n(f)$ es el polinomio de grado menor o igual que n que más se parece a f alrededor de a .

Ejemplo Consideremos la función $f(x) = x^{1/2}$ y $a = 1$. Calculemos sus derivadas:

Orden	$f^{(n)}(x)$	$f^{(n)}(1)$
0	$x^{1/2}$	1
1	$\frac{1}{2}x^{-1/2}$	$\frac{1}{2}$
2	$-\frac{1}{2} \cdot \frac{1}{2}x^{-3/2}$	$-\frac{1}{2} \cdot \frac{1}{2}$
3	$\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{3}{2}x^{-5/2}$	$\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{3}{2}$
4	$-\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{3}{2} \cdot \frac{5}{2}x^{-7/2}$	$-\frac{1}{2} \cdot \frac{1}{2} \cdot \frac{3}{2} \cdot \frac{5}{2}$

En general se prueba que las derivadas en 1 van alternando el signo, en el numerador tienen el producto de los primeros impares y en el denominador las sucesivas potencias de 2. Con esto podemos calcular cualquier polinomio de Taylor en 1:

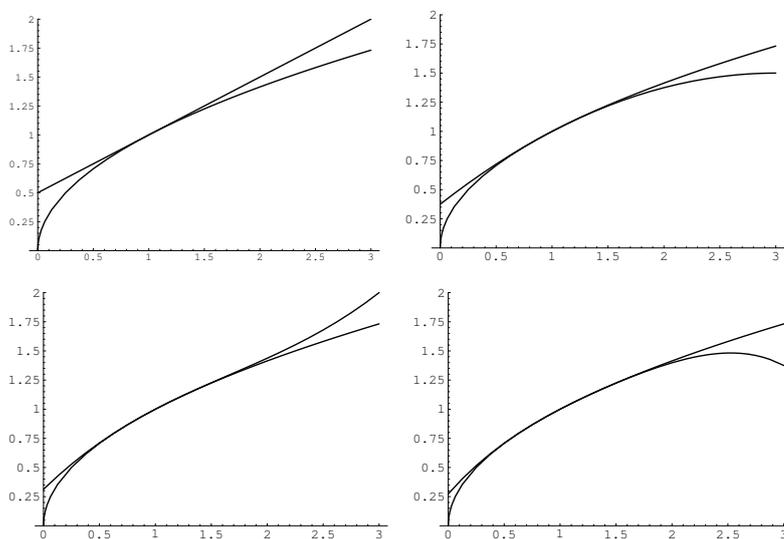
$$P_0(f)(x) = 1,$$

$$P_1(f)(x) = 1 + \frac{1}{2}(x - 1),$$

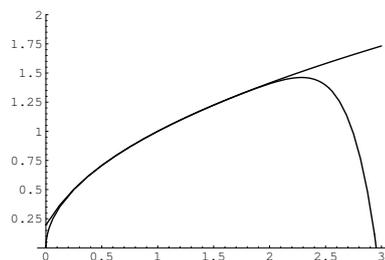
$$P_2(f)(x) = 1 + \frac{1}{2}(x - 1) - \frac{1}{8}(x - 1)^2,$$

$$P_3(f)(x) = 1 + \frac{1}{2}(x - 1) - \frac{1}{8}(x - 1)^2 + \frac{1}{16}(x - 1)^3.$$

Aquí están sus gráficas junto a la de la función. Vemos que el intervalo en que se confunden con la gráfica de f es cada vez mayor.



El polinomio de grado 8 es bastante representativo de lo que sucede cuando n es grande:



Vemos que la aproximación es cada vez mejor en el intervalo $[0, 2]$, pero a partir del 2 el polinomio se aleja bruscamente. Un ejemplo numérico:

$$P_8(f)(1,5) = 1,224729895, \quad \text{mientras que } \sqrt{1,5} = 1,224744871 \dots$$

La aproximación tiene 5 cifras exactas. ■

Los polinomios de Taylor plantean varios problemas importantes. La cuestión principal es si el error producido al aproximar una función de clase C^∞ por sus polinomios de Taylor puede reducirse arbitrariamente aumentando suficientemente el grado.

Definición 4.26 Si f es una función de clase C^∞ en un entorno de un punto a , llamaremos *serie de Taylor* de f en a a la serie funcional

$$\sum_{k=0}^{\infty} \frac{f^{(k)}(a)}{k!} (x-a)^k.$$

La cuestión es si la serie de Taylor de f converge a f . El ejemplo anterior sugiere que la serie de \sqrt{x} en 1 converge a la función en el intervalo $]0, 2[$, pero no parece converger más allá de 2. También hemos de tener presente la posibilidad de que la serie de Taylor de una función f converja a una función distinta de f . Para estudiar estos problemas introducimos el concepto de resto de Taylor:

Sea f una función derivable n veces en un punto a . Llamaremos *resto de Taylor* de grado n de f en a , a la función $R_n(f)(x) = f(x) - P_n(f)(x)$, donde $P_n(f)$ es el polinomio de Taylor de grado n de f en a .

Nuestro problema es determinar el comportamiento del resto de una función. Para ello contamos con el siguiente teorema, que es una generalización del teorema del valor medio.

Teorema 4.27 (Teorema de Taylor) Sea $f : A \rightarrow \mathbb{R}$ una función derivable $n+1$ veces en un intervalo abierto A y $a \in A$. Entonces para cada $x \in A$ existe un $\lambda \in]0, 1[$ tal que si $c = \lambda a + (1-\lambda)x$, se cumple

$$R_n(f)(x) = \frac{f^{(n+1)}(c)}{(n+1)!} (x-a)^{n+1}.$$

DEMOSTRACIÓN: Para $x = a$ es evidente, pues se cumple $P_n(f)(a) = f(a)$ y $R_n(f)(x) = 0$. Supongamos que $x \neq a$. Sea

$$Q(x) = \frac{1}{(x-a)^{n+1}} R_n(f)(x).$$

Sea $F : A \rightarrow \mathbb{R}$ dada por

$$F(t) = f(x) - \left(f(t) + \frac{x-t}{1!} f'(t) + \frac{(x-t)^2}{2!} f''(t) + \dots + \frac{(x-t)^n}{n!} f^n(t) + (x-t)^{n+1} Q(x) \right).$$

La función f y sus n primeras derivadas son continuas y derivables, luego F también es continua y derivable en A . Además $F(x) = f(x) - f(x) = 0$ y

$$F(a) = f(a) - (P_n(f)(a) + (x-a)^{n+1} Q(x)) = R_n(f)(x) - R_n(f)(x) = 0.$$

Por el teorema de Rolle existe un punto entre a y x , o sea, de la forma $c = \lambda a + (1-\lambda)x$, tal que $F'(c) = 0$. Calculemos en general $F'(t)$:

$$F'(t) = 0 - \left(f'(t) - f'(t) + \frac{x-t}{1!} f''(t) - \frac{2(x-t)}{2!} f''(t) + \frac{(x-t)^2}{2!} f'''(t) \dots - \frac{n(x-t)^{n-1}}{n!} f^n(t) + \frac{(x-t)^n}{n!} f^{n+1}(t) - (n+1)(x-t)^n Q(x) \right).$$

Los términos consecutivos se cancelan entre sí, y queda

$$F'(t) = -\frac{(x-t)^n}{n!} f^{n+1}(t) + (n+1)(x-t)^n Q(x).$$

Como $F'(c) = 0$, evaluando en c queda

$$Q(x) = \frac{f^{n+1}(c)}{(n+1)!},$$

y por definición de Q :

$$R_n(f)(x) = \frac{f^{n+1}(c)}{(n+1)!} (x-a)^{n+1}. \quad \blacksquare$$

Así pues, la diferencia entre $P_n(f)(x)$ y $f(x)$ tiene la forma de un monomio más del polinomio de Taylor salvo por el hecho de que la derivada $(n+1)$ -ésima no se evalúa en el punto a , sino en un punto intermedio entre a y x .

Por ejemplo, si las derivadas de f están uniformemente acotadas en un intervalo A , es decir, si existe una misma constante K tal que $|f^n(x)| \leq K$ para todo natural n y para todo $x \in A$, entonces

$$|f(x) - P_n(f)(x)| = \left| \frac{f^{n+1}(c)}{(n+1)!} (x-a)^{n+1} \right| \leq \frac{K|x-a|^{n+1}}{(n+1)!}.$$

En el ejemplo de la página 87 probamos que la sucesión $M^n/n!$ converge a 0, luego la sucesión $\{P_n(f)(x)\}_{n=0}^{\infty}$ tiende a $f(x)$. Por consiguiente:

Teorema 4.28 Si A es un intervalo abierto, $a \in A$, $f \in C^\infty(A)$ y las derivadas de f están uniformemente acotadas en A , entonces para cada punto $x \in A$ se cumple

$$f(x) = \sum_{n=0}^{\infty} \frac{f^{(n)}(a)}{n!} (x-a)^n.$$

Este teorema no es aplicable a \sqrt{x} . En muchos casos, entre ellos el de esta función, los problemas de convergencia de las series de Taylor se vuelven evidentes en el contexto de la teoría de funciones de variable compleja (véase el ejemplo de la página 167 o el teorema 4.34). Los resultados de la sección siguiente resultan de gran ayuda con frecuencia, como tendremos ocasión de comprobar más adelante.

4.6 Series de potencias

En toda esta sección \mathbb{K} será un cuerpo métrico completo. Los casos que más nos interesarán son $\mathbb{K} = \mathbb{C}$ y $\mathbb{K} = \mathbb{R}$.

Definición 4.29 Sea $a \in \mathbb{K}$ y $\{a_n\}_{n=0}^{\infty}$ una sucesión en \mathbb{K} . La *serie de potencias* de coeficientes $\{a_n\}_{n=0}^{\infty}$ y centro a es la serie funcional

$$\sum_{n=0}^{\infty} a_n (z-a)^n.$$

Las series de Taylor son, pues, series de potencias. En muchos casos es fácil determinar en qué puntos converge una serie de potencias. Para verlo necesitamos el concepto de límite superior de una sucesión de números reales. Se trata de lo siguiente:

Sea $\{a_n\}_{n=0}^{\infty}$ una sucesión de números reales. Su *límite superior* es el supremo (en \mathbb{R}) del conjunto de sus puntos adherentes. Lo representaremos mediante $\overline{\lim}_n a_n$. Se cumple que

$$\overline{\lim}_n a_n = \inf_{k \geq 0} \sup_{n \geq k} a_n.$$

En efecto, sea p un punto adherente de $\{a_n\}_{n=0}^{\infty}$. Dados $\epsilon > 0$ y $k \geq 0$, existe un $n \geq k$ tal que $a_n \in]p - \epsilon, p + \epsilon[$, luego $p - \epsilon \leq \sup_{n \geq k} a_n$. Esto vale para todo $\epsilon > 0$, luego $p \leq \sup_{n \geq k} a_n$ para todo $k \geq 0$, luego $p \leq \inf_{k \geq 0} \sup_{n \geq k} a_n$. Como el límite superior es el supremo de estos p , tenemos que $\overline{\lim}_n a_n \leq \inf_{k \geq 0} \sup_{n \geq k} a_n$.

Sea $L = \inf_{k \geq 0} \sup_{n \geq k} a_n$. Dado $\epsilon > 0$, existe un $k \geq 0$ tal que $L \leq \sup_{n \geq k} a_n \leq L + \epsilon$.

Si $L = \sup_{n \geq k} a_n$, entonces existe un $n \geq k$ tal que $L - \epsilon < a_n \leq L$. Si por el contrario $L < \sup_{n \geq k} a_n < L + \epsilon$, entonces existe un $n \geq k$ tal que $L \leq a_n < L + \epsilon$.

En cualquier caso existe un $n \geq k$ tal que $a_n \in]L - \epsilon, L + \epsilon[$. Esto significa que L es un punto adherente de la sucesión, luego $L \leq \overline{\lim}_n a_n$ y tenemos la igualdad. ■

Por la propia definición es claro que si una sucesión converge en $\overline{\mathbb{R}}$ entonces su límite, que es su único punto adherente, coincide con su límite superior. Ahora podemos probar:

Teorema 4.30 Sea $\sum_{n=0}^{\infty} a_n(z-a)^n$ una serie de potencias, sea $R = 1/\overline{\lim}_n \sqrt[n]{|a_n|}$ (entendiendo que $1/0 = +\infty$ y $1/(+\infty) = 0$). Entonces la serie converge absoluta y uniformemente en todo compacto contenido en $B_R(a)$ y diverge en todo punto de $\mathbb{K} \setminus B'_R(a)$ (con el convenio de que $B_{+\infty}(a) = \mathbb{K}$). En particular la serie converge absoluta y puntualmente en $B_R(a)$.

DEMOSTRACIÓN: Sea K un compacto en $B_R(a)$. Veamos que la serie converge absoluta y uniformemente en K . La función $|x-a|$ es continua en K , luego alcanza su máximo r en un punto $x \in K$, es decir, $|x-a| = r$ y para todo $y \in K$ se cumple $|y-a| \leq r$. Así $K \subset B'_r(a)$.

Como $x \in B_R(a)$ ha de ser $r < R$, luego $r \overline{\lim}_n \sqrt[n]{|a_n|} < 1$. Tomemos ρ tal que $r \overline{\lim}_n \sqrt[n]{|a_n|} < \rho < 1$. Como

$$\overline{\lim}_n \sqrt[n]{|a_n|} = \inf_{k \geq 0} \sup_{n \geq k} \sqrt[n]{|a_n|},$$

existe un $k \in \mathbb{N}$ tal que $\sup_{n \geq k} \sqrt[n]{|a_n|} < \rho/r$, luego si $n \geq k$ se cumple $\sqrt[n]{|a_n|} < \rho/r$, luego $|a_n|r^n < \rho^n$. Si $y \in K$ entonces $|y-a| \leq r$, luego $|y-a|^n \leq r^n$, luego

$$|a_n(y-a)^n| \leq |a_n|r^n < \rho^n.$$

Así pues, la serie $\sum_{n=k}^{\infty} a_n(y-a)^n$ está mayorada en K por $\sum_{n=k}^{\infty} \rho^n$, que es convergente por ser geométrica de razón menor que 1. El criterio de Weierstrass nos da que la serie de potencias converge absoluta y uniformemente a una función continua en K .

Ahora veamos que la serie diverge en $\mathbb{K} \setminus B'_R(a)$. Tomamos $x \in \mathbb{K}$ tal que $|x-a| > R$. Entonces $1 < |x-a| \overline{\lim}_n \sqrt[n]{|a_n|}$. Por lo tanto, para todo natural k se cumple $1/|x-a| < \sup_{n \geq k} \sqrt[n]{|a_n|}$, luego existe un $n \geq k$ tal que $\sqrt[n]{|a_n|}|x-a| > 1$, o sea, $|a_n(y-a)^n| > 1$. Esto significa que $a_n(y-a)^n$ no tiende a 0, luego la serie diverge. ■

El número R se llama *radio de convergencia* de la serie de potencias. La bola $B_R(a)$ se llama *disco de convergencia* (que en el caso en que $\mathbb{K} = \mathbb{R}$ será un intervalo). Tenemos, pues que una serie de potencias converge absolutamente en su disco de convergencia y diverge en los puntos exteriores a él (los puntos interiores de su complementario). En cada punto de la frontera del disco la serie puede converger absolutamente, condicionalmente o divergir, según los casos.

Nota En la prueba del teorema anterior hemos visto que si $K \subset B_R(a)$ es compacto, entonces la suma de la serie de potencias es continua en K . Si suponemos que \mathbb{K} es localmente compacto, como es el caso de \mathbb{R} o \mathbb{C} , entonces todo punto de $B_R(a)$ tiene un entorno compacto contenido en $B_R(a)$, luego la suma es continua en $B_R(a)$. ■

A la hora de determinar el radio de convergencia de una serie suele ser útil el teorema siguiente:

Teorema 4.31 Sea $\sum_{n=0}^{\infty} a_n(z-a)^n$ una serie de potencias tal que exista

$$\lim_n \frac{|a_{n+1}|}{|a_n|} = L.$$

Entonces su radio de convergencia es $1/L$.

DEMOSTRACIÓN: Por el teorema anterior, el radio de convergencia de la serie dada es el mismo que el de la serie $\sum_{n=0}^{\infty} |a_n|z^n$. Si $x > 0$, tenemos que

$$\lim_n \frac{|a_{n+1}|x^{n+1}}{|a_n|x^n} = Lx,$$

luego el criterio de D'Alembert implica que la serie converge cuando $Lx < 1$ y diverge si $Lx > 1$. Consecuentemente el radio de convergencia ha de ser $1/L$. ■

Las series de potencias se pueden derivar término a término. Conviene probar un resultado un poco más general:

Teorema 4.32 Sea $\{f_n\}_{n=0}^{\infty}$ una sucesión de funciones $f_n :]a, b[\rightarrow \mathbb{R}$ que converge uniformemente a una función f . Supongamos que todas ellas son derivables en $]a, b[$ y que la sucesión de derivadas converge uniformemente a una función g . Entonces f es derivable y $g = f'$.

DEMOSTRACIÓN: Fijemos un punto $x \in]a, b[$ y consideremos las funciones

$$F_n(y) = \begin{cases} \frac{f_n(x) - f_n(y)}{x - y} & \text{si } y \neq x \\ f'_n(x) & \text{si } y = x \end{cases}$$

Similarmente definimos $F :]a, b[\rightarrow \mathbb{R}$ mediante

$$F(y) = \begin{cases} \frac{f(x) - f(y)}{x - y} & \text{si } y \neq x \\ g(x) & \text{si } y = x \end{cases}$$

El hecho de que f_n sea derivable en x implica que F_n es continua en $]a, b[$. Basta probar que para todo $\epsilon > 0$ existe un número natural n_0 tal que para todo $m, n \geq n_0$ e $y \in]a, b[$ se cumple que $|F_m(y) - F_n(y)| < \epsilon$. En efecto, esto significa que $\{F_n\}_{n=0}^{\infty}$ es (uniformemente) de Cauchy, luego según 3.54 la sucesión converge uniformemente a una función continua, pero dado que converge puntualmente a F , de hecho convergerá uniformemente a F . Tenemos, pues, que F es continua, y a su vez esto implica que f es derivable en x y $f'(x) = g(x)$.

Existe un número natural n_0 tal que si $m, n > n_0$ entonces

$$|f'_n(u) - f'_m(u)| < \frac{\epsilon}{2} \quad \text{para todo } u \in]a, b[.$$

Entonces, dados $y \in]a, b[$ y $m, n \geq n_0$, si $y \neq x$ se cumple

$$\begin{aligned} |F_n(y) - F_m(y)| &= \left| \frac{f_n(x) - f_n(y)}{x - y} - \frac{f_m(x) - f_m(y)}{x - y} \right| \\ &\leq \left| \frac{1}{x - y} \right| |f_n(x) - f_m(x) - f_n(y) + f_m(y)| \end{aligned}$$

Aplicamos el teorema del valor medio a la función $f_n - f_m$ en el intervalo $[y, x]$ (suponemos, por ejemplo, $y < x$). Entonces existe un $u \in]a, b[$ tal que

$$f_n(x) - f_m(x) - f_n(y) + f_m(y) = (f'_n(u) - f'_m(u))(x - y).$$

Así pues,

$$|F_n(y) - F_m(y)| \leq \left| \frac{1}{x - y} \right| |f'_n(u) - f'_m(u)| |x - y| < \frac{\epsilon}{2}.$$

Además, si $y = x$ tenemos

$$|F_n(x) - F_m(x)| = |f'_n(x) - f'_m(x)| < \frac{\epsilon}{2} < \epsilon. \quad \blacksquare$$

Como consecuencia tenemos:

Teorema 4.33 Si $\sum_{n=0}^{\infty} a_n(x - a)^n$ es una serie de potencias en \mathbb{R} , la serie $\sum_{n=1}^{\infty} n a_n(x - a)^{n-1}$ tiene el mismo radio de convergencia y converge a la derivada de la primera.

DEMOSTRACIÓN: Sea R el radio de convergencia de la serie dada, de modo que $R = 1/\lim_n \sqrt[n]{|a_n|}$. Entonces el radio de convergencia de la segunda serie es $1/\lim_n \sqrt[n]{n} \sqrt[n]{|a_n|}$.

Vamos a aceptar que $\lim_n \sqrt[n]{n} = 1$. No es difícil demostrarlo directamente, pero en la sección siguiente² podremos dar una prueba mucho más simple. Teniendo esto en cuenta, es inmediato que $\lim_n \sqrt[n]{|a_n|} = \lim_n \sqrt[n]{n} \sqrt[n]{|a_n|}$, pues una subsucesión de una de las sucesiones converge a un número real si y sólo si la subsucesión correspondiente de la otra converge a ese mismo número real, luego ambas tienen los mismos puntos adherentes y, por consiguiente, el mismo límite superior.

²Véase el ejemplo de la página 171.

Llamemos $f(x)$ a la función definida sobre $]a - R, a + R[$ por la serie dada y $g(x)$ a la función definida por la segunda serie. Hemos de demostrar que $f'(x) = g(x)$ para todo x en el intervalo.

Sea $f_n(x) = \sum_{k=0}^n a_k(x-a)^k$. Se trata de una sucesión de polinomios cuyas derivadas $f'_n(x)$ son las sumas parciales de la segunda serie. Dado un punto $x \in]a - R, a + R[$, tomamos un intervalo cerrado $[x - \delta, x + \delta] \subset]a - R, a + R[$. Por el teorema 4.30, en este intervalo las sucesiones $f_n(x)$ y $f'_n(x)$ convergen absoluta y uniformemente a f y g respectivamente. Ahora basta aplicar el teorema anterior. ■

Nota Como en la demostración anterior falta justificar que $\lim_n \sqrt[n]{n} = 1$, conviene observar que, sin usar esto, el argumento dado prueba que si en un intervalo ambas series convergen, entonces la segunda converge a la derivada de la primera. ■

Teorema 4.34 Sea $f(x) = \sum_{n=0}^{\infty} a_n(x-a)^n$ la función definida en su intervalo de convergencia por una serie de potencias en \mathbb{R} . Entonces f es una función de clase C^∞ en dicho intervalo y $a_n = f^{(n)}(a)/n!$, de modo que la serie que define a f coincide con la serie de Taylor de f .

DEMOSTRACIÓN: El teorema anterior implica que f es una función de clase C^∞ en el intervalo de convergencia de la serie dada, y una simple inducción prueba que

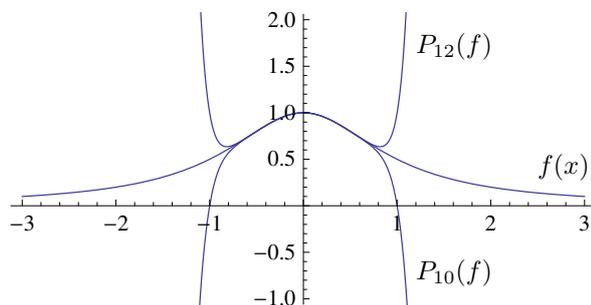
$$f^{(k)}(x) = \sum_{n=k}^{\infty} n(n-1)\cdots(n-k+1)a_n(x-a)^{n-k},$$

luego $f^{(k)}(a) = k! a_k$, ya que todos los términos de la serie valen 0 en a excepto el correspondiente a $n = k$. ■

Ejemplo Consideremos la función $f(x) = 1/(1+x^2)$, que claramente es de clase C^∞ en \mathbb{R} . No es fácil determinar su serie de Taylor en 0 calculando sus derivadas sucesivas, pero se cumple trivialmente que

$$f(x) = \sum_{n=0}^{\infty} (-1)^n x^{2n} \quad \text{para } |x| < 1$$

sin más que aplicar la fórmula para la suma de la serie geométrica de razón $-x^2$. Por el teorema anterior ésta es la serie de Taylor de f y su radio de convergencia es $R = 1$, pues la serie diverge claramente si $x = 1$, aunque también se obtiene inmediatamente por la definición o por el teorema 4.31.



La figura muestra la función f y sus polinomios de Taylor de grados 10 y 12 (el de grado 11 coincide con el de grado 10). A la vista de f , nada indica que la serie de Taylor tenga que dejar de converger en $x = \pm 1$, pero hay una razón por la que tenía que ser así necesariamente: podemos considerar la serie como serie de potencias en \mathbb{C} .

La fórmula para el radio de convergencia es la misma, pero el disco de convergencia es $B_1(0)$, y ahora sí que está claro que el radio de convergencia no podía ser mayor, porque la serie tiene que converger a una función continua en su disco de convergencia, y la función $1/(1+x^2)$ tiende a infinito en $x = \pm i$. Así pues, el obstáculo a la convergencia de la serie en un intervalo mayor que $[-1, 1]$ no está en \mathbb{R} , sino en \mathbb{C} . ■

4.7 La función exponencial

Vamos a aplicar las ideas de las secciones precedentes a la construcción de las funciones más importantes del análisis: la exponencial, la logarítmica y las trigonométricas. En esta sección nos ocuparemos de la función exponencial.

Hasta ahora tenemos definido a^r cuando a es un número real positivo y r es un número racional. No es difícil probar que la función $r \mapsto a^r$ admite una única extensión continua a \mathbb{R} que sigue conservando la propiedad $a^{x+y} = a^x a^y$. Además esta función es infinitamente derivable y coincide en todo punto con su serie de Taylor en 0. En lugar de probar todos estos hechos, lo que haremos será definir la función exponencial a partir de su serie de Taylor, para lo cual no necesitaremos siquiera el hecho de que ya la tenemos definida sobre \mathbb{Q} . No obstante, ahora vamos a suponer la existencia de la función a^x , así como que es derivable, y vamos a calcular su serie de Taylor. Así obtendremos la serie que deberemos tomar como definición.

Sea $f(x) = a^x$. Entonces,

$$f'(x) = \lim_{h \rightarrow 0} \frac{a^{x+h} - a^x}{h} = a^x \lim_{h \rightarrow 0} \frac{a^h - 1}{h} = a^x f'(0).$$

Llamemos $k = f'(0)$. Entonces hemos probado que $f'(x) = kf(x)$, luego por inducción concluimos que f es infinitamente derivable y $f^{(n)}(x) = k^n f(x)$. No puede ser $k = 0$, o de lo contrario f sería constante. Sea $e = a^{1/k} = f(1/k)$.

Entonces la función $g(x) = e^x = a^{x/k} = f(x/k)$ cumple $g'(x) = f'(x/k)(1/k) = f(x/k) = g(x)$, es decir, escogiendo adecuadamente la base e obtenemos una función exponencial que coincide con su derivada. Su serie de Taylor en 0 es entonces fácil de calcular, pues todas las derivadas valen $e^0 = 1$, lo cual nos lleva a la definición siguiente:

Definición 4.35 Llamaremos *función exponencial* a la definida³ por la serie de potencias en \mathbb{C}

$$e^z = \sum_{n=0}^{\infty} \frac{z^n}{n!}.$$

Puesto que

$$\lim_n \frac{1/(n+1)!}{1/n!} = \lim_n \frac{1}{n} = 0,$$

el radio de convergencia es infinito, luego la exponencial está definida sobre todo número complejo z . El teorema 4.33 implica que la exponencial real es derivable,⁴ y su derivada en un punto x es

$$\sum_{n=1}^{\infty} \frac{nx^{n-1}}{n!} = \sum_{n=1}^{\infty} \frac{x^{n-1}}{(n-1)!} = \sum_{n=0}^{\infty} \frac{x^n}{n!} = e^x.$$

Es claro que $e^0 = 1$. Definimos el número

$$e = e^1 = \sum_{n=0}^{\infty} \frac{1}{n!} = 2.7182818284590452353602874 \dots$$

Ahora probamos la ecuación que caracteriza a la función exponencial:

Teorema 4.36 Si $z_1, z_2 \in \mathbb{C}$, entonces $e^{z_1+z_2} = e^{z_1}e^{z_2}$.

DEMOSTRACIÓN: Usamos la fórmula del producto de Cauchy:

$$\begin{aligned} e^{z_1}e^{z_2} &= \sum_{n=0}^{\infty} \frac{z_1^n}{n!} \sum_{n=0}^{\infty} \frac{z_2^n}{n!} = \sum_{n=0}^{\infty} \sum_{k=0}^n \frac{1}{k!(n-k)!} z_1^k z_2^{n-k} \\ &= \sum_{n=0}^{\infty} \sum_{k=0}^n \frac{1}{n!} \binom{n}{k} z_1^k z_2^{n-k} = \sum_{n=0}^{\infty} \frac{(z_1+z_2)^n}{n!} = e^{z_1+z_2}. \end{aligned}$$

■

De aquí obtenemos muchas consecuencias. Por una parte, si n es un número natural no nulo entonces $e^n = e^{1+\dots+1} = e \cdot \dots \cdot e$, es decir, la función exponencial sobre números naturales (incluido el 0) coincide con la exponenciación usual con base e . Así mismo, $1 = e^0 = e^{x-x} = e^x e^{-x}$, luego $e^{-x} = 1/e^x$, con lo que la función exponencial coincide también con la usual cuando el exponente es entero.

³En realidad la habíamos introducido ya en el ejemplo de la página 87.

⁴Notemos que es inmediato que el radio de convergencia de la serie derivada es también infinito, por lo que no estamos usando el hecho no probado de 4.33.

Como los coeficientes de la serie exponencial son positivos, vemos que si $x \geq 0$ entonces $e^x > 0$, y si $x < 0$ entonces $e^x = 1/e^{-x} > 0$. Así pues, $e^x > 0$ para todo número real x .

Como, $(e^{1/n})^n = e^{1/n+\dots+1/n} = e^1 = e$, resulta que $e^{1/n} = \sqrt[n]{e}$. Es fácil ver ahora que $e^{p/q} = \sqrt[q]{e^p}$, luego la función exponencial coincide con la que teníamos definida para exponentes racionales.

Puesto que la derivada es positiva en todo punto, vemos que la función exponencial es estrictamente creciente en \mathbb{R} . En particular es inyectiva. Separando los dos primeros términos de la serie vemos que si $x \geq 0$ entonces $1 + x \leq e^x$, luego $\lim_{x \rightarrow +\infty} e^x = +\infty$. A su vez esto implica que

$$\lim_{x \rightarrow -\infty} e^x = \lim_{x \rightarrow +\infty} e^{-x} = \lim_{x \rightarrow +\infty} \frac{1}{e^x} = 0.$$

Por el teorema de los valores intermedios, la función exponencial biyecta \mathbb{R} con el intervalo $]0, +\infty[$.

Ejemplo Aplicando n veces la regla de L'Hôpital se concluye claramente que

$$\lim_{x \rightarrow +\infty} \frac{x^n}{e^x} = 0,$$

y cambiando x por $1/x$ llegamos a que

$$\lim_{x \rightarrow 0^+} \frac{e^{-1/x}}{x^n} = 0, \quad n = 0, 1, 2, \dots \quad (4.5)$$

Esto implica que la función $h : \mathbb{R} \rightarrow \mathbb{R}$ dada por

$$h(x) = \begin{cases} e^{-1/x} & \text{si } x > 0 \\ 0 & \text{si } x \leq 0 \end{cases}$$

es de clase C^∞ en \mathbb{R} . En efecto, una simple inducción prueba que las derivadas de h para $x > 0$ son de la forma

$$\frac{e^{-1/x}}{x^n} P(x),$$

donde $P(x)$ es un polinomio. De aquí que las derivadas sucesivas de h en 0 existen y valen todas 0. En efecto, admitiendo que existe $h^{(k)}(0) = 0$ (para $k \geq 0$) la derivada $h^{(k+1)}(0)$ se obtiene por un límite cuando $\Delta x \rightarrow 0$ que por la izquierda es claramente 0 y por la derecha es de la forma

$$\lim_{\Delta x \rightarrow 0^+} \frac{e^{-1/\Delta x}}{\Delta x^n} P(\Delta x),$$

de modo que el primer factor tiende a 0 por (4.5) y el segundo está acotado en un entorno de 0. Por lo tanto existe $h^{(k+1)}(0) = 0$.

En particular, la función $h(x^2)$ es de clase C^∞ , sus derivadas son todas nulas en 0 pero es no nula en todo punto distinto de 0. Tenemos así un ejemplo de función de clase C^∞ cuya serie de Taylor en 0 sólo converge (a ella) en 0. ■

La función h del ejemplo anterior permite construir una familia de funciones que nos serán de gran utilidad más adelante:

Teorema 4.37 *Dados números reales $0 \leq a < b$ existe una función $f : \mathbb{R} \rightarrow \mathbb{R}$ de clase C^∞ tal que $f(x) > 0$ si $x \in]a, b[$ y $f(x) = 0$ en caso contrario.*

DEMOSTRACIÓN: En efecto, la función $h_1(x) = h(x - a)$ se anula sólo en los puntos $x \leq a$ y la función $h_2(b - x)$ se anula sólo en los puntos $x \geq b$. Su producto se anula sólo en los puntos exteriores al intervalo $x \in]a, b[$. ■

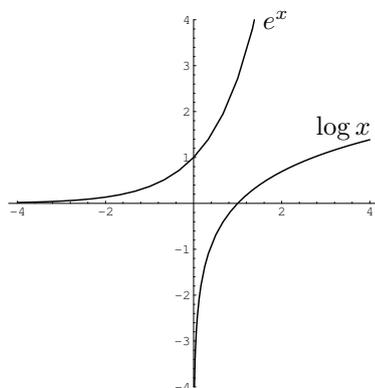
Definición 4.38 Llamaremos *función logarítmica* a la inversa de la función exponencial, $\log :]0, +\infty[\rightarrow \mathbb{R}$.

El teorema de la función inversa nos da que $y = \log x$ es derivable, y su derivada es $y' = 1/(e^y)' = 1/e^y = 1/x$.

De las propiedades de la función exponencial se deducen inmediatamente las de la función logarítmica. Obviamente es una función estrictamente creciente, además verifica la ecuación funcional $\log(xy) = \log x + \log y$. También es claro que $\log 1 = 0$, $\log e = 1$ y

$$\lim_{x \rightarrow +\infty} \log x = +\infty, \quad \lim_{x \rightarrow 0} \log x = -\infty.$$

He aquí las gráficas de las funciones exponencial y logarítmica:



Ejemplo La regla de L'Hôpital nos da inmediatamente que

$$\lim_{x \rightarrow +\infty} \frac{\log x}{x} = \lim_{x \rightarrow +\infty} \frac{1}{x} = 0.$$

En particular, $\lim_n \frac{\log n}{n} = 0$, luego

$$\lim_n \sqrt[n]{n} = \lim_n \sqrt[n]{e^{\log n}} = \lim_n e^{(\log n)/n} = e^0 = 1,$$

que era lo que faltaba para completar la demostración del teorema 4.33. ■

Veamos cómo puede calcularse en la práctica un logaritmo. Es decir, vamos a calcular el desarrollo de Taylor de la función \log . Obviamente no hay un desarrollo en serie sobre todo $]0, +\infty[$. Si desarrollamos alrededor del 1, a lo sumo podemos obtener una serie convergente en $]0, 2[$.

Como las series de potencias centradas en 0 son más fáciles de manejar, vamos a desarrollar en 0 la función $\log(1+x)$. Sus derivadas son

$$(1+x)^{-1}, \quad -(1+x)^{-2}, \quad 2(1+x)^{-3}, \quad -2 \cdot 3(1+x)^{-4}, \dots$$

y en general, la derivada n -sima es $(-1)^{n+1}(n-1)!(1+x)^{-n}$. Puesto que $\log(1+0) = 0$, la serie de Taylor queda:

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} x^n.$$

Como

$$\lim_n \frac{1/(n+1)}{1/n} = \lim_n \frac{n}{n+1} = 1,$$

el radio de convergencia es 1 (como era de esperar), luego la serie converge en $] -1, 1[$. De hecho en $x = -1$ obtenemos una serie divergente, pero en $x = 1$ queda $\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n}$, que es convergente luego, con exactitud, la serie converge en el intervalo $] -1, 1[$.

Según hemos visto en la sección anterior, la función definida por la serie es derivable, y su derivada se obtiene derivando cada monomio, es decir, se trata de la serie geométrica

$$\sum_{n=1}^{\infty} (-1)^{n+1} x^{n-1} = \sum_{n=0}^{\infty} (-x)^n = \frac{1}{1+x}.$$

Resulta, pues, que la serie de Taylor y la función $\log(1+x)$ tienen la misma derivada en $] -1, 1[$. Por lo tanto la diferencia entre ambas funciones es una constante, pero como ambas toman el valor 0 en 0, se concluye que

$$\log(1+x) = \sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} x^n$$

para todo número $x \in] -1, 1[$.

Ejercicio: Estudiando el resto de Taylor de la función $\log(1+x)$, probar que

$$\sum_{n=1}^{\infty} \frac{(-1)^{n+1}}{n} = 1 - \frac{1}{2} + \frac{1}{3} - \frac{1}{4} + \frac{1}{5} - \frac{1}{6} + \dots = \log 2.$$

Para calcular el logaritmo de un número $x > 2$ podemos usar la relación $\log(1/x) = -\log x$.

Ahora podemos definir a^x para cualquier base $a > 0$. La forma más fácil de hacerlo es la siguiente:

Definición 4.39 Sea $a > 0$ y $x \in \mathbb{R}$. Definimos $a^x = e^{x \log a}$.

Notemos que, como $\log e = 1$, en el caso $a = e$ la exponencial que acabamos de definir coincide con la que ya teníamos definida. Sin embargo la función e^x se diferencia de las otras exponenciales en que está definida sobre todo el plano complejo y no sólo sobre la recta real. Más adelante interpretaremos esta extensión compleja. Las funciones exponenciales verifican las propiedades siguientes:

$$\begin{aligned} a^{x+y} &= e^{(x+y) \log a} = e^{x \log a} e^{y \log a} = a^x a^y, \\ \log a^x &= \log e^{x \log a} = x \log a, \\ (a^x)^y &= e^{y \log a^x} = e^{xy \log a} = a^{xy} \\ a^0 &= 1, \quad a^1 = a, \quad a^{-x} = 1/a^x. \end{aligned}$$

Así mismo es claro que a^x coincide con la exponenciación usual cuando x es un número entero y que sobre números racionales es $a^{p/q} = \sqrt[q]{a^p}$. La función y^x considerada como función de dos variables en $]0, +\infty[\times \mathbb{R}$ es continua.

La derivada de a^x es $(\log a)a^x$, la derivada de x^b es

$$e^{b \log x} \frac{b}{x} = \frac{1}{x} b x^b = b x^{b-1}.$$

Finalmente, puesto que la derivada de a^x es siempre positiva si $a > 1$ y siempre negativa si $a < 1$, tenemos que a^x es monótona y biyecta \mathbb{R} con el intervalo $]0, +\infty[$. Por lo tanto tiene una inversa, que representaremos por $\log_a x$ y se llama *logaritmo* en base a de x . Las propiedades algebraicas de estos logaritmos son las mismas que las de la función \log y se demuestran igual. A estas hay que añadir las siguientes, ambas elementales:

$$\log_a x^b = b \log_a x, \quad \log_b x = \frac{\log_a x}{\log_a b}.$$

En particular

$$\log_a x = \frac{\log x}{\log a}.$$

El número e está estrechamente relacionado con los límites de la forma 1^∞ :

Teorema 4.40 Sean $f(x)$, $g(x)$ funciones tales que

$$\lim_{x \rightarrow a} f(x) = 1, \quad \lim_{x \rightarrow a} g(x) = \infty,$$

donde $a \in \mathbb{R}$ o bien $a = \pm\infty$. Si existe $\lim_{x \rightarrow a} (f(x) - 1)g(x) = L$, entonces

$$\lim_{x \rightarrow a} f(x)^{g(x)} = e^L.$$

En particular:

$$e^x = \lim_n \left(1 + \frac{x}{n}\right)^n, \quad e = \lim_n \left(1 + \frac{1}{n}\right)^n.$$

DEMOSTRACIÓN: Aplicando la regla de L'Hôpital vemos que

$$\lim_{t \rightarrow 1} \frac{\log t}{t-1} = 1,$$

luego

$$\lim_{x \rightarrow a} \frac{\log f(x)}{f(x)-1} = 1, \quad \lim_{x \rightarrow a} g(x) \log f(x) = L,$$

luego

$$\lim_{x \rightarrow a} f(x)^{g(x)} = \lim_{x \rightarrow a} e^{g(x) \log f(x)} = e^L.$$

En particular

$$\lim_{t \rightarrow +\infty} \left(1 + \frac{x}{t}\right)^t = e^x,$$

y lo mismo vale si restringimos el límite a la sucesión de los números naturales. ■

Terminamos la sección con una aplicación de los logaritmos junto a técnicas analíticas.

Ejemplo La *media aritmética* y la *media geométrica* de n números reales x_1, \dots, x_n se definen respectivamente como

$$\frac{x_1 + \dots + x_n}{n} \quad \text{y} \quad \sqrt[n]{x_1 \cdots x_n}$$

Vamos a probar que la media geométrica siempre es menor o igual que la media aritmética. A su vez, deduciremos esto de la desigualdad $\log t \leq t - 1$, válida para todo $t > 0$.

Para probar esta desigualdad vemos que la derivada de $f(t) = t - 1 - \log t$ es $1 - 1/t$, que es negativa si $t < 1$ y positiva si $t > 1$. Por el teorema 4.13, f es decreciente en $]0, 1]$ y creciente en $[1, +\infty[$. Puesto que $f(1) = 0$, es claro entonces que $f(t) \geq 0$ para todo $t > 0$.

Respecto a la desigualdad entre las medias, si uno de los números es nulo la media geométrica es nula y el resultado es obvio. Supongamos que son todos no nulos y sea $x = x_1 \cdots x_n$. Entonces

$$\frac{x_i}{\sqrt[n]{x}} - 1 \geq \log \frac{x_i}{\sqrt[n]{x}},$$

luego sumando obtenemos

$$\frac{\sum_{i=1}^n x_i}{\sqrt[n]{x}} - n \geq \log \frac{x_1 \cdots x_n}{x} = 0,$$

con lo que

$$\frac{\sum_{i=1}^n x_i}{n} \geq \sqrt[n]{x}. \quad \blacksquare$$

4.8 Las funciones trigonométricas

En geometría se definen varias funciones de interés, entre las que destacan las funciones seno y coseno. Si llamamos R a la medida de un ángulo recto, entonces la función $\operatorname{sen} x$ está definida sobre \mathbb{R} y tiene periodo $4R$, es decir, $\operatorname{sen}(x + 4R) = \operatorname{sen} x$ para todo $x \in \mathbb{R}$. Así definido, el valor de R es arbitrario, pues podemos tomar cualquier número real como medida de un ángulo recto. Si queremos que existan ángulos unitarios deberemos exigir que $R > 1/4$. Por ejemplo, si tomamos como unidad de ángulo el grado sexagesimal, entonces $R = 90$. Supongamos que las funciones seno y coseno son derivables, así como sus propiedades algebraicas,⁵ y vamos a calcular su serie de Taylor. Con ello obtendremos una definición analítica alternativa.

En primer lugar, como el coseno tiene un máximo en 0, ha de ser $\cos' 0 = 0$. En cualquier otro punto tenemos

$$\begin{aligned} \cos' x &= \lim_{h \rightarrow 0} \frac{\cos(x+h) - \cos x}{h} = \lim_{h \rightarrow 0} \frac{\cos x \cos h - \operatorname{sen} x \operatorname{sen} h - \cos x}{h} \\ &= \cos x \lim_{h \rightarrow 0} \frac{\cos h - 1}{h} - \operatorname{sen} x \lim_{h \rightarrow 0} \frac{\operatorname{sen} h}{h} \\ &= \cos x \cos' 0 - \operatorname{sen} x \operatorname{sen}' 0 = -\operatorname{sen} x \operatorname{sen}' 0. \end{aligned}$$

Si llamamos $k = \operatorname{sen}' 0$ concluimos que $\cos' x = -k \operatorname{sen} x$, y similarmente llegamos a que $\operatorname{sen}' x = k \cos x$.

Del mismo modo que hicimos con la exponencial, podemos normalizar las funciones seno y coseno cambiándolas por $\operatorname{sen}(x/k)$ y $\cos(x/k)$. Geométricamente esto significa fijar una unidad de ángulos. Entonces tenemos $\operatorname{sen}' x = \cos x$ y $\cos' x = -\operatorname{sen} x$.

Vemos entonces que las funciones seno y coseno son infinitamente derivables, y sus derivadas están uniformemente acotadas por 1, luego las series de Taylor deben converger en \mathbb{R} a las funciones respectivas. Puesto que $\operatorname{sen} 0 = 0$ y $\cos 0 = 1$, las series han de ser las que consideramos en la definición siguiente:

Definición 4.41 Llamaremos *seno* y *coseno* a las funciones definidas por las series de potencias

$$\operatorname{sen} z = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n+1)!} z^{2n+1}, \quad \cos z = \sum_{n=0}^{\infty} \frac{(-1)^n}{(2n)!} z^{2n}.$$

No es difícil probar directamente la convergencia de estas series sobre todo el plano complejo. El teorema siguiente muestra una sorprendente conexión entre las funciones trigonométricas así definidas y la función exponencial. Observemos que la prueba contiene otra demostración alternativa de la convergencia de estas series.

Teorema 4.42 Para todo $z \in \mathbb{C}$ se cumple

$$\operatorname{sen} z = \frac{e^{iz} - e^{-iz}}{2i}, \quad \cos z = \frac{e^{iz} + e^{-iz}}{2}.$$

⁵Véase el teorema [G 4.46].

DEMOSTRACIÓN:

$$\frac{e^{iz} + e^{-iz}}{2} = \frac{\sum_{n=0}^{\infty} i^n \frac{z^n}{n!} + \sum_{n=0}^{\infty} (-i)^n \frac{z^n}{n!}}{2} = \sum_{n=0}^{\infty} \frac{i^n + (-i)^n}{2} \frac{z^n}{n!}.$$

Ahora bien, la sucesión $(i^n + (-i)^n)/2$ es simplemente $1, 0, -1, 0, 1, 0, -1, 0, \dots$, luego queda la serie del coseno. Similarmente se razona con el seno. ■

Derivando término a término las series de Taylor se concluye fácilmente que

$$\text{sen}' x = \cos x, \quad \text{cos}' x = -\text{sen } x.$$

Las fórmulas siguientes son todas consecuencias sencillas del teorema anterior:

$$\begin{aligned} \text{sen}^2 z + \text{cos}^2 z &= 1 \\ \text{sen}(x + y) &= \text{sen } x \cos y + \cos x \text{sen } y, \\ \text{cos}(x + y) &= \cos x \cos y - \text{sen } x \text{sen } y, \\ e^{iz} &= \cos z + i \text{sen } z. \end{aligned}$$

La primera fórmula implica que si $x \in \mathbb{R}$ entonces $-1 \leq \text{sen } x, \cos x \leq 1$. De la última se sigue que para todo $x, y \in \mathbb{R}$ se cumple

$$e^{x+iy} = e^x (\cos y + i \text{sen } y),$$

con lo que tenemos descrita la exponencial compleja en términos de la exponencial real y de las funciones seno y coseno reales.

El hecho de que $\text{sen}' 0 = \cos 0 = 1$ equivale a

$$\lim_{x \rightarrow 0} \frac{\text{sen } x}{x} = 1,$$

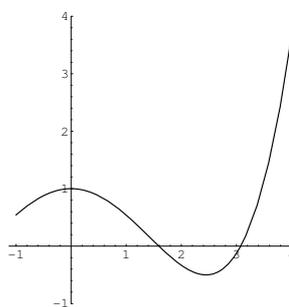
que es otra propiedad del seno que conviene recordar.

Vamos a probar ahora la periodicidad de las funciones trigonométricas reales. El punto más delicado es demostrar que $\cos x$ se anula en algún $x \neq 0$. Para ello probaremos que el coseno es menor o igual que los cuatro primeros términos de su serie de Taylor:

$$\cos x \leq 1 - \frac{x^2}{2} + \frac{x^4}{24}.$$

Esto equivale a probar que $1 - x^2/2 + x^4/24 - \cos x \geq 0$ para $x \geq 0$. Puesto que esta función vale 0 en 0, basta probar que su derivada es positiva. Dicha derivada es $\text{sen } x - x + x^3/6$. Esta función vale también 0 en 0, luego para probar que es positiva (para $x \geq 0$) basta ver que su derivada lo es. Dicha derivada es $\cos x - 1 + x^2/2$. Por el mismo argumento derivamos una vez más y obtenemos $x - \text{sen } x$. Al derivar una vez más llegamos a $1 - \cos x$, que sabemos que es positiva.

Vemos que la gráfica del polinomio de Taylor de grado 4 en 0 de $\cos x$ toma valores negativos. De hecho un simple cálculo nos da que en $\sqrt{3}$ toma el valor $-1/8$, luego $\cos \sqrt{3} \leq -1/8$. Como $\cos 0 = 1$, por continuidad existe un punto $0 < x < \sqrt{3}$ tal que $\cos x = 0$.



Sea $A = \{x > 0 \mid \cos x = 0\}$. El conjunto A es la antiimagen de $\{0\}$ por la aplicación coseno restringida a $[0, +\infty[$. Como $\{0\}$ es cerrado y \cos es continua, A es un cerrado. El ínfimo de un conjunto está en su clausura, luego $F = \inf A \in A$ y así $\cos F = 0$. Es obvio que $F \geq 0$, y como $\cos 0 \neq 0$, ha de ser $0 < F < \sqrt{3}$.

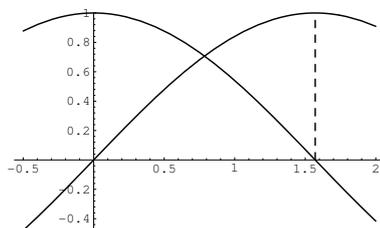
Es costumbre llamar $\pi = 2F$. Así, se cumple $0 < \pi < 2\sqrt{3}$, $\cos(\pi/2) = 0$, pero $\cos x > 0$ en el intervalo $[0, \pi/2[$.

Ahora, como el coseno es la derivada del seno, resulta que $\sen x$ es estrictamente creciente en el intervalo $[0, \pi/2]$. Como $\sen 0 = 0$, resulta que $\sen x \geq 0$ en $[0, \pi/2]$. Concretamente, $\sen(\pi/2)$ es un número positivo que cumple

$$\sen^2(\pi/2) + \cos^2(\pi/2) = \sen^2(\pi/2) + 0 = 1,$$

luego $\sen(\pi/2) = 1$.

Además, $\cos' x = -\sen x \leq 0$ en $[0, \pi/2]$, luego el coseno es estrictamente decreciente en $[0, \pi/2]$. En resumen, tenemos demostrado lo que refleja la gráfica siguiente:



Incidentalmente hemos probado una desigualdad que a veces es de interés: si $x \geq 0$ entonces $\sen x \leq x$. Más en general, $|\sen x| \leq |x|$.

El comportamiento de las funciones seno y coseno fuera del intervalo $[0, \pi/2]$ se deduce de las relaciones trigonométricas que ya hemos probado. Por ejemplo, expresando $\pi = \pi/2 + \pi/2$ obtenemos $\sen \pi = 0$, $\cos \pi = -1$, y a su vez de aquí $\sen 2\pi = 0$, $\cos 2\pi = 1$. Ahora

$$\sen(x + 2\pi) = \sen x, \quad \cos(x + 2\pi) = \cos x,$$

lo que prueba que ambas funciones son periódicas y basta estudiarlas en el intervalo $[0, 2\pi]$.

Teorema 4.43 Sea $z \in \mathbb{C}$, $z \neq 0$ y sea $a \in \mathbb{R}$. Entonces existe un único número real $\theta \in [a, a + 2\pi[$ tal que $z = |z|e^{i\theta} = |z|(\cos \theta + i \sen \theta)$.

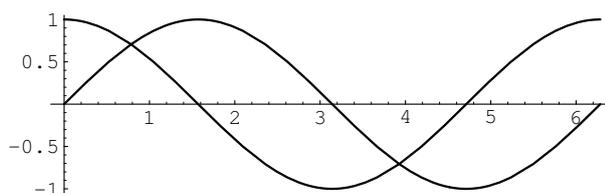
DEMOSTRACIÓN: Sea $z/|z| = x + iy$. Entonces $x^2 + y^2 = 1$. Distinguiamos cuatro casos: según el signo de x e y . Todos son análogos, así que supondremos por ejemplo $x \leq 0$, $y \geq 0$. Más concretamente tenemos $-1 \leq x \leq 0$. En el intervalo $[\pi, 3\pi/2]$ se cumple $\cos \pi = -1$, $\cos 3\pi/2 = 0$, luego por continuidad existe un número $\phi \in [\pi, 3\pi/2]$ tal que $\cos \phi = x$. Entonces $1 = x^2 + y^2 = \cos^2 \phi + \sin^2 \phi$, por lo que $y^2 = \sin^2 \phi$ y, como ambos son negativos, ha de ser $y = \sin \phi$. Así pues, $z = |z|(\cos \phi + i \sin \phi) = |z|e^{i\phi}$.

Existe un número entero p tal que $\theta = 2p\pi + \phi \in [a, a + 2\pi[$. Entonces, teniendo en cuenta que $e^{2p\pi i} = 1$, resulta que $z = |z|e^{i\phi}e^{2p\pi i} = |z|e^{i\theta}$.

La unicidad se debe a que si $|z|e^{i\theta_1} = |z|e^{i\theta_2}$, entonces $e^{i(\theta_1 - \theta_2)} = 1$, luego $\cos(\theta_1 - \theta_2) = 1$ y $\sin(\theta_1 - \theta_2) = 0$, ahora bien, $\cos x = 1$ y $\sin x = 0$ sólo ocurre en $x = 0$ en el intervalo $[0, 2\pi[$, luego sólo ocurre en los números reales de la forma $2k\pi$, con $k \in \mathbb{Z}$. Así pues $\theta_1 - \theta_2 = 2k\pi$, y si ambos están en el intervalo $[a, a + 2\pi[$, ha de ser $\theta_1 = \theta_2$. ■

Un *argumento* de un número complejo $z \neq 0$ es un número real θ tal que $z = |z|e^{i\theta}$. Hemos probado que cada número complejo no nulo tiene un único argumento en cada intervalo $[a, a + 2\pi[$. En particular en el intervalo $[0, 2\pi[$.

A partir de lo que ya hemos probado es fácil obtener todos los resultados que se demuestran en geometría.⁶ He aquí sus gráficas en $[0, 2\pi[$:

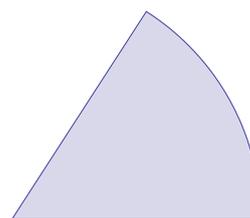


Recordemos que para conseguir que la derivada del seno fuera el coseno hemos tenido que fijar una medida de ángulos concreta. El ángulo de medida 1 respecto a esta unidad, es decir, el ángulo que forman los vectores $(1, 0)$ y $(\cos 1, \sin 1)$, recibe el nombre de *radián*⁷. La figura muestra un radián.

Las funciones seno y coseno nos permiten mostrar algunos ejemplos de interés sobre derivabilidad:

Ejemplo La función $f : \mathbb{R} \rightarrow \mathbb{R}$ dada por

$$f(x) = \begin{cases} x^2 \operatorname{sen} \frac{1}{x} & \text{si } x \neq 0 \\ 0 & \text{si } x = 0 \end{cases}$$



⁶Concretamente, todas las propiedades enunciadas en [G 4.47], luego [G 4.47] nos da que las funciones que acabamos de construir son las mismas construidas en el capítulo IV de [G].

⁷El nombre se debe a que, según probaremos al final de esta sección, si un arco de circunferencia mide un radián, entonces su longitud es igual al radio. Sería más adecuado llamarlo ángulo "radiante".

es derivable en \mathbb{R} . El único punto donde esto no es evidente es $x = 0$, pero

$$f'(0) = \lim_{h \rightarrow 0} h \operatorname{sen} \frac{1}{h} = 0.$$

Para probar esto observamos en general que el producto de una función acotada por otra que tiende a 0 tiende a 0 (basta aplicar la definición de límite).

Sin embargo la derivada no es continua, pues en puntos distintos de 0 vale

$$f'(x) = 2x \operatorname{sen} \frac{1}{x} - \cos \frac{1}{x},$$

y es fácil ver que el primer sumando tiende a 0 en 0 (como antes), mientras que el segundo no tiene límite, luego no existe $\lim_{x \rightarrow 0} f'(x)$.

Los mismos cálculos que acabamos de realizar prueban una limitación de la regla de L'Hôpital. Consideremos la función $g(x) = x$. Entonces

$$\lim_{x \rightarrow 0} \frac{f(x)}{g(x)} = 0,$$

pero si intentamos calcular el límite por la regla de L'Hôpital nos encontramos con

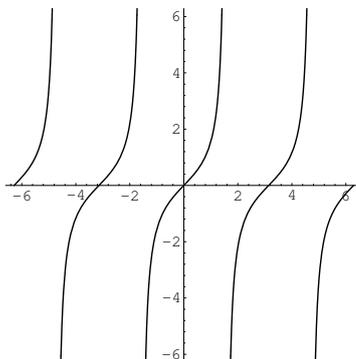
$$\lim_{x \rightarrow 0} \frac{f'(x)}{g'(x)} = \lim_{x \rightarrow 0} f'(x),$$

y ya hemos visto que este límite no existe. La regla de L'Hôpital sólo afirma que si existe el límite del cociente de derivadas también existe el límite original y ambos coinciden, pero es importante recordar que si el segundo límite no existe de ahí no podemos deducir que el primero tampoco exista. ■

Otra función trigonométrica importante es la tangente, definida como

$$\tan z = \frac{\operatorname{sen} z}{\operatorname{cos} z}.$$

No es difícil probar que la función coseno se anula únicamente sobre los múltiplos enteros de $\pi/2$ (no tiene ceros imaginarios). En efecto, si $\cos z = 0$, por definición $e^{iz} + e^{-iz} = 0$, luego $e^{2iz} = -1$. Si $z = a + bi$, queda $e^{2ia}e^{-2b} = -1$ y tomando módulos, $e^{-2b} = 1$, luego $b = 0$ y z es real. Igualmente ocurre con el seno. Por lo tanto la tangente está definida sobre todos los números complejos que no son múltiplos enteros de $\pi/2$. Claramente su restricción a \mathbb{R} es derivable en su dominio, y su derivada es $1/\cos^2 x = 1 + \tan^2 x$. En particular es siempre positiva, luego la tangente es creciente. Su gráfica es:



Es fácil ver que la función tangente biyecta el intervalo $]-\pi/2, \pi/2[$ con la recta real. Junto con ésta, tenemos también las biyecciones siguientes:

$$\text{sen} : \left[-\frac{\pi}{2}, \frac{\pi}{2}\right] \longrightarrow [-1, 1], \quad \text{cos} : [0, \pi] \longrightarrow [-1, 1].$$

Por lo tanto podemos definir las funciones inversas, llamadas respectivamente, *arco seno*, *arco coseno* y *arco tangente*:

$$\text{arcsen} : [-1, 1] \longrightarrow \left[-\frac{\pi}{2}, \frac{\pi}{2}\right], \quad \text{arccos} : [-1, 1] \longrightarrow [0, \pi],$$

$$\text{arctan} : \mathbb{R} \longrightarrow]-\pi/2, \pi/2[.$$

El teorema de la función inversa permite calcular sus derivadas:

$$\text{arcsen}' x = \frac{1}{\sqrt{1-x^2}}, \quad \text{arccos}' x = \frac{-1}{\sqrt{1-x^2}}, \quad \text{arctan}' x = \frac{1}{1+x^2}.$$

Por ejemplo, si $y = \text{arcsen } x$, entonces $x = \text{sen } y$, luego $\frac{dx}{dy} = \cos y$, luego

$$\frac{dy}{dx} = \frac{1}{\cos y} = \frac{1}{\sqrt{1-\text{sen}^2 y}} = \frac{1}{\sqrt{1-x^2}}.$$

Vamos a calcular la serie de Taylor de la función arco tangente. No podemos calcular directamente las derivadas, pues las expresiones que se obtienen son cada vez más complicadas y no permiten obtener una fórmula general. En su lugar emplearemos la misma técnica que hemos usado en la sección anterior para calcular la serie del logaritmo. Claramente

$$\frac{1}{1+x^2} = \frac{1}{1-(-x^2)} = \sum_{n=0}^{\infty} (-1)^n x^{2n}, \quad \text{para } |x| < 1.$$

Ahora es fácil obtener una serie cuya derivada sea la serie anterior, a saber:

$$\sum_{n=0}^{\infty} \frac{(-1)^n}{2n+1} x^{2n+1}.$$

Es fácil ver que su radio de convergencia es 1.

Por lo tanto esta serie se diferencia en una constante de la función $\arctan x$ en el intervalo $] -1, 1[$, pero como ambas funciones toman el valor 0 en $x = 0$, concluimos que son iguales, o sea:

$$\arctan x = \sum_{n=0}^{\infty} \frac{(-1)^n}{2n+1} x^{2n+1}, \quad \text{para } |x| < 1.$$

Cuando $x = \pm 1$ la serie se convierte en una serie alternada cuyo término general es decreciente y tiende a 0, luego por el criterio de Leibniz también converge. No vamos a demostrarlo aquí (véase el ejemplo final de la sección siguiente), pero el límite resulta ser $\arctan(\pm 1)$. Puesto que $\arctan 1 = \pi/4$, esto nos lleva a la conocida fórmula de Leibniz para el cálculo de π :

$$\pi = 4 \left(1 - \frac{1}{3} + \frac{1}{5} - \frac{1}{7} + \frac{1}{9} - \frac{1}{11} + \dots \right).$$

Esta fórmula converge muy lentamente a π , en el sentido de que es necesario calcular muchos términos para obtener pocas cifras exactas. Hay otras expresiones más complicadas pero más eficientes. Veamos una de ellas. Es fácil probar la fórmula de la tangente del ángulo doble:

$$\tan 2\alpha = \frac{2 \tan \alpha}{1 - \tan^2 \alpha}.$$

De aquí se sigue que

$$\tan \left(2 \arctan \frac{1}{5} \right) = \frac{5}{12}, \quad \tan \left(4 \arctan \frac{1}{5} \right) = \frac{120}{119}.$$

Teniendo en cuenta que $\arctan(-x) = -\arctan x$ resulta

$$\tan \left(4 \arctan \frac{1}{5} - \arctan \frac{1}{239} \right) = \frac{\frac{120}{119} - \frac{1}{239}}{1 + \frac{120}{119} \frac{1}{239}} = 1,$$

con lo que, finalmente,

$$\pi = 4 \left(4 \arctan \frac{1}{5} - \arctan \frac{1}{239} \right),$$

o sea,

$$\pi = \sum_{n=0}^{\infty} \frac{(-1)^n 4}{2n+1} \left(\frac{4}{5^{2n+1}} - \frac{1}{239^{2n+1}} \right).$$

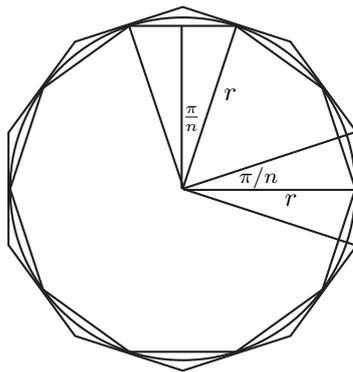
Esta serie converge muy rápidamente. Veamos sus primeras sumas parciales:

0	3.18326359832635983263598326360
1	3.14059702932606031430453110658
2	3.14162102932503442504683251712
3	3.14159177218217729501821229111
4	3.14159268240439951724025983607
5	3.14159265261530860814935074767
6	3.14159265362355476199550459382
7	3.14159265358860222866217126049
8	3.14159265358983584748570067225
9	3.14159265358979169691727961962
10	3.14159265358979329474737485772
11	3.14159265358979323639184094467
12	3.14159265358979323853932459267
13	3.14159265358979323845978816126
14	3.14159265358979323846275020768
15	3.14159265358979323846263936981

Vemos que con 6 sumas tenemos ya una aproximación con diez cifras exactas, y con 15 superamos las 20 cifras.

La longitud de la circunferencia Vamos a calcular la longitud de una circunferencia de radio r . La idea intuitiva de lo que pretendemos es clara: se trata de determinar, por ejemplo, qué longitud de hilo hace falta para rodear una columna cilíndrica de radio r , pero desde un punto de vista técnico, no podemos calcular, sino más bien definir, la longitud de una circunferencia, puesto que no tenemos ninguna definición de la que partir. La cuestión es dar una definición que podamos considerar que se corresponde con lo que queremos calcular.

Nos basamos en un argumento clásico: consideremos una circunferencia de radio r junto con un polígono regular de n lados inscrito en ella (es decir, con sus vértices sobre la circunferencia) y otro circunscrito (es decir, con los puntos medios de sus lados sobre la circunferencia).



Los griegos calcularon la longitud L de la circunferencia postulando que la longitud (o perímetro) I_n del polígono inscrito y la longitud C_n del polígono

circunscrito deben cumplir $I_n < L < C_n$. Concretamente:

$$I_n = 2nr \operatorname{sen} \frac{\pi}{n} < L < C_n = 2nr \tan \frac{\pi}{n}.$$

Vamos a probar que existen los límites $I = \lim_n I_n$ y $C = \lim_n C_n$, lo que obliga a que $I \leq L \leq C$. Para ello podemos sustituir n por una variable real y aplicar la regla de L'Hôpital:

$$\begin{aligned} I &= 2r \lim_{x \rightarrow +\infty} \frac{\operatorname{sen}(\pi/x)}{1/x} = 2r \lim_{x \rightarrow +\infty} \frac{\pi \cos(\pi/x)(-1/x^2)}{-1/x^2} \\ &= 2\pi r \lim_{x \rightarrow +\infty} \cos(\pi/x) = 2\pi r. \\ C &= 2r \lim_{x \rightarrow +\infty} \frac{\tan(\pi/x)}{1/x} = 2r \lim_{x \rightarrow +\infty} \frac{\pi(1 + \tan^2(\pi/x))(-1/x^2)}{-1/x^2} \\ &= 2\pi r \lim_{x \rightarrow +\infty} (1 + \tan^2(\pi/x)) = 2\pi r. \end{aligned}$$

Por lo tanto tiene que ser también $L = 2\pi r$. De hecho, los griegos llamaron π a esta constante por ser la proporción entre el *perímetro* de la circunferencia y su diámetro. ■

4.9 Las funciones hiperbólicas

Vamos a estudiar los valores de las funciones trigonométricas sobre el eje imaginario. Observemos que, según el teorema 4.42,

$$\operatorname{sen} iz = i \frac{e^z - e^{-z}}{2}, \quad \cos iz = \frac{e^z + e^{-z}}{2}.$$

Así pues, la función coseno toma valores reales sobre el eje imaginario, mientras que la función seno toma valores imaginarios puros. Corregimos esto eliminando el factor i en la definición siguiente:

Definición 4.44 Las funciones *seno hiperbólico*, *coseno hiperbólico* y *tangente hiperbólica* se definen como las funciones $\operatorname{senh}, \operatorname{cosh}, \operatorname{tanh} : \mathbb{C} \rightarrow \mathbb{C}$ dadas por

$$\begin{aligned} \operatorname{senh} z &= -i \operatorname{sen} iz = \frac{e^z - e^{-z}}{2}, & \operatorname{cosh} z &= \cos iz = \frac{e^z + e^{-z}}{2}, \\ \operatorname{tanh} z &= \frac{\operatorname{senh} z}{\operatorname{cosh} z} = -i \tan iz = \frac{e^z - e^{-z}}{e^z + e^{-z}}. \end{aligned}$$

De la expresión en términos de las funciones trigonométricas se sigue inmediatamente la relación fundamental

$$\operatorname{cosh}^2 z - \operatorname{senh}^2 z = 1.$$

Aquí vamos a estudiar únicamente la restricción a la recta real de estas funciones. De la definición se sigue inmediatamente que son funciones infinitamente derivables y que

$$\sinh' x = \cosh x, \quad \cosh' x = \sinh x, \quad \tanh' x = \frac{1}{\cosh^2 x} = 1 - \tanh^2 x.$$

Como claramente $\cosh x > 0$, para todo $x \in \mathbb{R}$, vemos que la función seno hiperbólico es estrictamente creciente en \mathbb{R} . También es inmediato que

$$\lim_{x \rightarrow +\infty} \sinh x = +\infty, \quad \lim_{x \rightarrow -\infty} \sinh x = -\infty,$$

luego $\sinh : \mathbb{R} \rightarrow \mathbb{R}$ es biyectiva y creciente.

Como $\sinh 0 = 0$, resulta que el seno hiperbólico es negativo sobre los números negativos y positivo sobre los números positivos. Esto implica a su vez que el coseno hiperbólico es decreciente en $]-\infty, 0[$ y creciente en $]0, +\infty[$. Por lo tanto toma su valor mínimo en 0, donde $\cosh 0 = 1$. Como

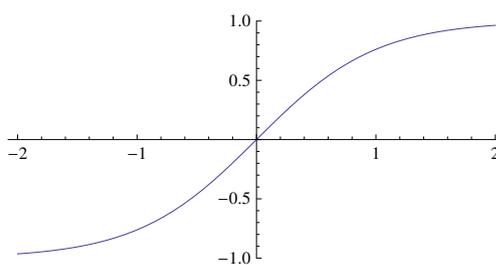
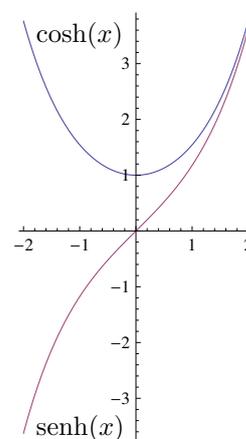
$$\lim_{x \rightarrow -\infty} \cosh x = \lim_{x \rightarrow +\infty} \cosh x = +\infty,$$

resulta que el coseno hiperbólico se restringe a dos biyecciones

$$]-\infty, 0] \rightarrow [1, +\infty[, \quad]0, +\infty[\rightarrow [1, +\infty[,$$

la primera decreciente, la segunda creciente. En suma, hemos probado que las gráficas del seno y el coseno hiperbólico tienen el aspecto que muestra la figura.

También es fácil probar que la tangente hiperbólica es como muestra la figura siguiente:



Es creciente, pues su derivada es trivialmente positiva, y tiende a -1 en $-\infty$ y a 1 en $+\infty$.

Ahora podemos entender el nombre de “funciones hiperbólicas”: igual que $t \mapsto (\cos t, \sin t)$ es una parametrización de la circunferencia unitaria, la curva

$t \mapsto (\cosh t, \sinh t)$ parametriza una de las ramas de la hipérbola⁸ de ecuación $x^2 - y^2 = 1$.

Las fórmulas siguientes se demuestran sin más que introducir unidades imaginarias en las fórmulas análogas para el seno y el coseno:

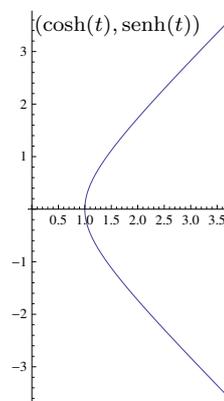
$$\cosh(x + y) = \cosh x \cosh y + \sinh x \sinh y,$$

$$\sinh(x + y) = \cosh x \sinh y + \sinh x \cosh y.$$

También es claro que

$$\cosh 0 = 1, \quad \sinh 0 = 0,$$

$$\cosh(-\alpha) = \cosh \alpha, \quad \sinh(-\alpha) = -\sinh \alpha.$$



Introduciendo la unidad imaginaria en las series de Taylor del seno y el coseno obtenemos inmediatamente las del seno y el coseno hiperbólicos:

$$\sinh x = \sum_{n=0}^{\infty} \frac{x^{2n+1}}{(2n+1)!}, \quad \cosh x = \sum_{n=0}^{\infty} \frac{x^{2n}}{(2n)!}.$$

Observemos que

$$\lim_{x \rightarrow 0} \frac{\sinh x}{x} = 1.$$

Esto puede deducirse de la serie de Taylor, o bien mediante la regla de L'Hôpital.

Las inversas de las funciones seno, coseno y tangente hiperbólicas se llaman, respectivamente, *argumento del seno*, *del coseno* y *de la tangente hiperbólica*:

$$\arg \sinh : \mathbb{R} \longrightarrow \mathbb{R}, \quad \arg \cosh : [1, +\infty[\longrightarrow [0, +\infty[, \quad \arg \tanh :]-1, 1[\longrightarrow \mathbb{R}$$

Podemos obtener una expresión explícita de ambas funciones sin más que plantear

$$y = \arg \sinh x \leftrightarrow x = \sinh y = \frac{e^y - e^{-y}}{2} \leftrightarrow e^y - 2x - e^{-y} = 0$$

$$\leftrightarrow (e^y)^2 - 2xe^y - 1 = 0 \leftrightarrow e^y = x + \sqrt{x^2 + 1}.$$

⁸Más en general, (véase [G, sección 10.1]) toda elipse admite, respecto del sistema de referencia adecuado, la ecuación

$$\left(\frac{x}{a}\right)^2 + \left(\frac{y}{b}\right)^2 = 1$$

y toda hipérbola admite la ecuación

$$\left(\frac{x}{a}\right)^2 - \left(\frac{y}{b}\right)^2 = 1.$$

Por lo tanto, toda elipse puede parametrizarse en la forma $(\frac{\cos t}{a}, \frac{\sin t}{b})$, y toda hipérbola puede parametrizarse en la forma $(\pm \frac{\cosh t}{a}, \frac{\sinh t}{b})$.

Por lo tanto:

$$\arg \sinh x = \log(x + \sqrt{x^2 + 1}),$$

e igualmente se obtiene que

$$\arg \cosh x = \log(x + \sqrt{x^2 - 1}), \quad \arg \tanh x = \frac{1}{2} \log \frac{1+x}{1-x}.$$

Aplicando el teorema sobre la derivada de funciones inversas resulta inmediatamente que

$$\frac{d \arg \sinh x}{dx} = \frac{1}{\sqrt{x^2 + 1}}, \quad \frac{d \arg \cosh x}{dx} = \frac{1}{\sqrt{x^2 - 1}}, \quad \frac{d \arg \tanh x}{dx} = \frac{1}{1 - x^2}.$$

4.10 Primitivas

El teorema 4.15 afirma que una función derivable está completamente determinada por su derivada y su valor en un punto. A menudo se plantea el problema práctico de determinar una función a partir de estos datos.

Definición 4.45 Diremos que una función $F : I \rightarrow \mathbb{R}$ definida en un intervalo abierto I es una *primitiva* de una función $f : I \rightarrow \mathbb{R}$ si F es derivable en I y $F' = f$.

En estos términos, lo que afirma el teorema 4.15 es que si una función f admite una primitiva F en un intervalo I , entonces el conjunto de todas las primitivas de f en I está formado por las funciones de la forma $F + c$, para cada $c \in \mathbb{R}$. Representaremos por

$$\int f(x) dx$$

al conjunto de todas las primitivas de f en un intervalo prefijado. Si F es una de ellas, es habitual expresarlo en la forma

$$\int f(x) dx = F(x) + c.$$

Esta notación supone un abuso de lenguaje, pues estamos igualando el conjunto de todas las primitivas con la expresión genérica para una de ellas, pero esto no debería dar lugar a ningún malentendido.

Nota En principio, los signos $\int dx$ deben considerarse como una única notación “indivisible”, es decir, que en lugar de $\int f(x) dx$ podríamos haber usado la notación $I(f)$ y no habría ninguna diferencia. Simplemente estamos asociando a cada función un conjunto de funciones (que puede ser vacío). No obstante, veremos que la notación $\int f(x) dx$ es consistente con el concepto de diferencial de una función de una variable, es decir, $dy = y' dx$.

Por ejemplo, podemos convenir en que a partir de esta igualdad podemos pasar a

$$\int dy = \int y' dx$$

(con lo que las diferenciales dy , dx han dejado de ser diferenciales para convertirse en mera “notación”) si entendemos la igualdad como la definición del miembro izquierdo, para toda función derivable $y(x)$. Trivialmente entonces:

$$\int dy = y(x) + c.$$

Todas las manipulaciones que hagamos pasando de diferenciales a integrales pueden justificarse trivialmente a partir de las propiedades elementales de las diferenciales y las primitivas. No obstante, en el capítulo X veremos que existe un fundamento teórico que explica por qué los signos \int y dx “combinan tan bien”. ■

Cada regla de derivación da lugar a una regla de integración. Por ejemplo, el hecho de que la derivada de una suma es la suma de las derivadas implica que una suma tiene por primitiva a la suma de las primitivas. Más en general, dadas dos funciones u , v y números reales α , β ,

$$\int (\alpha u + \beta v) dx = \alpha \int u dx + \beta \int v dx$$

Esto hay que entenderlo como que las primitivas de $\alpha u + \beta v$ son las funciones que se obtienen multiplicando por α y β respectivamente una primitiva de u y otra de v . Uniendo este hecho con la regla evidente:

$$\int x^n = \frac{x^{n+1}}{n+1} + c, \quad n \neq -1,$$

podemos integrar cualquier polinomio. El caso exceptuado es claro:

$$\int x^{-1} dx = \log x + c.$$

Aunque una función f tiene infinitas primitivas F en un intervalo I (en caso de tener al menos una), el hecho de que dos de ellas se diferencien en una constante hace que, para todo par de números reales $a, b \in I$, la diferencia $[F(x)]_a^b = F(b) - F(a)$ sea independiente de la elección de F . Por ello definimos

$$\int_a^b f(x) dx = [F(x)]_a^b = F(b) - F(a),$$

para toda función f definida en un intervalo abierto I en el que admita una primitiva F y todo par de números $a, b \in I$.

En realidad, podemos trabajar en un contexto ligeramente más general: Si $F, f : [a, b] \rightarrow \mathbb{R}$, diremos que F es una primitiva de f en $[a, b]$ si F es continua en $[a, b]$, derivable en $]a, b[$ y $F' = f$ en $]a, b[$.

Así, si F y G son primitivas de una misma función f en $[a, b]$, tenemos que $F - G$ es continua en $[a, b]$ y constante en $]a, b[$, lo que implica claramente que $F - G$ es constante en $[a, b]$, luego $F(b) - F(a) = G(b) - G(a)$ y la integral

$$\int_a^b f(x) dx$$

es independiente de la elección de la primitiva F (supuesto que exista alguna).

A las expresiones de este tipo se les llama *integrales definidas*, y son números reales, por oposición a las *integrales indefinidas* que hemos considerado hasta ahora, y que son conjuntos de funciones. Los números a y b se llaman *límites* o *extremos* de la integral definida.

Nota La notación $\int_a^b f(x) dx$ expresa la interpretación siguiente de la integral: si $F : [a, b] \rightarrow \mathbb{R}$ es una primitiva de f , entonces el incremento ΔF que experimenta F cuando la variable x se incrementa en una cantidad Δx es aproximadamente $\Delta F \approx f(x)\Delta x$. La igualdad de diferenciales $dF = f(x) dx$ formaliza el hecho de que la aproximación es “exacta” cuando consideramos “incrementos infinitesimales”, es decir, $f(x) dx$ representa el incremento infinitesimal que experimenta f cuando la variable x se incrementa infinitesimalmente en dx , y entonces, la igualdad $\int_a^b f(x) dx = F(b) - F(a)$ expresa que el incremento total que F experimenta cuando x varía desde a hasta b es la “suma” de los infinitos incrementos infinitesimales $f(x) dx$ que experimenta F cuando la x varía desde a hasta b en pasos infinitesimales. De hecho, el signo \int es una deformación de la S de suma (en realidad de la inicial de la palabra latina *summa*).

Más precisamente, consideremos una función continua $f : [a, b] \rightarrow \mathbb{R}$, con primitiva F , podemos dividir el intervalo $[a, b]$ en un número grande de subintervalos, no infinitesimales, pero sí muy pequeños,

$$a = x_0 < x_1 < x_2 < \cdots < x_n = b,$$

digamos de longitud $\Delta x = x_i - x_{i-1} = (b - a)/n$. Vamos a probar que

$$\int_a^b f(x) dx = \lim_n \sum_{i=1}^n f(x_{i-1})\Delta x,$$

lo cual es una formulación rigurosa de las afirmaciones precedentes sobre incrementos infinitesimales. Así, $\int dx$ puede verse como “un resumen” del proceso consistente en partir el intervalo $[a, b]$ en intervalos “casi infinitesimales”, calcular las sumas $f(x_{i-1})\Delta x$ y tomar el límite para obtener “lo que hubiéramos obtenido si hubiéramos podido tomar intervalos infinitesimales”.

Para probar la igualdad anterior usamos que, de hecho, f es uniformemente continua en $[a, b]$, por lo que, dado $\epsilon > 0$, existe un δ tal que si $|x - y| < \delta$, entonces $|f(x) - f(y)| < \epsilon/(b - a)$. Tomamos entonces cualquier n tal que $\Delta x = (b - a)/n < \delta$. Por el teorema del valor medio,

$$F(x_i) - F(x_{i-1}) = f(y_i)(x_i - x_{i-1}) = f(y_i)\Delta x,$$

para cierto $y_i \in]x_{i-1}, x_i[$, de modo que $|y_i - x_{i-1}| < \delta$. Por lo tanto

$$\left| \int_a^b f(x) dx - \sum_{i=1}^n f(x_{i-1})\Delta x \right| = \left| \sum_{i=1}^n (F(x_i) - F(x_{i-1})) - \sum_{i=1}^n f(x_{i-1})\Delta x \right|$$

$$= \left| \sum_{i=1}^n (f(y_i) - f(x_{i-1})) \Delta x \right| \leq \sum_{i=1}^n |f(y_i) - f(x_{i-1})| \Delta x < \sum_{i=1}^n \frac{\epsilon}{b-a} \frac{b-a}{n} = \epsilon.$$

Hemos probado este resultado usando que f es continua y tiene primitiva en $[a, b]$. Demostraremos más adelante que toda función continua en un intervalo cerrado tiene primitiva,⁹ por lo que la interpretación de la integral que acabamos de dar vale para funciones continuas arbitrarias. ■

El teorema siguiente recoge las propiedades básicas de las integrales definidas:

Teorema 4.46 Sean f y g dos funciones con primitiva en un intervalo abierto I y sean $a, b, c \in I$. Entonces:

$$a) \int_a^b (\alpha f(x) + \beta g(x)) dx = \alpha \int_a^b f(x) dx + \beta \int_a^b g(x) dx.$$

$$b) \int_a^c f(x) dx = \int_a^b f(x) dx + \int_b^c f(x) dx.$$

$$c) \int_a^b f(x) dx = - \int_b^a f(x) dx.$$

$$d) \int_a^a f(x) dx = 0.$$

$$e) \text{ Si } f(x) \leq g(x) \text{ en } [a, b], \text{ entonces } \int_a^b f(x) dx \leq \int_a^b g(x) dx.$$

$$f) \text{ Si } |f(x)| \text{ tiene primitiva en } I, \text{ entonces } \left| \int_a^b f(x) dx \right| \leq \int_a^b |f(x)| dx.$$

DEMOSTRACIÓN: Las cuatro primeras propiedades son inmediatas. Para probar e) observamos que si F y G son primitivas de f y g , respectivamente, entonces $(G - F)'(x) = g(x) - f(x) \geq 0$, luego la función $G - F$ es monótona creciente en $[a, b]$, luego $G(a) - F(a) \leq G(b) - F(b)$, luego concluimos que $F(b) - F(a) \leq G(b) - G(a)$.

La propiedad f) se deduce de e), pues $-|f(x)| \leq f(x) \leq |f(x)|$, luego

$$- \int |f(x)| dx \leq \int f(x) dx \leq \int |f(x)| dx,$$

de donde se sigue la relación del enunciado. ■

⁹Teorema 8.59, la prueba consiste esencialmente en construir la integral sin apoyarse en la existencia de primitiva

La traducción de la regla de derivación del producto a una regla de cálculo de primitivas requiere más atención: dadas dos funciones derivables u y v tenemos que $d(uv) = u dv + v du$, luego integrando tenemos

$$\int u dv = uv - \int v du.$$

Esta fórmula se conoce como “regla de integración por partes”. Obviamente de aquí se sigue a su vez la versión con límites:

$$\int_a^b u dv = [uv]_a^b - \int_a^b v du.$$

Ejemplo Vamos a calcular $\int x e^x dx$. Para ello llamamos $u = x$ y $dv = e^x dx$. Claramente entonces $du = dx$ y $v = \int dv = \int e^x dx = e^x$. La fórmula anterior nos da

$$\int x e^x dx = x e^x - \int e^x dx = x e^x - e^x + c.$$

Observar que hemos omitido la constante al calcular $v = e^x$. Es claro que para aplicar la regla podemos tomar como v una primitiva fija cualquiera. ■

Ejemplo Sea $I_{m,n} = \int \operatorname{sen}^m x \cos^n x dx$. Vamos a encontrar unas expresiones recurrentes que nos permitan calcular estas integrales. Integramos por partes tomando $u = \cos^{n-1} x$, $dv = \operatorname{sen}^m x \cos x dx$. De este modo:

$$\begin{aligned} I_{m,n} &= \frac{\operatorname{sen}^{m+1} x \cos^{n-1} x}{m+1} + \frac{n-1}{m+1} \int \operatorname{sen}^{m+1} x \cos^{n-2} x \operatorname{sen} x dx \\ &= \frac{\operatorname{sen}^{m+1} x \cos^{n-1} x}{m+1} + \frac{n-1}{m+1} \int \operatorname{sen}^m x \cos^{n-2} x (1 - \cos^2 x) dx \\ &= \frac{\operatorname{sen}^{m+1} x \cos^{n-1} x}{m+1} + \frac{n-1}{m+1} (I_{m,n-2} - I_{m,n}). \end{aligned}$$

Despejando llegamos a

$$I_{m,n} = \frac{\operatorname{sen}^{m+1} x \cos^{n-1} x}{m+n} + \frac{n-1}{m+n} I_{m,n-2}.$$

Similarmente se prueba

$$I_{m,n} = -\frac{\operatorname{sen}^{m-1} x \cos^{n+1} x}{m+n} + \frac{m-1}{m+n} I_{m-2,n}.$$

Estas fórmulas reducen el cálculo de cualquier integral $I_{m,n}$ al cálculo de las cuatro integrales

$$\int dx, \quad \int \operatorname{sen} x dx, \quad \int \cos x dx, \quad \int \operatorname{sen} x \cos x dx,$$

y todas ellas son inmediatas (para la última aplicamos la fórmula del seno del ángulo doble). ■

Veamos ahora en qué se traduce la regla de la cadena. Supongamos que tenemos una integral $\int u(x) dx$, que $F(x)$ es una primitiva de $u(x)$ (la que queremos calcular) y que $x = x(t)$ es una función derivable con derivada no nula (que por consiguiente tiene inversa derivable $t = t(x)$). Entonces por la regla de la cadena $F(x(t))' = F'(x(t)) x'(t) = u(x(t))x'(t)$, luego

$$\int u(x(t)) x'(t) dt = F(x(t)) + c.$$

Así pues, para calcular $\int u(x) dx$ podemos sustituir x por $x(t)$ y dx por $x'(t) dt$ y calcular una primitiva $G(t)$ de la función resultante, ésta será $F(x(t)) + c$, luego si sustituimos $G(t(x)) = F(x(t(x))) + c = F(x) + c$, obtenemos una primitiva de la función original. A esta técnica se la llama *integración por sustitución*. La versión con límites es:

$$\int_{t(a)}^{t(b)} u(x(t)) x'(t) dt = F(x(t(b))) - F(x(t(a))) = F(b) - F(a) = \int_a^b u(x) dx,$$

o sea:

$$\int_a^b u(x) dx = \int_{t(a)}^{t(b)} u(x(t)) x'(t) dt.$$

Nota Notemos que aquí tenemos otro ejemplo de “buena relación” entre el signo \int y la notación para diferenciales, pues la igualdad precedente puede interpretarse formalmente como que en el miembro izquierdo podemos sustituir $x = x(t)$ y $dx = x'(t) dt$. ■

Ejemplo Vamos a calcular $\int \sqrt{1-x^2} dx$. El integrando está definido en el intervalo $] -1, 1[$. Consideramos la función $x = \cos t$, definida y biyectiva en $]0, \pi[$. Entonces $dx = -\sin t dt$ y se cumple

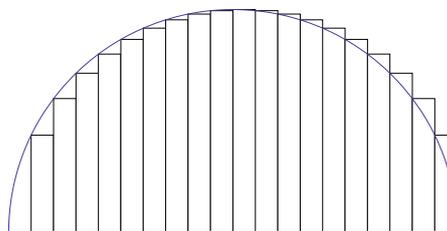
$$\begin{aligned} \int \sqrt{1-x^2} dx &= - \int \sqrt{1-\cos^2 t} \sin t dt = - \int \sin^2 t dt = - \int \frac{1-\cos 2t}{2} dt \\ &= -\frac{1}{2} \int dt + \frac{1}{4} \int 2 \cos 2t dt = -\frac{t}{2} + \frac{1}{4} \sin 2t + c = -\frac{t}{2} + \frac{1}{2} \sin t \cos t + c \\ &= -\frac{\arccos x}{2} + \frac{x\sqrt{1-x^2}}{2} + c. \end{aligned}$$

Aquí tenemos un ejemplo en el que la primitiva se extiende a una función continua en el intervalo $[-1, 1]$ (pero no en un intervalo mayor), por lo que podemos afirmar que

$$\int_{-1}^1 \sqrt{1-x^2} dx = \frac{\pi}{2} \approx 1.5708 \dots$$

La figura muestra la gráfica de la función $f(x)$ y una partición del intervalo $[-1, 1]$ en 20 subintervalos de longitud $\Delta x = 0.1$. Podemos calcular

$$\sum_{i=1}^{20} f(x_i) \Delta x \approx 1.5522 \dots,$$



con lo que obtenemos una aproximación al valor de la integral con un error inferior al 2%. Con 200 intervalos obtenemos el valor $1.5702 \dots$, con un error inferior al 0.04%.

Observemos que, gráficamente, las sumas que aproximan a la integral son las sumas de las áreas de los rectángulos que muestra la figura, y es claro que cuanto menores son los subintervalos más se parece dicha suma al área de la región situada bajo la gráfica de f . Esta interpretación geométrica de las integrales como áreas la justificaremos más adelante en un contexto más general. ■

Ejemplo Vamos a demostrar el resultado que habíamos dejado pendiente en la sección anterior, a saber, la convergencia de la serie de Taylor del arco tangente incluso en los puntos frontera ± 1 . Para ello partimos de la suma parcial de la serie geométrica de razón $-t^2$:

$$\frac{1}{1+t^2} = 1 - t^2 + t^4 - \dots + (-1)^n t^{2n} + (-1)^{n+1} \frac{t^{2n+2}}{1+t^2}.$$

Integrando ambos miembros resulta¹⁰

$$\begin{aligned} \arctan x = \int_0^x \frac{1}{1+t^2} dt &= x - \frac{x^3}{3} + \frac{x^5}{5} - \dots + (-1)^n \frac{x^{2n+1}}{2n+1} \\ &+ (-1)^{n+1} \int_0^x \frac{t^{2n+2}}{1+t^2} dt. \end{aligned}$$

El polinomio de la derecha es el polinomio de Taylor del arco tangente, luego

$$|R_{2n+1}(x)| = \left| \int_0^x \frac{t^{2n+2}}{1+t^2} dt \right| \leq \left| \int_0^x t^{2n+2} dt \right| = \frac{|x|^{2n+3}}{2n+3},$$

de donde se sigue que el resto tiende a cero incluso si $x = \pm 1$, como había que probar. ■

El teorema de Taylor 4.27 proporciona una expresión para el resto de Taylor $R_n(f)(x)$ que depende de un cierto número c cuya relación con x no es sencilla. Ahora podemos dar una expresión alternativa que no requiere seleccionar ningún número de ese modo:

¹⁰Notemos que todos los términos tienen primitiva salvo a lo sumo el último, pero precisamente por ello, el último también tiene primitiva.

Teorema 4.47 Sea $f : A \rightarrow \mathbb{R}$ una función derivable $n + 1$ veces en un intervalo abierto A y $a \in A$. Si $f^{(n+1)}(x)$ es continua en el intervalo de extremos a y x , entonces

$$R_n(f)(x) = \int_a^x \frac{f^{(n+1)}(t)}{n!} (x-t)^n dt.$$

DEMOSTRACIÓN: Notemos que para que la integral del enunciado esté bien definida hace falta que el integrando tenga primitiva. Aquí —al igual que vamos a hacer en la prueba— usamos que toda función continua en un intervalo cerrado tiene primitiva, hecho que —como ya hemos indicado— demostraremos más adelante (teorema 8.59).

Razonamos por inducción sobre n . Para $n = 0$ el teorema se reduce a que, por definición de integral,

$$f(x) = f(a) + \int_a^x f'(t) dt.$$

Supuesto cierto para n , es decir, si

$$f(x) = \sum_{k=0}^n \frac{f^{(k)}(a)}{k!} (x-a)^k + \int_a^x \frac{f^{(n+1)}(t)}{n!} (x-t)^n dt,$$

integrando por partes obtenemos que

$$\begin{aligned} f(x) &= \sum_{k=0}^n \frac{f^{(k)}(a)}{k!} (x-a)^k + \\ &\left[-\frac{f^{(n+1)}(t)}{n!} \frac{(x-t)^{n+1}}{k+1} \right]_a^x + \int_a^x \frac{(x-t)^{n+1}}{n+1} \frac{f^{(n+2)}(t)}{n!} dt, \\ &= \sum_{k=0}^{n+1} \frac{f^{(k)}(a)}{k!} (x-a)^k + \int_a^x \frac{f^{(n+2)}(t)}{(n+1)!} (x-t)^n dt. \end{aligned}$$

■

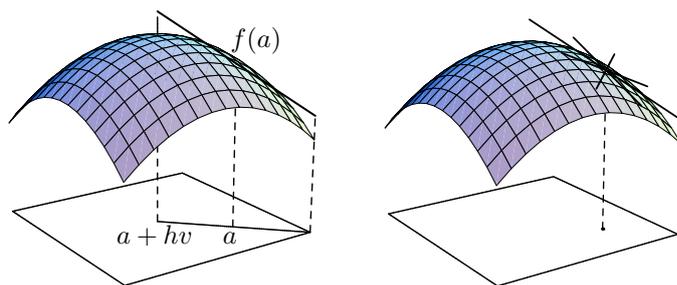
Capítulo V

Cálculo diferencial de varias variables

Este capítulo está dedicado a generalizar a funciones de varias variables las ideas que introdujimos en el capítulo anterior. Básicamente se trata de estudiar cómo varía una función de varias variables cuando incrementamos éstas infinitesimalmente. Más concretamente, estudiaremos funciones $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$, donde A es un abierto, aunque a efectos de interpretar la teoría nos centraremos de momento en el caso en que $m = 1$.

5.1 Diferenciación

Pensemos por ejemplo en una función $f : \mathbb{R}^2 \rightarrow \mathbb{R}$. Mientras una función de una variable derivable en un punto tiene asociada una única recta tangente que la aproxima, una función de dos variables tiene (o puede tener) una tangente distinta para cada dirección. La figura muestra (a la izquierda) una tangente a la gráfica de f en un punto, y a la derecha vemos varias tangentes distintas en ese mismo punto.



Intuitivamente está claro qué es la recta tangente a una superficie en un punto y en una dirección. Ahora vamos a caracterizar matemáticamente este

concepto. Tenemos un punto $a \in \mathbb{R}^n$ y una recta que pasa por a . Ésta queda determinada por un vector $v \in \mathbb{R}^n$ no nulo. Podemos suponer $\|v\| = 1$ (mientras no se indique lo contrario, todas las normas que consideraremos serán euclídeas). Los puntos de la recta son los de la forma $a + hv$, con $h \in \mathbb{R}$. Más concretamente, el punto $a + hv$ es el que se encuentra a una distancia $|h|$ de a sobre dicha recta (el signo de h distingue los dos puntos en estas condiciones).

Buscamos una recta que se parece a la gráfica de f alrededor del punto a . Si la gráfica fuera rectilínea en la dirección considerada, su pendiente vendría dada por

$$\frac{f(a + hv) - f(a)}{h},$$

para cualquier $h \neq 0$. Si no es así, entonces esta expresión se parecerá más a la pendiente que buscamos cuanto menor sea h . Ello nos lleva a la definición siguiente:

Definición 5.1 Dada una función $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ definida en un abierto, un punto $a \in A$ y un vector $v \in \mathbb{R}^n$ no nulo, llamaremos *derivada direccional* de f en a y en la dirección de v al vector

$$f'(a; v) = \lim_{h \rightarrow 0} \frac{f(a + hv) - f(a)}{h} \in \mathbb{R}^m.$$

Es fácil ver que si existe $f'(a; v)$ entonces existe $f'(a; \lambda v) = \lambda f'(a; v)$ para todo $\lambda \in \mathbb{R} \setminus \{0\}$. Por lo tanto no perdemos generalidad si suponemos $\|v\| = 1$. Si existe $f'(a; v)$, para valores pequeños de h tenemos la aproximación

$$f(a + hv) \approx f(a) + h f'(a; v),$$

con lo que la expresión $h f'(a; v)$ aproxima el incremento que experimenta $f(a)$ cuando la variable se incrementa h unidades en la dirección de v .

En el caso $m = 1$ la función $a + hv \mapsto f(a) + h f'(a; v)$ se llama *recta tangente* a la gráfica de f en a . Es claro que se trata de la recta que pretendíamos caracterizar.

Como en el caso de una variable, si una función $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ admite derivada direccional en la dirección de v y en todo punto de A , entonces tenemos definida una función

$$f'(\cdot; v) : A \rightarrow \mathbb{R}^m.$$

Respecto al cálculo de derivadas direccionales, las propiedades de los límites nos dan en primer lugar que si $f(x) = (f_1(x), \dots, f_m(x))$, entonces

$$f'(a; v) = (f'_1(a; v), \dots, f'_m(a; v)),$$

entendiendo que la derivada de f existe si y sólo si existen las derivadas de todas las funciones coordenadas f_i . Por consiguiente el cálculo de derivadas direccionales se reduce al caso de funciones $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}$. A su vez éstas se reducen al cálculo de derivadas de funciones de una variable. Efectivamente, basta considerar la función $\phi(h) = f(a + hv)$. El hecho de que A sea abierto implica claramente que ϕ está definida en un entorno de 0 y comparando las definiciones es claro que $f'(a; v) = \phi'(0)$.

Ejemplo Vamos a calcular la derivada de $f(x, y) = x^2y^2$ en el punto $(2, 1)$ y en la dirección $(-1, 1)$. Para ello consideramos

$$\phi(h) = f(2-h, 1+h) = (2-h)^2(1+h)^2.$$

Entonces $\phi'(h) = -2(2-h)(1+h)^2 + 2(2-h)^2(1+h)$ y $\phi'(0) = 4$. ■

Existen unas derivadas direccionales especialmente simples de calcular y especialmente importantes en la teoría. Se trata de las derivadas en las direcciones de la base canónica de \mathbb{R}^n .

Definición 5.2 Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ una función definida en un abierto y $a \in A$. Se define la *derivada parcial* de f respecto a la i -ésima variable en el punto a como

$$D_i f(a) = \frac{\partial f}{\partial x_i}(a) = f'(a; e_i),$$

donde $e_i = (0, \dots, 1, \dots, 0)$ es el vector con un 1 en la posición i -ésima.

Explícitamente:

$$D_i f(a) = \lim_{h \rightarrow 0} \frac{f(a_1, \dots, a_i + h, \dots, a_n) - f(a)}{h},$$

luego si h es pequeño

$$f(a_1, \dots, a_i + h, \dots, a_n) \approx f(a) + h D_i f(a).$$

En otras palabras, la expresión $h D_i f(a)$ aproxima el incremento que experimenta $f(a)$ cuando incrementamos h unidades la variable x_i .

Si f admite derivada parcial i -ésima en todos los puntos de A entonces tenemos definida la función $D_i f : A \rightarrow \mathbb{R}^m$.

El cálculo de las derivadas parciales de una función $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}$ es mucho más simple que el de las derivadas direccionales en general, pues podemos considerar la función $\phi(x_i) = f(a_1, \dots, x_i, \dots, a_n)$ y entonces es claro que $D_i f(a) = \phi'(a_i)$. Notar que si f es una función de una variable, entonces $\phi = f$, luego la derivada parcial respecto de la única variable coincide con la derivada de f en el sentido del capítulo anterior. Un poco más en general, si $f : A \subset \mathbb{R} \rightarrow \mathbb{R}^m$ llamaremos también *derivada* de f a su única derivada parcial en un punto a , y la representaremos también por $f'(a)$.

Ejemplo Las derivadas parciales de la función $f(x, y) = x^2y^3$ son

$$\frac{\partial f}{\partial x} = 2xy^3, \quad \frac{\partial f}{\partial y} = 3x^2y^2.$$

En efecto, para calcular la derivada parcial respecto de x en un punto (x_0, y_0) hay que derivar la función $x \mapsto x^2y_0^3$ en x_0 . La derivada es obviamente $2x_0y_0^3$. En la práctica podemos ahorrarnos los subíndices: para derivar x^2y^3 respecto de x basta considerar a y como una constante (la segunda coordenada del punto en que derivamos) y derivar respecto de x . Lo mismo vale para y . ■

Es claro que todas las reglas de derivación de funciones de una variable pueden ser usadas en el cálculo de derivadas parciales.

Ejercicio: Sean dos funciones derivables $f, g : I \subset \mathbb{R} \rightarrow \mathbb{R}^n$. Probar la regla de derivación $(fg)' = f'g + fg'$. Si $n = 3$ se cumple también $(f \times g)' = f' \times g + f \times g'$.

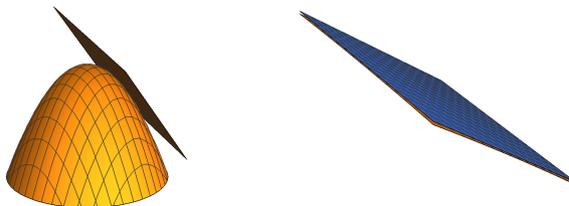
El hecho de que una función admita derivadas direccionales en un punto no puede ser equiparado a la derivabilidad de una función de una variable. Por ejemplo, la existencia de derivadas direccionales no implica siquiera la continuidad de la función en el punto. La generalización adecuada del concepto de función derivable es el concepto de “función diferenciable”, que vamos a introducir ahora.

Recordemos que si una función de una variable f tiene derivada en un punto a , entonces podemos definir su diferencial en a , que es una aplicación lineal $df(a) : \mathbb{R} \rightarrow \mathbb{R}$ con la propiedad de que $f(a) + df(a)(x - a)$ es una recta que “se confunde” con f alrededor de a . Una función de varias variables será diferenciable cuando exista una aplicación lineal que represente un papel análogo.

Ahora bien, si $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}$ es una función definida en un abierto A y $a \in A$, ¿qué tiene que cumplir una aplicación lineal $\phi : \mathbb{R}^n \rightarrow \mathbb{R}$ para que podamos decir que la función $g(x) = f(a) + \phi(x - a)$ “se confunde” con $f(x)$ alrededor de a ? Llamando por conveniencia $x = a + v$, no podemos pedir únicamente que

$$\lim_{v \rightarrow 0} (f(a + v) - f(a) - \phi(v)) = 0,$$

porque esto lo cumple cualquier función lineal ϕ . Que f y g “se confundan” alrededor de a significa que si ampliamos suficientemente la gráfica alrededor de a , la diferencia entre f y g se vuelve inapreciable, como muestra la figura siguiente:



A la izquierda vemos la función $f(x, y) = 10 - x^2 - y^2$ y su plano tangente en el punto $a = (1, 2)$. A la derecha vemos ambas gráficas en el intervalo $[0.9, 1.1] \times [1.9, 2.1]$. Observamos que la diferencia es casi inapreciable, y si ampliáramos aún más las gráficas (cubriendo, por consiguiente, un entorno menor) la diferencia se volvería completamente inapreciable.

Ampliar una figura significa multiplicar todas las coordenadas (y, por consiguiente, todas las distancias) por una misma cantidad α . En general, si ampliamos las gráficas de las dos funciones f y g alrededor del punto a de modo que un incremento pequeño v de las variables pase a verse de longitud 1, eso significa dividir todas las coordenadas (y todas las distancias) entre $\|v\|$. En la

gráfica ampliada, la distancia entre $f(a+v)$ y $g(a+v) = f(a) + \phi(v)$ sería

$$\frac{f(a+v) - f(a) - \phi(v)}{\|v\|}.$$

En estos términos, decir que la gráfica de g se confunde con la de f significa que, cuanto menor sea $\|v\|$ y, por consiguiente, mayor sea la ampliación de la gráfica, esta distancia se hace más pequeña. Con esto llegamos a la definición de función diferenciable. En la definición siguiente consideramos, más en general, una función con valores en \mathbb{R}^m :

Definición 5.3 Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ una función definida en un abierto A . Sea $a \in A$. Diremos que f es *diferenciable* en A si existe una aplicación lineal $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ tal que

$$\lim_{v \rightarrow 0} \frac{f(a+v) - f(a) - \phi(v)}{\|v\|} = 0.$$

Esta definición expresa la idea que hemos estado discutiendo. Cuando esto sucede, dado $\epsilon > 0$, existe un $\delta > 0$ tal que si $\|v\| < \delta$, entonces

$$\frac{\|f(a+v) - f(a) - \phi(v)\|}{\delta} < \frac{\|f(a+v) - f(a) - \phi(v)\|}{\|v\|} < \epsilon,$$

y esto significa que, si representamos las funciones a una escala $1 : \delta$, en todos los puntos x con $\|x - a\| < \delta$, la diferencia entre $f(x)$ y $g(x) = f(a) + \phi(x - a)$ (escalada) es menor que ϵ . Si tomamos una unidad de medida “apreciable” y un ϵ “inapreciable” respecto de ella, el resultado es que las funciones ampliadas resultan indistinguibles. Así pues, la función $g(x)$ resulta ser la generalización natural del concepto de recta tangente (es la aplicación afín que más aproxima a f alrededor de a), y ello nos lleva a conjeturar que ϕ , en caso de existir, tiene que ser única. En efecto, así es:

Teorema 5.4 Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ una función diferenciable en un punto $a \in A$. Sea $\phi : \mathbb{R}^n \rightarrow \mathbb{R}^m$ una aplicación lineal que cumpla la definición anterior. Entonces, para cada vector $v \in \mathbb{R}^n$ no nulo existe $f'(a;v)$ y además $\phi(v) = f'(a;v)$.

DEMOSTRACIÓN: Obviamente hv tiende a 0 cuando h tiende a 0. Por lo tanto, restringiendo el límite de la definición de diferenciabilidad concluimos que

$$\lim_{h \rightarrow 0} \frac{f(a+hv) - f(a) - \phi(hv)}{\|hv\|} = 0.$$

Usando que ϕ y la norma son lineales vemos que

$$\lim_{h \rightarrow 0} \frac{1}{\|v\|} \left(\frac{f(a+hv) - f(a)}{|h|} - \frac{h}{|h|} \phi(v) \right) = 0.$$

Claramente podemos eliminar el factor $1/\|v\|$ sin que el límite varíe. Ahora multiplicamos por la función $\mathbb{R} \setminus \{0\} \rightarrow \{\pm 1\}$ dada por $h \mapsto |h|/h$ y usamos que el producto de una función que tiende a 0 por otra acotada tiende a 0:

$$\lim_{h \rightarrow 0} \frac{f(a + hv) - f(a)}{h} - \phi(v) = 0,$$

Por lo tanto existe

$$f'(a; v) = \lim_{h \rightarrow 0} \frac{f(a + hv) - f(a)}{h} = \phi(v).$$

■

En particular vemos que, si f es diferenciable en a , existe una única aplicación lineal ϕ que cumple la definición de diferenciabilidad, a saber, la dada por $\phi(v) = f'(a; v)$.

Definición 5.5 Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ una función diferenciable en un punto $a \in A$. Llamaremos *diferencial* de f en a a la única aplicación lineal, representada por $df(a) : \mathbb{R}^n \rightarrow \mathbb{R}^m$, que cumple

$$\lim_{v \rightarrow 0} \frac{f(a + v) - f(a) - df(a)(v)}{\|v\|} = 0.$$

El teorema anterior afirma que para todo $v \in \mathbb{R}^n \setminus \{0\}$ se cumple

$$df(a)(v) = f'(a; v).$$

Consideremos el caso $m = 1$. Entonces, para un vector unitario v , la función

$$a + hv \mapsto f(a) + df(a)(hv) = f(a) + h df(a)(v) = f(a) + h f'(a; v)$$

es la recta tangente a f por a y en la dirección de v . Cuando h varía en \mathbb{R} y v varía entre los vectores unitarios, el punto $x = a + hv$ varía en todo \mathbb{R}^n y la aplicación $x \mapsto f(a) + df(a)(x - a)$ recorre todos los puntos de todas las rectas tangentes a f por a . Puesto que se trata de una aplicación afín, dichas tangentes forman un hiperplano.

En resumen, hemos probado que para que una aplicación (con $m = 1$) sea diferenciable en un punto a es necesario que tenga rectas tangentes por a en todas las direcciones y que además éstas formen un hiperplano. A este hiperplano se le llama *hiperplano tangente* a la gráfica de f en a . De la propia definición de diferencial (haciendo $v = x - a$) se sigue que para puntos x cercanos a a se cumple

$$f(x) \approx f(a) + df(a)(x - a).$$

El miembro derecho es precisamente el hiperplano tangente a f en a . Esta expresión indica, pues, que dicho hiperplano aproxima a f alrededor de a , pero no sólo eso (cualquier hiperplano que pasa por $(a, f(a))$ cumple eso), sino que hemos razonado que el hiperplano tangente proporciona la mejor aproximación posible, en el sentido de que se confunde con la gráfica de la función cuando la ampliamos suficientemente.

No obstante, el hecho de que todas las tangentes a la gráfica de una función por un punto formen un plano no es suficiente para garantizar que la función es diferenciable.

Ejemplo Consideremos la función

$$f(x, y) = \begin{cases} \frac{x^5 y^5}{x^{12} + y^8} & \text{si } (x, y) \neq (0, 0), \\ 0 & \text{si } (x, y) = (0, 0). \end{cases}$$

Se cumple que $f(x, y)$ es continua en $(0, 0)$.

En efecto, por la desigualdad entre la media aritmética y la geométrica, tenemos que

$$\frac{x^{12} + y^8}{13} = \frac{\overbrace{\frac{x^{12}}{5} + \dots + \frac{x^{12}}{5}}^5 + \overbrace{\frac{y^8}{8} + \dots + \frac{y^8}{8}}^8}{13} \geq \frac{\sqrt[13]{x^{60}y^{64}}}{5^5 8^8},$$

luego, llamando $C = 5^5 8^8 / 13$,

$$C(x^{12} + y^8)|x|^{5/13}|y|^{1/13} \geq |x|^5|y|^5,$$

luego

$$\left| \frac{x^5 y^5}{x^{12} + y^8} \right| \leq C \sqrt[13]{|x|^5|y|},$$

y así la continuidad es inmediata. Por otra parte, si llamamos

$$\phi(h) = f(hu, hv) = h^2 \frac{u^5 v^5}{h^3 u^{12} + v^8},$$

es claro que $\lim_{h \rightarrow 0} \phi(h) = 0$, lo que significa que existen todas las derivadas direccionales de f en $(0, 0)$ y que todas valen 0. A su vez, esto significa que la gráfica de f tiene rectas tangentes en $(0, 0)$ en todas las direcciones, y que todas ellas forman el plano $z = 0$. Sin embargo, f no es diferenciable en $(0, 0)$. La diferenciableidad equivale a que

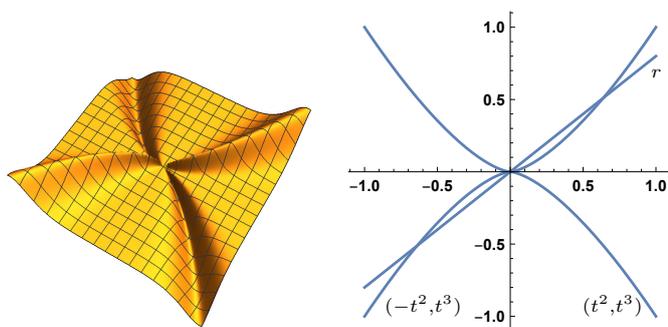
$$\lim_{(u,v) \rightarrow (0,0)} \frac{f(u, v)}{\sqrt{u^2 + v^2}} = 0,$$

pero esto no es cierto, porque si consideramos, en particular, puntos de la forma $(u, v) = (t^2, t^3)$, vemos que

$$\left| \frac{f(u, v)}{\sqrt{u^2 + v^2}} \right| = \frac{|t|^{25}/(t^{24} + t^{24})}{\sqrt{t^4 + t^6}} = \frac{1}{2|t|\sqrt{1 + t^2}},$$

luego vemos que hay puntos arbitrariamente próximos a $(0, 0)$ donde esta expresión es arbitrariamente grande.

La figura siguiente muestra la gráfica de f alrededor de $(0, 0)$. Vemos que es bastante plana salvo alrededor de las curvas (t^2, t^3) y $(-t^2, t^3)$, sobre cada una de las cuales tiene una cresta y un valle. Y por mucho que ampliemos la gráfica alrededor de $(0, 0)$ nunca dejaremos de ver esas crestas y esos valles que delatan que la superficie no es plana.



Sobre cualquier recta r que pasa por $(0, 0)$, en un entorno de $(0, 0)$, la gráfica tiende a parecerse a la propia recta r (es decir, tiene pendiente 0), y por eso las derivadas direccionales son 0, pero en cuanto la gráfica sobre r se aleja de 0 para acercarse a las curvas $(\pm t^2, t^3)$, crece o decrece para formar la cresta o el valle correspondiente. Como cualquier entorno de $(0, 0)$ contiene parte de las crestas y los valles, la función no es diferenciable en $(0, 0)$, pues las crestas y los valles distinguen a su gráfica del plano $z = 0$ en cualquier entorno de $(0, 0)$. ■

Pasamos ahora al cálculo de la diferencial de una función. Como primeras observaciones elementales notamos que si f es lineal entonces $df(a) = f$ y, si f es constante (alrededor de a), entonces $df(a) = 0$ (la aplicación nula). Ambos hechos se demuestran comprobando que con las elecciones indicadas para ϕ se cumple trivialmente la definición de diferencial.

Para una función $f(x) = (f_1(x), \dots, f_m(x))$, las propiedades de los límites nos dan que f es diferenciable en un punto a si y sólo si lo es cada función coordenada f_i , y en tal caso

$$df(a)(v) = (df_1(a)(v), \dots, df_m(a)(v)).$$

Para determinar $df(a)$ es suficiente conocer su matriz en las bases canónicas de \mathbb{R}^n y \mathbb{R}^m . Dicha matriz tiene por filas las imágenes de los vectores e_i de la base canónica.

Definición 5.6 Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ una función diferenciable en un punto $a \in A$. Llamaremos *matriz jacobiana* de f en a a la matriz $Jf(a)$ que tiene por filas a las derivadas parciales $D_i f(a)$.

Más concretamente, si $f(x) = (f_1(x), \dots, f_m(x))$, el coeficiente de la fila i , columna j de $Jf(a)$ es $D_i f_j(a)$.

Si $m = 1$, se llama *vector gradiente* de f en a al vector formado por las derivadas parciales de f en a . Se representa:

$$\nabla f(a) = (D_1 f(a), \dots, D_n f(a)).$$

Si e_i es el vector i -ésimo de la base canónica, sabemos que

$$df(a)(e_i) = f'(a; e_i) = D_i f(a),$$

luego la matriz jacobiana de f es simplemente la matriz de $df(a)$ en las bases canónicas de \mathbb{R}^n y \mathbb{R}^m . Así pues,

$$df(a)(v) = vJf(a).$$

Cuando $m = 1$, usando el producto escalar en lugar del producto de matrices tenemos también

$$df(a)(v) = \nabla f(a)v = \frac{\partial f}{\partial x_1}(a)v_1 + \dots + \frac{\partial f}{\partial x_n}(a)v_n.$$

Consideremos en particular la función polinómica x_i , es decir, la función $\mathbb{R}^n \rightarrow \mathbb{R}$ dada por $(x_1, \dots, x_n) \mapsto x_i$. Es claro que $\nabla x_i(a) = e_i$, luego $dx_i(a)(v) = v_i$. Por consiguiente, la ecuación anterior puede escribirse como

$$df(a)(v) = \frac{\partial f}{\partial x_1}(a) dx_1(a)(v) + \dots + \frac{\partial f}{\partial x_n}(a) dx_n(a)(v).$$

Como esto es válido para todo v , tenemos la ecuación funcional

$$df(a) = \frac{\partial f}{\partial x_1}(a) dx_1(a) + \dots + \frac{\partial f}{\partial x_n}(a) dx_n(a).$$

Si $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}$ es diferenciable en todos los puntos de A podemos considerar df , dx_i como funciones de A en el espacio de aplicaciones lineales de \mathbb{R}^n en \mathbb{R} y la ecuación anterior nos da

$$df = \frac{\partial f}{\partial x_1} dx_1 + \dots + \frac{\partial f}{\partial x_n} dx_n.$$

Esta fórmula expresa que si cada variable experimenta un incremento infinitesimal dx_i entonces la función f experimenta un incremento df de la forma que se indica. Como en el caso de una variable, la expresión ha de entenderse en realidad como una igualdad funcional que a cada vector de incrementos $(\Delta x_1, \dots, \Delta x_n)$ le asigna una aproximación del incremento Δf que experimenta la función.

Similarmente, en el caso $m > 1$ tenemos

$$df = (df_1, \dots, df_m) = (dx_1, \dots, dx_n) Jf.$$

Ejemplo Consideremos la función $]0, +\infty[\times]-\pi, \pi[\rightarrow \mathbb{R}^2$ dada por

$$\begin{aligned}x &= \rho \cos \theta, \\y &= \rho \operatorname{sen} \theta.\end{aligned}$$

Podríamos demostrar que es diferenciable aplicando la definición, pero más adelante será inmediato (teorema 5.12), así que vamos a aceptar que lo es y calcularemos su diferencial. Para ello calculamos la matriz jacobiana:

$$J(x, y)(\rho, \theta) = \begin{pmatrix} \frac{\partial x}{\partial \rho} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial \rho} & \frac{\partial y}{\partial \theta} \end{pmatrix} = \begin{pmatrix} \cos \theta & -\rho \operatorname{sen} \theta \\ \operatorname{sen} \theta & \rho \cos \theta \end{pmatrix}$$

Por consiguiente

$$\begin{aligned}(dx, dy) &= (d\rho, d\theta) \begin{pmatrix} \cos \theta & -\rho \operatorname{sen} \theta \\ \operatorname{sen} \theta & \rho \cos \theta \end{pmatrix} \\ &= (\cos \theta d\rho - \rho \operatorname{sen} \theta d\theta, \operatorname{sen} \theta d\rho + \rho \cos \theta d\theta),\end{aligned}$$

o más claramente:

$$\begin{aligned}dx &= \cos \theta d\rho - \rho \operatorname{sen} \theta d\theta, \\dy &= \operatorname{sen} \theta d\rho + \rho \cos \theta d\theta.\end{aligned}$$

También podríamos haber calculado dx y dy de forma independiente. ■

5.2 Propiedades de las funciones diferenciables

Vamos a estudiar las funciones diferenciables. Entre otras cosas obtendremos un criterio sencillo que justificará la diferenciable de la mayoría de funciones de interés. Comenzamos observando que en funciones de una variable la diferenciable equivale a la derivabilidad.

Teorema 5.7 Sea $f : A \subset \mathbb{R} \rightarrow \mathbb{R}^m$ una función definida en un abierto A y sea $a \in A$. Entonces f es diferenciable en a si y sólo si existe la derivada de f en a . Además en tal caso $df(a)(h) = f'(a)h$.

DEMOSTRACIÓN: Si f es derivable en a entonces existe

$$f'(a) = \lim_{h \rightarrow 0} \frac{f(a+h) - f(a)}{h} = k,$$

luego

$$\lim_{h \rightarrow 0} \frac{f(a+h) - f(a) - kh}{h} = 0,$$

y si multiplicamos por la función acotada $h/|h|$ el límite sigue siendo 0, es decir, tenemos

$$\lim_{h \rightarrow 0} \frac{f(a+h) - f(a) - kh}{|h|} = 0,$$

lo que indica que f es diferenciable y que $df(a)(h) = f'(a)h$.

El recíproco se prueba igualmente, partiendo de que $df(a)(h) = kh$ se llega a que existe $f'(a) = k$. ■

Teorema 5.8 Si $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ es diferenciable en un punto a , entonces f es continua en a .

DEMOSTRACIÓN: Tenemos que

$$\lim_{v \rightarrow 0} \frac{f(a+v) - f(a) - df(a)(v)}{\|v\|} = 0.$$

Multiplicamos por $\|v\|$, que también tiende a 0, con lo que

$$\lim_{v \rightarrow 0} f(a+v) - f(a) - df(a)(v) = 0.$$

La aplicación $df(a)$ es lineal, luego es continua, luego tiende a 0, luego

$$\lim_{v \rightarrow 0} f(a+v) - f(a) = 0,$$

y esto equivale a

$$\lim_{x \rightarrow a} f(x) = f(a),$$

luego f es continua en a . ■

Las propiedades algebraicas de la derivabilidad de funciones son válidas también para la diferenciable:

Teorema 5.9 Sean f y g funciones diferenciables en un punto a . Entonces

- $f + g$ es diferenciable en a y $d(f+g)(a) = df(a) + dg(a)$.
- Si $\alpha \in \mathbb{R}$ entonces αf es diferenciable en a y $d(\alpha f)(a) = \alpha df(a)$.
- fg es diferenciable en a y $d(fg)(a) = g(a)df(a) + f(a)dg(a)$.
- si $g(a) \neq 0$ entonces f/g es diferenciable en a y

$$d(f/g)(a) = \frac{g(a)df(a) - f(a)dg(a)}{g^2(a)}.$$

DEMOSTRACIÓN: Veamos por ejemplo la propiedad c). Llamemos

$$E(v) = \frac{f(a+v) - f(a) - df(a)(v)}{\|v\|}, \quad F(v) = \frac{g(a+v) - g(a) - dg(a)(v)}{\|v\|}.$$

Ambas funciones están definidas en un entorno de 0 y tienden a 0. Además

$$f(a+v) - f(a) = df(a)(v) + \|v\|E(v), \quad g(a+v) - g(a) = dg(a)(v) + \|v\|F(v).$$

Entonces

$$\begin{aligned} (fg)(a+v) - (fg)(a) &= f(a+v)g(a+v) - f(a)g(a+v) + f(a)g(a+v) - f(a)g(a) \\ &= (f(a+v) - f(a))g(a+v) + f(a)(g(a+v) - g(a)). \end{aligned}$$

Sustituimos $f(a+v) - f(a)$, $g(a+v)$ y $g(a+v) - g(a)$ usando las igualdades anteriores. Al operar queda

$$\begin{aligned} (fg)(a+v) - (fg)(a) - (g(a)df(a) + f(a)dg(a)) &= df(a)(v)dg(a)(v) \\ + \|v\|(df(a)(v)F(v) + E(v)g(a) + E(v)dg(a)(v) + f(a)F(v)) &+ \|v\|^2E(v)F(v). \end{aligned}$$

Hay que probar que el miembro derecho dividido entre $\|v\|$ tiende a 0. El único término para el que esto no es inmediato es

$$\frac{df(a)(v)dg(a)(v)}{\|v\|},$$

pero $\|df(a)(v)dg(a)(v)\| \leq \|df(a)\| \|dg(a)\| \|v\|^2$, luego la norma del cociente está mayorada por

$$\|df(a)\| \|dg(a)\| \|v\|,$$

que tiende a 0. ■

Veamos ahora la versión en varias variables de la regla de la cadena.

Teorema 5.10 (Regla de la cadena) *Consideremos $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ y $g : B \subset \mathbb{R}^m \rightarrow \mathbb{R}^k$ de modo que $f[A] \subset B$. Si f es diferenciable en un punto $a \in A$ y g es diferenciable en $f(a)$, entonces $f \circ g$ es diferenciable en a y*

$$d(f \circ g)(a) = df(a) \circ dg(f(a)).$$

DEMOSTRACIÓN: Llamemos $h = f \circ g$ y $b = f(a)$. Dado un $v \in \mathbb{R}^n$ tal que $a+v \in A$, tenemos

$$h(a+v) - h(a) = g(f(a+v)) - g(f(a)) = g(b+u) - g(b),$$

donde $u = f(a+v) - f(a)$. Consideremos las funciones

$$E(v) = \frac{f(a+v) - f(a) - df(a)(v)}{\|v\|}, \quad F(u) = \frac{g(b+u) - g(b) - dg(b)(u)}{\|u\|},$$

definidas en un entorno de 0 y con límite 0. Se cumple

$$\begin{aligned} h(a+v) - h(a) &= dg(b)(u) + \|u\|F(u) = dg(b)(df(a)(v) + \|v\|E(v)) + \|u\|F(u) \\ &= (df(a) \circ dg(f(a)))(v) + \|v\|dg(b)(E(v)) + \|u\|F(u). \end{aligned}$$

Basta probar que

$$\lim_{v \rightarrow 0} dg(b)(E(v)) + \frac{\|u\|}{\|v\|} F(u) = 0,$$

para lo cual basta a su vez probar que la función $\|u\|/\|v\|$ está acotada en un entorno de 0. Ahora bien,

$$\frac{\|u\|}{\|v\|} = \frac{\|df(a)(v) + \|v\|E(v)\|}{\|v\|} \leq \|df(a)\| + \|E(v)\|,$$

y, como E tiende a 0 en 0, está acotada en un entorno de 0. ■

Como consecuencia, $J(f \circ g)(a) = Jf(a)Jg(f(a))$.

Equivalentemente, supongamos que tenemos una función $z = z(y_1, \dots, y_m)$, donde a su vez $y_i = y_i(x_1, \dots, x_n)$. Entonces la regla de la cadena nos dice que, si las funciones son diferenciables,

$$\nabla z^t(x_1, \dots, x_n) = Jy(x_1, \dots, x_n)\nabla z^t(y_1, \dots, y_m),$$

luego

$$\frac{\partial z}{\partial x_i} = \frac{\partial z}{\partial y_1} \frac{\partial y_1}{\partial x_i} + \dots + \frac{\partial z}{\partial y_m} \frac{\partial y_m}{\partial x_i}.$$

Ésta es la forma explícita de la regla de la cadena.

En otros términos, si tenemos dz expresado como combinación lineal de dy_1, \dots, dy_m y cada dy_i como combinación lineal de dx_1, \dots, dx_n , es decir,

$$dz = (dy_1, \dots, dy_m)\nabla z(y_1, \dots, y_m)^t,$$

$$(dy_1, \dots, dy_m) = (dx_1, \dots, dx_n)Jy(x_1, \dots, x_n),$$

entonces, para expresar a dz como combinación lineal de dx_1, \dots, dx_n basta sustituir el segundo grupo de ecuaciones en la primera, pues así obtenemos

$$(dx_1, \dots, dx_n)Jy(x_1, \dots, x_n)\nabla z(y_1, \dots, y_m)^t = (dx_1, \dots, dx_n)\nabla z^t(x_1, \dots, x_n),$$

es decir, $dz(x_1, \dots, x_n)$.

Ejemplo Consideremos las funciones $z = \sqrt{x^2 + y^2}$, $x = \rho \cos \theta$, $y = \rho \sin \theta$. Suponemos $\rho > 0$, con lo que $(x, y) \neq (0, 0)$ y todas las funciones son diferenciables (ver el teorema 5.12, más abajo). Claramente

$$\begin{aligned} dz &= \frac{x}{\sqrt{x^2 + y^2}} dx + \frac{y}{\sqrt{x^2 + y^2}} dy, \\ dx &= \cos \theta d\rho - \rho \sin \theta d\theta, \\ dy &= \sin \theta d\rho + \rho \cos \theta d\theta. \end{aligned}$$

Entonces

$$dz = \frac{\rho \cos \theta}{\rho} (\cos \theta d\rho - \rho \operatorname{sen} \theta d\theta) + \frac{\rho \operatorname{sen} \theta}{\rho} (\operatorname{sen} \theta d\rho + \rho \cos \theta d\theta) = d\rho,$$

que es el mismo resultado que se obtiene si diferenciamos directamente la función compuesta $z(\rho, \theta) = \rho$.

Es importante comprender que el paso del primer grupo de ecuaciones a la expresión de dz en función de ρ y θ no es una mera manipulación algebraica, sino que se fundamenta en la regla de la cadena. En este caso particular, lo que dice la regla de la cadena es:

Si llamamos $z(\rho, \theta)$ a la función que resulta de sustituir x e y en $z(x, y)$ por sus valores en función de ρ y θ , entonces la diferencial de esta función es la que resulta de sustituir x, y, dx, dy en $dz(x, y)$ por sus valores en función de $\rho, \theta, d\rho, d\theta$, respectivamente.

Y esto no es evidente en absoluto. ■

Otra aplicación de la regla de la cadena nos da una regla para derivar determinantes:

Teorema 5.11 (Jacobi) *Si $A(t)$ es una matriz cuadrada cuyas coordenadas son funciones derivables de t , entonces*

$$\frac{d|A|}{dt} = \operatorname{Tr} \left(\operatorname{adj} A \frac{dA}{dt} \right).$$

En particular, si $|A|$ no se anula,

$$\frac{d|A|}{dt} = \operatorname{Tr} \left(A^{-1} \frac{dA}{dt} \right) |A|.$$

DEMOSTRACIÓN: Consideremos $|A|$ como función (polinómica) de los coeficientes a_{ij} de la matriz, de modo que

$$\frac{d|A|}{dt} = \sum_{ij} \frac{d|A|}{da_{ij}} \frac{da_{ij}}{dt}.$$

Para calcular la derivada respecto de a_{ij} desarrollamos el determinante por la fila i -ésima:

$$|A| = \sum_k a_{ik} \operatorname{adj}^t(A)_{ik},$$

así

$$\frac{d|A|}{da_{ij}} = \sum_k \frac{da_{ik}}{da_{ij}} \operatorname{adj}^t(A)_{ik} + \sum_k a_{ik} \frac{d \operatorname{adj}^t(A)_{ik}}{da_{ij}} = \operatorname{adj}^t(A)_{ij},$$

donde hemos usado que a_{ij} no aparece en $\operatorname{adj}^t(A)_{ik}$, pues es el determinante de la submatriz de A que resulta de eliminar la fila i y la columna k . Así pues:

$$\frac{d|A|}{dt} = \sum_{i,j} \operatorname{adj}^t(A)_{ij} \frac{da_{ij}}{dt} = \sum_{i,j} \operatorname{adj}(A)_{ji} \frac{da_{ij}}{dt} = \sum_j (\operatorname{adj} A \frac{dA}{dt})_j = \operatorname{Tr}(\operatorname{adj} A \frac{dA}{dt}).$$

■

El teorema siguiente es el único criterio de diferenciability que necesitaremos en la práctica:

Teorema 5.12 *Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$, donde A es un abierto en \mathbb{R}^n . Si f tiene derivadas parciales continuas en A entonces f es diferenciable en A .*

DEMOSTRACIÓN: Podemos suponer $m = 1$, pues si f tiene derivadas parciales continuas en A lo mismo vale para sus funciones coordenadas, y si éstas son diferenciables f también lo es.

Sea $a \in A$. Vamos a probar que f es diferenciable en a . Para ello basta probar que

$$\lim_{v \rightarrow 0} \frac{f(a+v) - f(a) - \sum_{i=1}^n D_i f(a) v_i}{\|v\|} = 0,$$

lo que a su vez equivale a que, dado $\epsilon > 0$, exista un $\delta > 0$ de modo que si $\|v\| < \delta$ entonces

$$\left| f(a+v) - f(a) - \sum_{i=1}^n D_i f(a) v_i \right| < \epsilon \|v\|.$$

Por la continuidad de las derivadas parciales tenemos que existe un $\delta > 0$ tal que si $\|y - a\| < \delta$ entonces $y \in A$ y $|D_i f(y) - D_i f(a)| < \epsilon/n$ para $i = 1, \dots, n$. Fijemos un v tal que $\|v\| < \delta$.

Definimos $F_i = f(a_1 + v_1, \dots, a_i + v_i, a_{i+1}, \dots, a_n)$. En particular, vemos que $f(a+v) = F_n$ y $f(a) = F_0$, luego

$$\begin{aligned} \left| f(a+v) - f(a) - \sum_{i=1}^n D_i f(a) v_i \right| &= \left| \sum_{i=1}^n (F_i - F_{i-1} - D_i f(a) v_i) \right| \\ &\leq \sum_{i=1}^n |F_i - F_{i-1} - D_i f(a) v_i|. \end{aligned}$$

Así pues, (teniendo en cuenta que $|v_i| \leq \|v\|$) basta probar que

$$|F_i - F_{i-1} - D_i f(a) v_i| < \frac{\epsilon}{n} |v_i|,$$

para $i = 1, \dots, n$. Esto resulta de aplicar el teorema del valor medio a la función de una variable dada por

$$g_i(t) = f(a_1 + v_1, \dots, a_{i-1} + v_{i-1}, a_i + tv_i, a_{i+1}, \dots, a_n).$$

Esta función está definida en un intervalo abierto que contiene al intervalo $[0, 1]$, y el hecho de que f tenga derivadas parciales implica que g_i es derivable en su dominio. En particular es derivable en $]0, 1[$ y continua en $[0, 1]$. Además, es claro que

$$g_i'(t) = D_i f(a_1 + v_1, \dots, a_{i-1} + v_{i-1}, a_i + tv_i, a_{i+1}, \dots, a_n) v_i.$$

El teorema del valor medio nos da que existe $0 < t_0 < 1$ tal que

$$F_i - F_{i-1} = g_i(1) - g_i(0) = g'_i(t_0)(1 - 0).$$

Notemos que $y = (a_1 + v_1, \dots, a_{i-1} + v_{i-1}, a_i + t_0 v_i, a_{i+1}, \dots, a_n)$ cumple

$$\|y - a\| = \|(v_1, \dots, v_{i-1}, t_0 v_i, 0, \dots, 0)\| \leq \|v\| < \delta,$$

luego

$$|F_i - F_{i-1} - D_i f(a)v_i| = |D_i f(y) - D_i f(a)| |v_i| < \frac{\epsilon}{n} |v_i|,$$

como había que probar. ■

Definición 5.13 Supongamos que una función $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ admite derivada parcial respecto a una variable x_i en todo el abierto A . Si a su vez la función $D_i f$ admite derivada parcial respecto a la variable x_j en A , a esta derivada se la representa por $D_{ij} f$. También se usa la notación

$$D_{ij} f = \frac{\partial^2 f}{\partial x_i \partial x_j}.$$

Cuando el índice es el mismo se escribe

$$D_{ii} f = \frac{\partial^2 f}{\partial x_i^2}.$$

Las funciones $D_{ij} f$ se llaman *derivadas segundas* de f . Más generalmente, la función $D_{i_1 \dots i_k} f$ será la función que resulta de derivar f respecto de i_1 , derivar dicha parcial respecto a i_2 , etc. Alternativamente,

$$\frac{\partial^k f}{\partial x_{i_1}^{k_1} \dots \partial x_{i_r}^{k_r}},$$

donde $k_1 + \dots + k_r = k$, representará la función que resulta de derivar k_1 veces f respecto a i_1 , luego k_2 veces la función resultante respecto a i_2 , etc. Estas funciones se llaman *derivadas parciales de orden k* de la función f .

Diremos que f es de clase C^k en A si existen todas sus derivadas parciales de orden k en A y todas ellas son continuas en A . En particular, las funciones de clase C^0 son las funciones continuas.

Obviamente una función es de clase C^{k+1} si y sólo si todas sus derivadas parciales son de clase C^k . El teorema anterior afirma que todas las funciones de clase C^1 son diferenciables. Si una función f es de clase C^2 , entonces sus derivadas parciales son de clase C^1 , luego son diferenciables, luego son continuas y por lo tanto f es también de clase C^1 . Por el mismo argumento se prueba en general que si $k \leq r$ entonces toda función de clase C^r es de clase C^k .

Las reglas de derivación justifican inmediatamente que la suma y el producto por un escalar de funciones de clase C^k es una función de clase C^k . El producto de funciones de clase C^k (con valores en \mathbb{R}) es de clase C^k . Lo mismo vale para el cociente si exigimos que el denominador no se anule.

Ejercicio: Probar que la composición de dos funciones de clase C^k es de clase C^k .

Ahora vamos a probar un teorema muy importante sobre derivadas sucesivas, pues nos dice que el orden de derivación no importa:

Teorema 5.14 (Teorema de Schwarz) Si $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ es una función de clase C^2 en el abierto A , entonces $D_{ij}f(a) = D_{ji}f(a)$.

DEMOSTRACIÓN: No perdemos generalidad si suponemos $m = 1$. También podemos suponer que $n = 2$, pues en general podemos trabajar con la función $F(x_i, x_j) = f(a_1, \dots, x_i, \dots, x_j, \dots, a_n)$.

Así pues, probaremos que $D_{12}f(a) = D_{21}f(a)$. Sea $a = (a_1, a_2)$. Consideremos la función

$$\Delta f(h) = f(a_1 + h, a_2 + h) - f(a_1 + h, a_2) - f(a_1, a_2 + h) + f(a_1, a_2),$$

definida en un entorno de 0. Vamos a probar que

$$D_{12}f(a) = \lim_{h \rightarrow 0^+} \frac{\Delta f(h)}{h^2}.$$

Por simetría este límite será también $D_{21}f(a)$ y el teorema estará probado.

Dado $\epsilon > 0$, la continuidad de $D_{12}f$ en a implica que existe un $h_0 > 0$ tal que si $0 < k, k' < h_0$, entonces $|D_{12}f(a_1 + k, a_2 + k') - D_{12}f(a_1, a_2)| < \epsilon$.

Podemos tomar h_0 suficientemente pequeño para que si $0 < h < h_0$ la función

$$G(t) = f(t, a_2 + h) - f(t, a_2)$$

esté definida en el intervalo $[a_1, a_1 + h]$. Por el teorema del valor medio existe un número $0 < k < h$ tal que

$$\begin{aligned} \Delta f(h) &= G(a_1 + h) - G(a_1) = hG'(a_1 + k) \\ &= h(D_1f(a_1 + k, a_2 + h) - D_1f(a_1 + k, a_2)). \end{aligned}$$

Ahora aplicamos el teorema del valor medio a la función

$$H(t) = D_1f(a_1 + k, t)$$

en el intervalo $[a_2, a_2 + h]$, que nos da un número $0 < k' < h$ tal que

$$D_1f(a_1 + k, a_2 + h) - D_1f(a_1 + k, a_2) = D_{12}f(a_1 + k, a_2 + k')h.$$

En total tenemos que

$$\Delta f(h) = D_{12}f(a_1 + k, a_2 + k')h^2,$$

luego

$$\left| \frac{\Delta f(h)}{h^2} - D_{12}f(a) \right| = |D_{12}f(a_1 + k, a_2 + k') - D_{12}f(a)| < \epsilon,$$

siempre que $0 < h < h_0$. ■

El teorema de Schwarz implica claramente que al calcular cualquier derivada parcial de orden k de una función de clase C^k es irrelevante el orden en que efectuemos las derivadas.

Ahora vamos a encaminarnos a probar el teorema de la función inversa. Esencialmente se trata de probar que las inversas de las funciones biyectivas y diferenciables son diferenciables. La situación es más complicada que en el caso de una variable, pues en el capítulo anterior vimos que toda función derivable cuya derivada no se anula es monótona, mientras que no hay ningún resultado análogo para el caso de varias variables. Por ello vamos a necesitar varios resultados previos. Entre otras cosas, nos apoyaremos en el concepto de extremo relativo y su relación con la diferenciabilidad. La situación en esto sí es análoga a la de una variable.

Definición 5.15 Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}$ y $a \in A$. Diremos que f tiene un *mínimo relativo* en a si existe un entorno G de a tal que para todo $p \in G$ se cumple $f(p) \geq f(a)$. Similarmente se define un *máximo relativo*. Diremos que a es un *extremo relativo* si es un máximo o un mínimo relativo.

Teorema 5.16 Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}$ una función diferenciable en un punto $a \in A$. Si f tiene un extremo relativo en a , entonces $df(a) = 0$.

DEMOSTRACIÓN: Sea $v \in \mathbb{R}^n$ y consideremos la función $\phi(h) = f(a + hv)$, para un cierto vector $v \in \mathbb{R}^n$, definida en un entorno de 0. Es claro que ϕ tiene un extremo relativo en 0. Sea $g(h) = a + hv$. Por la regla de la cadena

$$0 = \phi'(0) = d\phi(0)(1) = df(g(0))(dg(0)(1)) = df(a)(v). \quad \blacksquare$$

Teorema 5.17 Sea $f : \overline{B_\delta(a)} \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ una aplicación continua, diferenciable en $B_\delta(a)$, y tal que para todo $y \in B_\delta(a)$ se cumpla $|Jf(y)| \neq 0$. Supongamos además que para todo $x \in \partial B_\delta(a)$ se cumple $f(x) \neq f(a)$. Entonces $f[B_\delta(a)]$ es un entorno de $f(a)$.

DEMOSTRACIÓN: Sea $g : \partial B_\delta(a) \rightarrow \mathbb{R}$ la aplicación definida mediante $g(x) = \|f(x) - f(a)\|$. Por compacidad alcanza su mínimo en un punto x . Por hipótesis $m = g(x) > 0$. Tenemos, pues, que $g(z) \geq g(x)$, para todo z tal que $\|z - a\| = \delta$. Vamos a probar que $B_{m/2}(f(a)) \subset f[B_\delta(a)]$.

Sea $y \in B_{m/2}(f(a))$. Consideremos la función $h : \overline{B_\delta(a)} \rightarrow \mathbb{R}$ dada por $h(x) = \|f(x) - y\|$. Por compacidad alcanza su mínimo en un punto z y, puesto que $h(a) = \|f(a) - y\| < m/2$, vemos que $h(z) < m/2$.

Si $x \in \partial B_\delta(a)$, entonces

$$h(x) = \|f(x) - y\| \geq \|f(x) - f(a)\| - \|f(a) - y\| > g(x) - \frac{m}{2} \geq m - \frac{m}{2} = \frac{m}{2},$$

luego $h(x)$ no es el mínimo de h . En otros términos, $z \notin \partial B_\delta(a)$, luego $z \in B_\delta(a)$

Es claro que h^2 también alcanza su mínimo en z y

$$h^2(x) = \sum_{k=1}^n (f_k(x) - y_k)^2.$$

Por consiguiente,

$$D_j h^2(z) = \sum_{k=1}^n 2(f_k(z) - y_k) D_j f_k(z) = 0,$$

pues z es un extremo. Matricialmente tenemos $(f(z) - y)Jf(z)^t = 0$ y como el determinante de la matriz es no nulo por hipótesis, $y = f(z) \in f[B_\delta(a)]$. ■

Teorema 5.18 *Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ una aplicación inyectiva, diferenciable en el abierto A y tal que $|Jf(x)| \neq 0$ para todo $x \in A$. Entonces f es abierta, luego $f : A \rightarrow f[A]$ es un homeomorfismo.*

DEMOSTRACIÓN: Sea U un abierto en A . Veamos que $f[U]$ es abierto en \mathbb{R}^n . Tomemos un punto $a \in U$ y veamos que $f[U]$ es entorno de $f(a)$. Existe un $\delta > 0$ tal que $\overline{B_\delta(a)} \subset U$, y la restricción de f a esta bola cerrada está en las hipótesis del teorema anterior. Por consiguiente $f[B_\delta(a)] \subset f[U]$ es un entorno de $f(a)$. ■

Ahora ya podemos probar el teorema de la función inversa.

Teorema 5.19 (Teorema de la función inversa) *Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ una función inyectiva de clase C^k , con $k \geq 1$, en el abierto A y tal que se cumpla $|Jf(x)| \neq 0$ para todo $x \in A$. Entonces $B = f[A]$ es abierto y $f^{-1} : B \rightarrow A$ es de clase C^k en B .*

DEMOSTRACIÓN: Llamemos $g = f^{-1}$. Por el teorema anterior f y g son homeomorfismos. Veamos que tiene g parciales continuas. Sea e_i el i -ésimo vector de la base canónica de \mathbb{R}^n . Sea $b \in B$ y $a = g(b)$. Tomemos una bola abierta $B_\eta(a) \subset A$ y un $\alpha > 0$ suficientemente pequeño tal que el segmento de extremos b y $b + \alpha e_i$ esté contenido en $f[B_\eta(a)]$. Sea $a' = g(b + \alpha e_i)$. Entonces $a' \in B_\eta(a)$, luego el segmento de extremos a y a' está contenido en A .

Claramente $f(a') - f(a) = \alpha e_i$. Si f_j es la j -ésima función coordenada de f , tenemos

$$f_j(a') - f_j(a) = \alpha \delta_{ij},$$

donde (δ_{ij}) es la matriz identidad. Aplicamos el teorema del valor medio a la función $\phi(t) = f_j(a + t(a' - a))$, definida en $[0, 1]$, en virtud del cual existe $0 < t < 1$ tal que

$$\alpha \delta_{ij} = \phi(1) - \phi(0) = d\phi(t)(1) = df_j(a + t(a' - a))(a' - a).$$

Sea $z^j = a + t(a' - a)$. Así

$$\alpha \delta_{ij} = \sum_{k=1}^n D_k f_j(z^j)(a'_k - a_k).$$

Matricialmente tenemos

$$\alpha I = (a' - a)(D_k f_j(z^j)).$$

La función $h : A^n \rightarrow \mathbb{R}$ dada por $h(z^1, \dots, z^n) = |(D_k f_j(z^j))|$ es continua, pues las derivadas parciales son continuas y el determinante es un polinomio. Por hipótesis tenemos que $h(a, \dots, a) \neq 0$, luego existe un entorno de (a, \dots, a) en el cual h no se anula. Tomando α suficientemente pequeño podemos exigir que (a', \dots, a') esté en una bola de centro (a, \dots, a) contenida en dicho entorno, con lo que el punto (z^1, \dots, z^n) que hemos construido también está en dicho entorno, luego $|(D_k f_j(z^j))| \neq 0$ y podemos calcular la matriz inversa, cuyos coeficientes se expresan como un cociente de determinantes de matrices cuyos coeficientes son derivadas parciales. En definitiva obtenemos una expresión de la forma

$$\frac{a'_k - a_k}{\alpha} = \frac{P_k(D_k f_j(z^j))}{Q_k(D_k f_j(z^j))},$$

donde P_k y Q_k son polinomios. Por definición de a y a' tenemos

$$\frac{g_k(b + \alpha e_i) - g_k(b)}{\alpha} = \frac{P_k(D_k f_j(z^j))}{Q_k(D_k f_j(z^j))}.$$

Tomando α suficientemente pequeño podemos conseguir que $\|z^i - a\|$ se haga arbitrariamente pequeño. Por la continuidad de las derivadas parciales podemos exigir que $|D_k f_j(z^j) - D_k f_j(a)|$ se haga arbitrariamente pequeño y por la continuidad de los polinomios P_k y Q_k podemos hacer que el miembro derecho de la ecuación anterior se aproxime cuanto queramos al término correspondiente con a en lugar de los puntos z^i . En definitiva, existe

$$D_i g_k(b) = \lim_{\alpha \rightarrow 0} \frac{g_k(b + \alpha e_i) - g_k(b)}{\alpha} = \frac{P_k(D_k f_j(g(b)))}{Q_k(D_k f_j(g(b)))}.$$

Más aún, los polinomios P_k y Q_k son los que expresan las soluciones de una ecuación lineal en términos de sus coeficientes, luego no dependen de b , luego esta expresión muestra también que $D_i g_k$ es una composición de funciones continuas, luego es continua. Más en general, si f es de clase C^k , entonces g también lo es. ■

Si f está en las condiciones del teorema anterior tenemos que $f \circ f^{-1} = f^{-1} \circ f = 1$, luego por la regla de la cadena

$$df(a) \circ df^{-1}(b) = df^{-1}(b) \circ df(a) = 1,$$

luego las dos diferenciales son biyectivas y $df^{-1}(b) = df(a)^{-1}$.

Equivalentemente, si $y_i(x_1, \dots, x_n)$ es una transformación biyectiva y diferenciable con determinante jacobiano no nulo y llamamos $x_i(y_1, \dots, y_n)$ a la

función inversa, entonces el sistema de ecuaciones

$$\begin{aligned} dy_1 &= \frac{\partial y_1}{\partial x_1} dx_1 + \cdots + \frac{\partial y_1}{\partial x_n} dx_n \\ &\vdots \\ dy_n &= \frac{\partial y_n}{\partial x_1} dx_1 + \cdots + \frac{\partial y_n}{\partial x_n} dx_n \end{aligned}$$

tiene por matriz de coeficientes a la matriz jacobiana de la transformación, luego podemos despejar dx_1, \dots, dx_n en función de dy_1, \dots, dy_n y así obtenemos precisamente la diferencial de la función inversa.

Ejemplo Como ya advertíamos, al contrario de lo que sucede en una variable, el hecho de que una aplicación tenga diferencial no nula en todo punto (o incluso determinante jacobiano no nulo) no garantiza que sea biyectiva. Por ejemplo, consideremos $f(x, y) = (e^x \cos y, e^x \sin y)$ (vista como aplicación de \mathbb{C} en \mathbb{C} , se trata simplemente de la exponencial compleja). La matriz jacobiana de f es

$$\begin{pmatrix} e^x \cos y & e^x \sin y \\ -e^x \sin y & e^x \cos y \end{pmatrix}$$

y su determinante en cada punto (x, y) es $e^{2x} \neq 0$. Por otro lado es fácil ver que f no es biyectiva. ■

Lo máximo que podemos deducir del hecho de que el determinante jacobiano no se anule es que la función es localmente inyectiva. La prueba se basa en un argumento que hemos usado en la prueba del teorema de la función inversa.

Teorema 5.20 (Teorema de inyectividad local) Sea $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^n$ una función de clase C^1 en el abierto A y sea $a \in A$ tal que $|Jf(a)| \neq 0$. Entonces existe un entorno B de a tal que $|Jf(x)| \neq 0$ para todo $x \in B$ y f es inyectiva sobre B .

DEMOSTRACIÓN: La función $h : A^n \rightarrow \mathbb{R}$ dada por

$$h(z^1, \dots, z^n) = |(D_k f_j(z^j))|$$

es continua, pues las derivadas parciales son continuas y el determinante es un polinomio. Por hipótesis tenemos que $h(a, \dots, a) \neq 0$, luego existe un entorno de (a, \dots, a) en el cual h no se anula. Este entorno lo podemos tomar de la forma $B_\delta(a) \times \cdots \times B_\delta(a)$. Tomaremos $B = B_\delta(a)$.

Veamos que si $x, y \in B_\delta(a)$ y $f(x) = f(y)$ entonces $x = y$. Consideremos la función $\phi(t) = f_i(x + t(y - x))$, definida en $[0, 1]$. Claramente es derivable. Por el teorema del valor medio,

$$f_i(y) - f_i(x) = d\phi(t_i)(1) = df_i(x + t_i(y - x))(y - x) = df_i(z^i)(y - x),$$

donde z^i es un punto entre x e y , luego $z^i \in B_\delta(a)$. Por lo tanto

$$f_i(y) - f_i(x) = \sum_{k=1}^n D_k f_i(z^i)(y_k - x_k),$$

lo que matricialmente se expresa como

$$0 = f(y) - f(x) = (y - x)(D_k f_i(z^i)).$$

La matriz tiene determinante $h(z^1, \dots, z^n) \neq 0$, luego ha de ser $x = y$. ■

Para terminar generalizamos a varias variables un hecho que tenemos probado para el caso de una:

Teorema 5.21 Si $f : A \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ es una función diferenciable en un abierto conexo A y $df = 0$, entonces f es constante.

DEMOSTRACIÓN: Podemos suponer $m = 1$. Dados dos puntos $a, b \in A$, existe una poligonal contenida en A con extremos a y b . Basta probar que f toma el mismo valor en los vértices de la poligonal, luego en definitiva basta probar que si a y b son los extremos de un segmento contenido en A entonces $f(a) = f(b)$. Consideramos la función $\phi(t) = f(a + t(b - a))$ en $[0, 1]$ y le aplicamos el teorema del valor medio. Concluimos que

$$f(b) - f(a) = \phi'(t) = df(a + t(b - a))(b - a) = 0. \quad \blacksquare$$

5.3 Curvas parametrizables

Las técnicas de este capítulo nos capacitan para estudiar curvas más eficientemente que las del capítulo anterior. En efecto, allí estudiábamos curvas considerándolas como gráficas de funciones de una variable, pero esto no permite trabajar con curvas cualesquiera. Por ejemplo, una elipse no es la gráfica de ninguna función. Ahora podemos aplicar la derivabilidad al estudio de curvas en el sentido de aplicaciones $x : I \rightarrow \mathbb{R}^n$, donde I es un intervalo en \mathbb{R} . Para que una tal aplicación x pueda ser llamada “curva” razonablemente, deberemos exigir al menos que sea continua. Aquí nos centraremos en las curvas que además son derivables. Si una curva está definida en un intervalo cerrado $[a, b]$, exigiremos que sea derivable en $]a, b[$. A estas curvas las llamaremos *arcos*. Si $x : [a, b] \rightarrow \mathbb{R}^n$, los puntos $x(a)$ y $x(b)$ se llaman *extremos* del arco. Concretamente, $x(a)$ es el *extremo inicial* y $x(b)$ es el *extremo final*. La gráfica de una función continua $f : [a, b] \rightarrow \mathbb{R}$ puede identificarse con el arco $x(t) = (t, f(t))$. De este modo, el tratamiento de los arcos que estamos dando aquí generaliza al del capítulo anterior.

Según lo dicho, una curva no es un mero conjunto de puntos en \mathbb{R}^n , sino un conjunto de puntos recorridos de un modo en concreto. Si $x(t)$ es una curva, la variable t se llama *parámetro* de la misma. Conviene imaginarse a t como una variable temporal, de modo que $x(t)$ es la posición en el instante t de un punto móvil que recorre la curva. La imagen de x es la trayectoria del móvil.

Vamos a interpretar la derivabilidad de una curva $x(t)$ en un punto t . Si existe $x'(t) = v \neq 0$, entonces para valores pequeños de h tenemos

$$\frac{x(t+h) - x(t)}{h} \approx v, \quad (5.1)$$

luego $x(t+h) \approx x(t)+hv$. La curva $x(t)+hv$, cuando varía h , recorre los puntos de una recta, y estamos diciendo que para valores pequeños de h la curva x se parece a dicha recta. Así pues, la recta de dirección $x'(t)$ se confunde con la curva alrededor de $x(t)$, y por ello la llamamos *recta tangente a x en $x(t)$* . Esto interpreta la dirección de $x'(t)$. Es fácil ver que su sentido es el sentido de avance al recorrer la curva. Observemos que si $f : I \rightarrow \mathbb{R}$ es una función derivable en un punto t , entonces la tangente del arco $x(t) = (t, f(t))$ es la recta de dirección $(1, f'(t))$, luego su pendiente es $f'(t)$, luego coincide con la tangente tal y como la definimos en el capítulo anterior.

Ya tenemos interpretados la dirección y el sentido de $x'(t)$. Falta interpretar su módulo. Claramente se trata del límite del módulo de (5.1). La cantidad $\|x(t+h) - x(t)\|$ es la distancia que recorreremos en h unidades de tiempo desde el instante t , luego al dividir entre $|h|$ obtenemos la distancia media recorrida por unidad de tiempo en el intervalo de extremos t y $t+h$, es decir, la velocidad media en dicho intervalo. El límite cuando $h \rightarrow 0$ es, pues, la velocidad con que recorreremos la curva en el instante t . En realidad los físicos prefieren llamar *velocidad* a todo el vector $x'(t)$, de modo que la dirección indica hacia dónde nos movemos en el instante t y el módulo indica con qué rapidez lo hacemos.

Conviene precisar estas ideas. La *primera ley de Newton* afirma que si un cuerpo está libre de toda acción externa, permanecerá en reposo o se moverá en línea recta a velocidad constante. Todo esto se resume en que su velocidad en sentido vectorial permanece constante.

Ejemplo Consideremos la curva $x(t) = (\cos t, \sin t)$, definida en el intervalo $[0, +\infty[$. Esta curva recorre infinitas veces la circunferencia de centro $(0, 0)$ y radio 1. Su velocidad es $x'(t) = (-\sin t, \cos t)$, cuyo módulo es constante igual a 1, esto significa que recorreremos la circunferencia a la misma velocidad, digamos de un metro por segundo. Para que un objeto siga esta trayectoria es necesario que una fuerza lo obligue a mantenerse a la misma distancia del centro. Por ejemplo, el móvil podría ser un cuerpo que gira atado a una cuerda. El hecho de que $x'(2\pi) = (0, 1)$ significa que si en el instante 2π se cortara la cuerda entonces el cuerpo, libre ya de la fuerza que lo retenía, saldría despedido hacia arriba a razón de un metro por segundo.

Consideremos ahora la curva $x(t) = (\cos t^2, \sin t^2)$. Su trayectoria es la misma, pero ahora la velocidad es $x'(t) = (-2t \sin t^2, 2t \cos t^2)$, cuyo módulo es $2t$, lo que significa que ahora el móvil gira cada vez más rápido. Comienza a velocidad 0, al dar una vuelta alcanza la velocidad de 2 metros por segundo, a la segunda vuelta de 4, etc. ■

Las consideraciones anteriores muestran que la derivabilidad de una curva se traduce en la existencia de una recta tangente salvo que la derivada sea nula. La existencia de una recta tangente en $x(t)$ significa que el arco se confunde con una recta alrededor de $x(t)$, con lo que el arco no puede formar un “pico” en $x(t)$. Esto ya no es cierto si $x'(t) = 0$. Por ejemplo, la curva $x(t) = (t^3, |t^3|)$ es derivable en 0, pues su derivada por la izquierda coincide con la de $(t^3, -t^3)$ y

su derivada por la derecha coincide con la de (t^3, t^3) , y ambas son nulas, pero su gráfica es la misma que la de $(t, |t|)$, es decir, la de la función $|x|$, que tiene un pico en 0. Esto nos lleva a descartar las curvas con derivada nula:

Definición 5.22 Una *curva parametrizada regular* $x : I \rightarrow \mathbb{R}^n$ es una aplicación definida sobre un intervalo abierto $I \subset \mathbb{R}$ derivable y con derivada no nula en ningún punto.

El vector

$$T(t) = \frac{x'(t)}{\|x'(t)\|}$$

se llama *vector tangente* a x en el punto $x(t)$. La recta que pasa por $x(t)$ con dirección $T(t)$ se llama *recta tangente* a x por $x(t)$.

Hemos visto un ejemplo de una misma trayectoria recorrida a velocidades distintas. Desde un punto de vista geométrico, lo que importa de una curva es su forma, y no la velocidad con la que se recorre. Vamos a explicitar esta distinción.

Dado un arco parametrizado $x : [a, b] \rightarrow \mathbb{R}^n$, un *cambio de parámetro* es una aplicación $t : [u, v] \rightarrow [a, b]$ que se extiende a un intervalo abierto mayor donde es derivable y la derivada no se anula. Por consiguiente t es biyectiva y t' tiene signo constante. Diremos que t es un cambio de parámetro *directo* o *inverso* según si $t' > 0$ o $t' < 0$. El arco parametrizado $y(s) = x(t(s))$ se llama *reparametrización* de x mediante el cambio de parámetro t .

Diremos que dos arcos parametrizados regulares x e y son (*estrictamente*) *equivalentes* si existe un cambio de parámetro (directo) que transforma uno en otro.

Es claro que la identidad es un cambio de parámetro directo, la inversa de un cambio de parámetro (directo) es un cambio de parámetro (directo) y la composición de dos cambios de parámetro (directos) es un cambio de parámetro (directo). De aquí se sigue que la equivalencia y la equivalencia estricta son relaciones de equivalencia entre los arcos parametrizados.

Llamaremos *arcos regulares* a las clases de equivalencia estricta de arcos parametrizados regulares, de modo que dos elementos de la misma clase se considerarán dos *parametrizaciones* de un mismo arco.

Todas las parametrizaciones de un arco tienen la misma imagen, a la que podemos llamar *imagen* del arco. Los cambios de parámetro directos son crecientes, por lo que dos parametrizaciones de un mismo arco tienen los mismos extremos, a los que podemos llamar *extremos* del arco.

Consideremos dos parametrizaciones $x : [a, b] \rightarrow \mathbb{R}^n$, $y : [c, d] \rightarrow \mathbb{R}^n$ de un mismo arco. Entonces $y(s) = x(t(s))$ para un cierto cambio de parámetro directo t . Consideremos ahora dos cambios de parámetro inversos

$$u : [a', b'] \rightarrow [a, b], \quad v : [c', d'] \rightarrow [c, d]$$

y las reparametrizaciones $x(u(r))$, $y(v(r))$. Entonces $y(v(r)) = x(t(v(r))) = x(u(u^{-1}(t(v(r))))$ y $v \circ t \circ u^{-1}$ es un cambio de parámetro directo, luego las

dos reparametrizaciones corresponden a un mismo arco. Por lo tanto, podemos definir el *arco inverso* de un arco x al arco resultante de componer con un cambio de parámetro inverso cualquiera de las parametrizaciones de x . Lo representaremos por $-x$. Es claro que x y $-x$ tienen la misma imagen, pero el extremo inicial de x es el extremo final de $-x$ y viceversa.

De este modo, la noción de arco como clase estricta de arcos parametrizados recoge el concepto geométrico de arco regular independiente del modo en que se recorre, pero conservando el sentido del recorrido. Si consideramos clases no estrictas identificamos cada arco con su inverso, y con ello hacemos abstracción incluso del sentido en que lo recorremos. Todos estos conceptos se aplican igualmente a curvas cualesquiera.

Longitud de un arco En el capítulo anterior hemos “calculado–definido” la longitud de una circunferencia de radio r , que ha resultado ser igual a $2\pi r$. Vamos a dar otro argumento que conduce al mismo resultado, pero que motiva una definición aplicable a una curva cualquiera, no necesariamente una circunferencia:

Consideremos el arco $x(t) = (r \cos t, r \sin t)$, con $r > 0$. Su derivada tiene módulo r , lo que se interpreta como que el arco recorre la circunferencia de centro $(0, 0)$ y radio r a una velocidad constante de r unidades de longitud por unidad de tiempo. Puesto que recorre la circunferencia completa en 2π unidades de tiempo, concluimos que el espacio que recorre, es decir, la longitud de la circunferencia, es $2\pi r$. Más en general, mediante esta parametrización recorreremos un arco de α radianes en α unidades de tiempo, luego la longitud de un arco de α radianes es αr .

La ventaja de este argumento es que, como hemos indicado, se puede generalizar:

Sea $x : [a, b] \rightarrow \mathbb{R}^n$ un arco y vamos a definir la función $s : [a, b] \rightarrow \mathbb{R}$ tal que $s(t)$ es la longitud de arco entre $x(a)$ y $x(t)$. Obviamente ha de cumplir $s(a) = 0$. Supongamos que es derivable y vamos a calcular su derivada en un punto t . El arco x se confunde con una recta alrededor de $x(t)$. Esto significa que si h es suficientemente pequeño el arco entre $x(t)$ y $x(t+h)$ es indistinguible del segmento que une ambos puntos, luego sus pendientes serán también indistinguibles, es decir:

$$\frac{s(t+h) - s(t)}{h} \approx \left\| \frac{x(t+h) - x(t)}{h} \right\|.$$

Estos dos cocientes se parecerán más cuanto menor sea h , lo que se traduce en que los límites cuando $h \rightarrow 0$ han de coincidir. Así:

$$\frac{ds}{dt} = \|x'(t)\|, \quad (5.2)$$

luego

$$s(t) = \int_a^t \|x'(u)\| du.$$

En particular, la longitud del arco completo será

$$L(x) = \int_a^b \|x'(t)\| dt.$$

Definición 5.23 Diremos que un arco parametrizado regular $x : [a, b] \rightarrow \mathbb{R}^n$ es *rectificable* si la función $\|x'(t)\|$ tiene primitiva en $[a, b]$ (es decir, si tiene primitiva en el intervalo abierto y ésta se extiende continuamente al intervalo cerrado), y entonces llamaremos *longitud* de x al número real

$$L(x) = \int_a^b \|x'(t)\| dt.$$

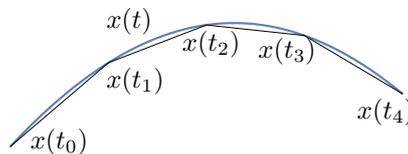
Notemos que, como el integrando es positivo, su primitiva ha de ser estrictamente creciente, luego $L(x) > 0$.

Conviene conocer la caracterización siguiente:

Teorema 5.24 Sea $x : [a, b] \rightarrow \mathbb{R}^n$ una función de clase C^1 en $[a, b]$ (en el sentido de que su derivada se extiende a una función continua en $[a, b]$). Para cada natural $N > 0$ consideramos la partición $a = t_0 < t_1 < \dots < t_N = b$ en subintervalos de longitud $\Delta t = (b - a)/N$. Entonces

$$\int_a^b \|x'(t)\| dt = \lim_N \sum_{i=1}^N \|x(t_i) - x(t_{i-1})\|.$$

Esto significa que la longitud de una curva se puede aproximar por la suma de las longitudes de los segmentos de una poligonal que una varios puntos de la curva, y la aproximación es mejor cuanto más próximos estén los puntos.



DEMOSTRACIÓN: Por el teorema del valor medio, existen $t_{i-1} < t_{ij} < t_i$ tales que

$$x_j(t_i) - x_j(t_{i-1}) = x'_j(t_{ij})\Delta t,$$

luego

$$\|x(t_i) - x(t_{i-1})\| = \sqrt{\sum_{j=1}^n (x_j(t_i) - x_j(t_{i-1}))^2} = \sqrt{\sum_{j=1}^n x'_j(t_{ij})^2} \Delta t = f(\bar{t}_i) \Delta t,$$

donde $\bar{t}_i = (t_{i1}, \dots, t_{in})$ y $f(u) = \sqrt{\sum_{j=1}^n x'_j(u_j)^2}$.

La función f es uniformemente continua en el compacto $[a, b]^n$ luego, dado $\epsilon > 0$, existe un $\delta > 0$ tal que si $|u_j - v_j| < \delta$ para todo $j = 1, \dots, n$, entonces

$|f(u) - f(v)| < \epsilon/2(b-a)$. Así, si tomamos N suficientemente grande como para que $\Delta t = (b-a)/N < \delta$, tenemos que $|t_{ij} - t_{i-1}| < \delta$, luego

$$|f(\bar{t}_i) - \|x'(t_{i-1})\|| = |f(\bar{t}_i) - f(t_{i-1}, \dots, t_{i-1})| < \frac{\epsilon}{2(b-a)}.$$

Por consiguiente,

$$\|x'(t_{i-1})\|\Delta t - \frac{\epsilon}{2(b-a)}\Delta t < \|x(t_i) - x(t_{i-1})\| < \|x'(t_{i-1})\|\Delta t + \frac{\epsilon}{2(b-a)}\Delta t,$$

luego

$$\sum_{i=1}^N \|x'(t_{i-1})\|\Delta t - \frac{\epsilon}{2} < \sum_{i=1}^N \|x(t_i) - x(t_{i-1})\| < \sum_{i=1}^N \|x'(t_{i-1})\|\Delta t + \frac{\epsilon}{2}$$

o, lo que es lo mismo,

$$\left| \sum_{i=1}^N \|x'(t_{i-1})\|\Delta t - \sum_{i=1}^N \|x(t_i) - x(t_{i-1})\| \right| < \frac{\epsilon}{2}.$$

Por otra parte (por la nota de la página 188), también podemos tomar N suficientemente grande como para que

$$\left| \int_a^b \|x'(t)\| dt - \sum_{i=1}^N \|x'(t_{i-1})\|\Delta t \right| < \frac{\epsilon}{2},$$

con lo que llegamos a que, si N es suficientemente grande,

$$\left| \int_a^b \|x'(t)\| dt - \sum_{i=1}^N \|x(t_i) - x(t_{i-1})\| \right| < \epsilon. \quad \blacksquare$$

Por ejemplo, del teorema anterior se deduce que $l(x) \geq \|x(b) - x(a)\|$, es decir, que la longitud de un arco es siempre mayor o igual que la distancia entre sus extremos o, en otros términos, que el camino mas corto que une dos puntos es el segmento que los tiene por extremos. En efecto, basta usar que

$$\|x(b) - x(a)\| = \left\| \sum_{i=1}^N x(t_i) - x(t_{i-1}) \right\| \leq \sum_{i=1}^N \|x(t_i) - x(t_{i-1})\|$$

y hacer que N tienda a infinito.

La longitud de un arco parametrizado coincide con la de cualquiera de sus reparametrizaciones. En efecto, si $y(s) = x(t(s))$, entonces $y'(s) = x'(t(s))t'(s)$, luego

$$L(x) = \int_a^b \|x'(t)\| dt = \int_{s(a)}^{s(b)} \|x'(t(s))\|t'(s) ds = \int_{s(a)}^{s(b)} \|y'(s)\| ds = L(y).$$

Ejercicio: Comprobar que la longitud de un arco coincide con la de su opuesto, y que la longitud es invariante por isometrías.

Ahora es inmediato comprobar que la longitud de un arco de circunferencia de radio r y amplitud α es αr . En particular la longitud de una circunferencia de radio r es $2\pi r$, tal y como ya sabíamos.

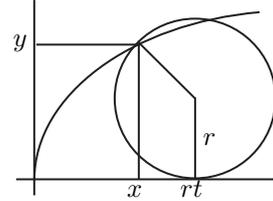
Sea $x : [a, b] \rightarrow \mathbb{R}^n$ un arco rectificable y sea $s(t)$ la longitud de arco entre $x(a)$ y $x(t)$. Hemos visto que se trata de una función derivable y que $s'(t) = \|x'(t)\|$. Por lo tanto s es creciente y biyectiva. Su inversa $t : [0, L] \rightarrow [a, b]$ es un cambio de parámetro, la función $x(s) = x(t(s))$ es una reparametrización del arco y

$$x'(s) = x'(t)t'(s) = \frac{x'(t)}{\|x'(t)\|},$$

luego $\|x'(s)\| = 1$ y así $x'(s) = T(s)$. Más concretamente, $x(s)$ se caracteriza por que la longitud de arco entre $x(0)$ y $x(s)$ es s . A esta parametrización del arco la llamaremos *parametrización natural*. También diremos entonces que x está *parametrizado por el arco*. Desde un punto de vista cinemático, la parametrización natural se interpreta como la que recorre el arco a velocidad constante igual a 1 (constante en módulo).

Ejemplo La trayectoria de un clavo de una rueda que gira se conoce con el nombre de *cicloide*. Vamos a calcular la longitud de una vuelta de cicloide.

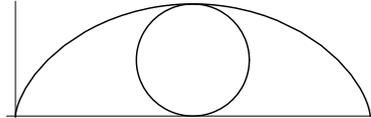
Primeramente necesitamos encontrar una parametrización de la curva. Supongamos que la rueda gira sobre el eje $y = 0$ y que el clavo parte de la posición $(0, 0)$. Si llamamos r al radio de la circunferencia, vemos que cuando la rueda ha girado t radianes su centro se encuentra en el punto (rt, r) , luego el clavo se encuentra en $x(t) = (rt - r \sin t, r - r \cos t)$.



Por consiguiente $x'(t) = r(1 - \cos t, \sin t)$,

$$\|x'(t)\| = r\sqrt{2 - 2\cos t} = 2r \sin \frac{t}{2}.$$

Vemos que los múltiplos de 2π son puntos singulares de la cicloide. Corresponden a los momentos en que el clavo toca el suelo. Entonces se para y cambia de sentido.

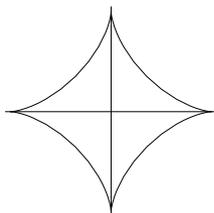


Es fácil calcular

$$L = \int_0^{2\pi} 2r \sin \frac{t}{2} dt = 4r \left[-\cos \frac{t}{2} \right]_0^{2\pi} = 8r.$$

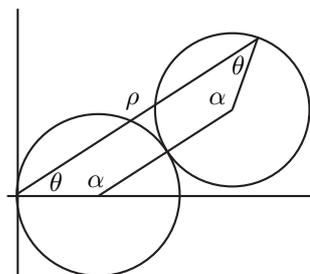
■

Ejercicio: Calcular la longitud de la *astroide*, dada por $x(t) = (a \cos^3 t, a \sin^3 t)$.



Ejemplo La trayectoria de un clavo de una rueda que gira sobre una circunferencia se llama *epicloide*. El caso más simple lo tenemos cuando ambas circunferencias tienen el mismo radio. La curva se llama entonces *cardioide*, porque su forma recuerda a un corazón.

Supongamos que la circunferencia fija tiene centro en $(a/4, 0)$ y radio $a/4$ y que el clavo parte de la posición $(0, 0)$. Cuando la rueda ha girado α radianes la situación es la que indica la figura. El cuadrilátero tiene dos ángulos y dos lados iguales, por lo que los otros dos ángulos también tienen la misma amplitud θ . Es claro entonces que

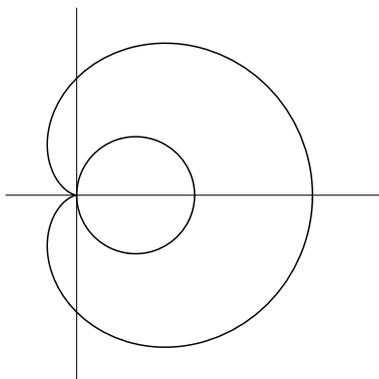


$$\rho = \frac{a}{2} + 2\frac{a}{4} \cos \theta,$$

luego

$$\rho = \frac{a}{2}(1 + \cos \theta),$$

donde $\theta \in]-\pi, \pi[$. Ésta es la ecuación de la cardioides en coordenadas polares.



Vamos a ver en general cuál es la expresión de la longitud de una curva parametrizada en coordenadas polares $(\rho(t), \theta(t))$. Notemos que (5.2), para el caso de dos variables puede escribirse también como

$$ds^2 = dx^2 + dy^2.$$

Se dice que ésta es la expresión del elemento de longitud (es decir, de una longitud infinitesimal) en coordenadas cartesianas. Se trata de la versión infinitesimal del teorema de pitágoras. Si diferenciamos las relaciones $x = \rho \cos \theta$,

$y = \rho \operatorname{sen} \theta$ obtenemos

$$dx = \cos \theta d\rho - \rho \operatorname{sen} \theta d\theta, \quad dy = \operatorname{sen} \theta d\rho + \rho \cos \theta d\theta.$$

Sustituyendo queda

$$ds^2 = d\rho^2 + \rho^2 d\theta^2. \quad (5.3)$$

Ésta es la expresión del elemento de longitud de un arco en coordenadas polares. Aplicado a la cardioide resulta

$$ds^2 = \frac{a^2}{2}(1 + \cos \theta) d\theta^2,$$

de donde

$$ds = a \cos \frac{\theta}{2} d\theta.$$

Así pues, la longitud de la cardioide es

$$L = \int_{-\pi}^{\pi} a \cos \frac{\theta}{2} d\theta = 2a \left[\operatorname{sen} \frac{\theta}{2} \right]_{-\pi}^{\pi} = 4a.$$

■

Ejemplo Consideremos un cuerpo puntual situado en $(l, 0)$ atado a una cuerda de longitud l con su otro extremo en $(0, 0)$. Estiramos de la cuerda de modo que su extremo suba por el eje Y . La trayectoria del cuerpo arrastrado por la cuerda recibe el nombre de *tractriz*. Vamos a obtener una parametrización de la tractriz. Un cuerpo estirado por una cuerda se mueve en la dirección de la cuerda, luego ésta ha de ser tangente a la trayectoria. Si llamamos $y = f(x)$ a la tractriz, definida para $0 < x < l$, entonces su recta tangente es

$$Y = f(x) + f'(x)(X - x).$$

Cortará al eje Y en el punto $(0, f(x) - f'(x)x)$. Éste es el punto donde está el extremo de la cuerda cuando el otro extremo está en $(x, f(x))$. La distancia entre ambos debe ser, pues, igual a l . Por consiguiente:

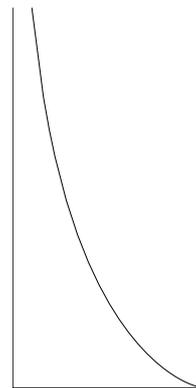
$$x^2 + f'(x)^2 x^2 = l^2.$$

Despejando obtenemos

$$dy = -\sqrt{\frac{l^2}{x^2} - 1} dx.$$

El signo negativo se debe a que, tal y como hemos planteado el problema, la función y ha de ser decreciente. El cambio $x = l \operatorname{sen} \theta$ biyecta los números $0 < x \leq l$ con los números $\pi/2 \leq \theta < \pi$. La igualdad anterior se transforma en

$$dy = -l \sqrt{\frac{1}{\operatorname{sen}^2 \theta} - 1} \cos \theta d\theta = l \frac{\cos^2 \theta}{\operatorname{sen} \theta} d\theta = \frac{l d\theta}{\operatorname{sen} \theta} - l \operatorname{sen} \theta d\theta.$$



La posición inicial corresponde a $x = l$, luego a $\theta = \pi/2$, es decir, $y(\pi/2) = 0$, luego

$$y(\theta) = l \int_{\pi/2}^{\theta} \frac{dt}{\operatorname{sen} t} - l \int_{\pi/2}^{\theta} \operatorname{sen} t \, dt.$$

Existen reglas de integración que permiten calcular metódicamente la primera primitiva. Como no nos hemos ocupado de ellas nos limitaremos a dar el resultado. El lector puede comprobar sin dificultad que es correcto derivando la solución que presentamos.

$$y(\theta) = l \left[\log \tan \frac{t}{2} \right]_{\pi/2}^{\theta} + l [\cos t]_{\pi/2}^{\theta} = l \log \tan \frac{\theta}{2} + l \cos \theta.$$

Así pues, la tractriz viene dada por

$$T(\theta) = (l \operatorname{sen} \theta, l \log \tan \frac{\theta}{2} + l \cos \theta), \quad \theta \in [\pi/2, \pi[.$$

Su derivada es

$$T'(\theta) = \left(l \cos \theta, l \frac{\cos^2 \theta}{\operatorname{sen} \theta} \right), \quad \|T'(\theta)\| = -l \frac{\cos \theta}{\operatorname{sen} \theta}.$$

La tractriz es regular en $] \pi/2, \pi[$. La longitud de un arco de tractriz es

$$s(\theta) = -l \int_{\pi/2}^{\theta} \frac{\cos t}{\operatorname{sen} t} \, dt = -l [\log \operatorname{sen} t]_{\pi/2}^{\theta} = -l \log \operatorname{sen} \theta. \quad \blacksquare$$

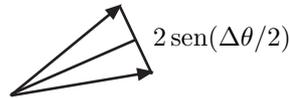
Vector normal y curvatura Sea $x(s)$ un arco parametrizado por la longitud de arco. Supongamos que admite derivada segunda. Derivando la igualdad $x'(s)x'(s) = 1$ obtenemos que $2x''(s)x'(s) = 0$, luego $x''(s) \perp x'(s)$. Supuesto que $x''(s) \neq 0$, podemos definir el *vector normal* del arco como

$$N(s) = \frac{x''(s)}{\|x''(s)\|},$$

y la *curvatura* del arco como $\kappa(s) = \|x''(s)\|$.

Para interpretar la curvatura llamemos $\Delta\theta$ al ángulo entre los vectores $x'(s)$ y $x'(s + \Delta s)$, donde $\Delta\theta$ es una función de Δs . Puesto que x' tiene módulo constante igual a 1, la trigonometría nos da que

$$\|x'(s + \Delta s) - x'(s)\| = 2 \operatorname{sen} \frac{\Delta\theta}{2}.$$



Por consiguiente,

$$\frac{\|x'(s + \Delta s) - x'(s)\|}{|\Delta s|} = \frac{\operatorname{sen} \frac{\Delta\theta}{2}}{\Delta\theta/2} \frac{\Delta\theta}{|\Delta s|}.$$

Es claro que $\Delta\theta \rightarrow 0$ cuando $\Delta s \rightarrow 0$, luego

$$\kappa(s) = \lim_{\Delta s \rightarrow 0} \frac{\Delta\theta}{|\Delta s|}.$$

Así pues, la curvatura mide la variación del ángulo del vector tangente por unidad de arco recorrido. Es claro que cuanto mayor sea la curvatura “más curvado” estará el arco.

Ejemplo La parametrización natural de una circunferencia de radio r es

$$x(s) = (r \cos(s/r), r \sin(s/r)).$$

El vector tangente es, por lo tanto, $x'(s) = (-\sin(s/r), \cos(s/r))$. De aquí obtenemos

$$x''(s) = \left(-\frac{1}{r} \cos \frac{s}{r}, -\frac{1}{r} \sin \frac{s}{r} \right),$$

con lo que el vector normal es $N(s) = -(\cos(s/r), \sin(s/r))$ y la curvatura es $\kappa(s) = 1/r$. Intuitivamente es claro que una circunferencia está menos curvada cuanto mayor es su radio, tal y como se pone de manifiesto en la fórmula que hemos obtenido. ■

Dada una curva x de curvatura no nula en un punto dado $x(s)$, la circunferencia de radio $r = 1/\kappa(s)$ y centro $x(s) - \kappa N(s)$ pasa por $x(s)$ y tiene el mismo vector tangente, el mismo vector normal y la misma curvatura en este punto. Esto hace que sea la circunferencia que más se parece a x en un entorno de $x(s)$ y se la llama *circunferencia oscultriz* a x por $x(s)$. El radio $r(s) = 1/\kappa(s)$ se llama *radio de curvatura* de x en $x(s)$.

Veamos ahora fórmulas explícitas para calcular el vector normal y la curvatura cuando la parametrización no es la natural. Conviene usar el lenguaje de la cinemática. Supongamos que $x(t)$ representa la posición de un móvil puntual en función del tiempo. Llamaremos $x(s)$ a la parametrización natural. Entonces la velocidad del móvil es $V = x'(t)$. Si llamamos $v = \|V\|$, entonces sabemos que $v = s'(t)$, con lo que v es la *velocidad sobre la trayectoria*, que mide la distancia recorrida por unidad de tiempo. Además $V = x'(t) = x'(s)s'(t)$, es decir,

$$V = vT.$$

Se define el vector *aceleración* como $A = V' = X''(t)$. Derivando en la relación anterior tenemos

$$A = v'(t)T(t) + v(t)T'(t).$$

Llamaremos $a(t) = v'(t)$, que es la *aceleración sobre la trayectoria*, es decir la variación de la velocidad sobre la trayectoria por unidad de tiempo. Aplicando la regla de la cadena a $T(t) = T(s(t))$ obtenemos $T'(t) = T'(s)s'(t) = \kappa vN$, luego

$$A = aT + \kappa v^2 N = aT + \frac{v^2}{r} N.$$

Vemos, pues, que la aceleración se descompone de forma natural en una componente *tangencial*, que mide la variación del módulo de la velocidad, y una componente *normal*, que determina la curvatura.

Ahora multiplicamos esta igualdad por sí misma: $\|A\|^2 = a^2 + \kappa^2 v^4$, de donde

$$\kappa^2 = \frac{\|A\|^2 - a^2}{v^4}.$$

Ahora bien,

$$a = v' = \|V\|' = \frac{VV'}{\|V\|} = \frac{VV'}{v},$$

luego

$$\kappa^2 = \frac{\|A\|^2 v^2 - (VV')^2}{v^6} = \frac{\|V \times A\|^2}{v^6},$$

lo que nos da

$$\kappa = \frac{\|V \times A\|}{v^3}.$$

Las dos últimas igualdades suponen que la imagen de x está en \mathbb{R}^3 . En el lenguaje geométrico hemos obtenido las fórmulas siguientes:

$$T = \frac{x'}{\|x'\|}, \quad N = \frac{\|x'\|^2 x'' - (x'x'')x'}{\|x'\| \|x' \times x''\|}, \quad \kappa = \frac{\|x' \times x''\|}{\|x'\|^3}. \quad (5.4)$$

En el caso de curvas planas es conveniente modificar como sigue el vector normal y la curvatura. Fijada una orientación en el plano, definimos el vector normal de una curva $x(s)$ parametrizada por el arco como el vector unitario $N(s)$ que es perpendicular a $T(s)$ y de modo que la base $(T(s), N(s))$ esté orientada positivamente. Esta definición puede diferir de la anterior en cuanto al signo de $N(s)$, por lo que redefinimos la curvatura de modo que $x''(s) = \kappa(s)N(s)$. Notemos que ahora el vector normal está definido incluso en los puntos donde la curvatura es nula.

Con el convenio usual de orientación, el vector N apunta hacia la izquierda si miramos en el sentido de T . La curvatura es positiva si cuando la curva avanza se desvía hacia la izquierda y negativa si se desvía hacia la derecha (o nula si no se desvía). Diremos que la curva gira en sentido positivo o en sentido negativo según el signo de su curvatura. El sentido de giro positivo es el contrario a las agujas del reloj. De este modo, la orientación distingue los dos sentidos de giro.

Vector binormal y torsión La teoría de curvas en \mathbb{R}^3 se completa con la introducción del vector binormal y la torsión. El vector *binormal* de una curva $x(s)$ tres veces derivable en un punto de curvatura no nula es $B(s) = T(s) \times N(s)$. La base formada por los vectores $(T(s), N(s), B(s))$ se conoce como *triedro de Frenet*. Definimos la *torsión* de x en cada punto como $\tau(s) = -N'(s)B(s)$.

Para interpretar la torsión empezaremos por determinar N' . Digamos que

$$N' = aT + bN + cB. \quad (5.5)$$

Multiplicando por T obtenemos $a = N'T$, y como $TN = 0$, derivando resulta $T'N + TN' = 0$, luego $a = -T'N = -\kappa NN = -\kappa$.

Si multiplicamos (5.5) por N resulta $b = N'N$, pero al derivar en $NN = 1$ resulta $2NN' = 0$, luego $b = 0$. Por último, es claro que $c = N'B = -\tau$. Por consiguiente

$$N' = -\kappa T - \tau B.$$

La misma técnica nos da una expresión para B' . Sea $B' = aT + bN + cB$. Multiplicando por T obtenemos $a = TB'$, pero de $BT = 0$ se concluye que $TB' = -T'B = -\kappa NB = 0$. Similarmente $b = NB' = -N'B = \tau$. Al multiplicar por B llegamos a $c = B'B = 0$, pues $BB = 1$.

Tenemos así las llamadas *fórmulas de Frenet*:

$$\begin{aligned} T' &= \kappa N, \\ N' &= -\kappa T - \tau B, \\ B' &= \tau N. \end{aligned}$$

Vemos que si $\tau = 0$ en todo punto entonces B es constante, luego $(xB)' = x'B = TB = 0$ implica que xB es constante, luego x está contenido en un plano perpendicular a B , luego la curva es plana. El recíproco es claro. Así pues, las curvas sin torsión son exactamente las curvas planas. En general, el plano $x(s) + \langle T(s), N(s) \rangle$ recibe el nombre de *plano osculante* a la curva. Si la curva es plana, su plano osculante es el mismo en todo punto, y la curva está contenida en él. En caso contrario, puede probarse que el plano osculante en un punto es el plano más próximo a la curva en un entorno del punto. La tercera fórmula de Frenet muestra que la torsión mide la rapidez con que varía B o, lo que es lo mismo, la rapidez con la que varía el plano osculante.

A continuación derivamos una fórmula explícita para la torsión de una curva. Si x está parametrizada por la longitud de arco, entonces

$$N = \frac{x''}{\kappa}, \quad N' = \frac{x''' \kappa - x'' \kappa'}{\kappa^2}.$$

luego,

$$\tau = -BN' = (T \times N)N' = -\frac{1}{\kappa} (x' \times x'')N' = -\frac{(x' \times x'')x'''}{\kappa^2} = -\frac{(x', x'', x''')}{\kappa^2},$$

donde $(u, v, w) = u(v \times w)$ es el *producto mixto* de vectores. Si la parametrización no es la natural tenemos evidentemente

$$B = \frac{x' \times x''}{\|x' \times x''\|},$$

y un cálculo rutinario nos lleva de la expresión que tenemos para $\tau(s)$ a

$$\tau = -\frac{(x', x'', x''')}{\|x' \times x''\|^2}.$$

Ejercicio: Calcular la curvatura y la torsión de la *hélice* $(r \sin t, r \cos t, kt)$.

Para acabar probaremos que la curvatura y la torsión determinan una curva salvo por su posición en el espacio.

Teorema 5.25 Sean $x, \bar{x} : I \rightarrow \mathbb{R}^3$ dos curvas con las mismas funciones κ y τ . Entonces existe una isometría $f : \mathbb{R}^3 \rightarrow \mathbb{R}^3$ tal que $x(s) = f(\bar{x}(s))$ para todo $s \in I$.

DEMOSTRACIÓN: Es fácil comprobar que los vectores del triedro de Frenet, así como la curvatura y la torsión se conservan por isometrías, en el sentido de que, por ejemplo, si f es una isometría se cumple $T_{x \circ f}(s) = T_x \circ \vec{f}$, donde \vec{f} es la isometría lineal asociada a f . Igualmente $\kappa_{x \circ f}(s) = \kappa(s)$, etc.

Sea $s_0 \in I$. Aplicando una isometría a \bar{x} podemos exigir que $x(s_0) = \bar{x}(s_0)$, $T(s_0) = \bar{T}(s_0)$, $N(s_0) = \bar{N}(s_0)$, $B(s_0) = \bar{B}(s_0)$, $\kappa(s_0) = \bar{\kappa}(s_0)$, $\tau(s_0) = \bar{\tau}(s_0)$. Probaremos que en estas condiciones $x = \bar{x}$.

Es claro que

$$\begin{aligned} & \frac{1}{2} \frac{d}{ds} (\|T - \bar{T}\|^2 + \|N - \bar{N}\|^2 + \|B - \bar{B}\|^2) \\ &= (T - \bar{T})(T' - \bar{T}') + (N - \bar{N})(N' - \bar{N}') + (B - \bar{B})(B' - \bar{B}') \end{aligned}$$

Aplicando las fórmulas de Frenet queda

$$\kappa(T - \bar{T})(N - \bar{N}) - \kappa(N - \bar{N})(T - \bar{T}) - \tau(N - \bar{N})(B - \bar{B}) + \tau(B - \bar{B})(N - \bar{N}) = 0,$$

luego $T = \bar{T}$, $N = \bar{N}$, $B = \bar{B}$, pues las diferencias son constantes y se anulan en s_0 . La primera igualdad es $x' = \bar{x}'$, y como ambas funciones coinciden en s_0 ha de ser $x = \bar{x}$. ■

Como consecuencia, si una curva plana tiene curvatura constante κ , entonces es un arco de circunferencia de radio $r = 1/\kappa$, pues su curvatura y su torsión (nula) coinciden con las de la circunferencia. Las curvas de curvatura nula son las rectas.

Cambios de sistema de referencia Trabajar en \mathbb{R}^3 equivale a trabajar en un espacio euclídeo tridimensional E en el que hemos fijado un sistema de referencia ortonormal, de modo que cada punto $P \in E$ se puede identificar con sus coordenadas $x \in \mathbb{R}^3$ respecto al sistema de referencia fijado.

Si $\gamma : I \rightarrow E$ es una curva, donde $I \subset \mathbb{R}$ es un intervalo abierto, sus coordenadas respecto al sistema de referencia fijado son una curva $x : I \rightarrow \mathbb{R}^3$. Si consideramos otro sistema de referencia ortonormal, las coordenadas de γ pasarán a ser $y = a + xM$, para un cierto $a \in \mathbb{R}^3$ y una cierta matriz ortogonal M , es decir, tal que $MM^t = I$. Es claro entonces que la función $x(t)$ es derivable si y sólo si lo es $y(t)$, y en tal caso $y'(t) = x'(t)M$.

Por lo tanto, podemos definir la derivada $\gamma' : I \rightarrow \vec{E}$ como la función que a cada $t \in I$ le asigna el vector cuyas coordenadas en un sistema de referencia ortonormal son $x'(t)$, donde $x(t)$ son las coordenadas de γ . La relación anterior muestra que la definición de γ' no depende del sistema de referencia elegido. Similarmente podemos definir la segunda derivada $\gamma''(t)$.

Del mismo modo podemos definir las derivadas sucesivas de una curva vectorial $\vec{\gamma} : I \rightarrow \vec{E}$ (el único cambio es que ahora la relación entre sus coordenadas en dos sistemas de referencia es $y = xM$, sin la constante a).

Consideremos ahora un sistema de referencia ortonormal $(O, \vec{u}_1, \vec{u}_2, \vec{u}_3)$ que dependa del tiempo, de decir, en el que $O : I \rightarrow E$, $\vec{u}_i : I \rightarrow \vec{E}$ son funciones derivables. El hecho de que los vectores \vec{u}_i formen una base ortonormal en todo instante implica que $\vec{u}_i \vec{u}_j$ es constante (igual a 0 o a 1), luego derivando queda que $\vec{u}_i' \vec{u}_j + \vec{u}_i \vec{u}_j' = 0$. Más explícitamente, $\vec{u}_i \vec{u}_i' = 0$, y además podemos definir las funciones

$$\omega_1 = \vec{u}_2' \vec{u}_3 = -\vec{u}_3' \vec{u}_2, \quad \omega_2 = \vec{u}_3' \vec{u}_1 = -\vec{u}_1' \vec{u}_3, \quad \omega_3 = \vec{u}_1' \vec{u}_2 = -\vec{u}_2' \vec{u}_1. \quad (5.6)$$

Puesto que las derivadas no dependen de la elección ningún sistema de referencia, lo mismo vale para las funciones ω_i . A su vez, definimos la función $\vec{\omega} : I \rightarrow \vec{E}$ dada por

$$\vec{\omega} = \omega_1 \vec{u}_1 + \omega_2 \vec{u}_2 + \omega_3 \vec{u}_3,$$

que tampoco depende de ninguna elección de sistema de referencia (es decir, está determinada por el sistema de referencia móvil).

Consideremos ahora una curva vectorial arbitraria $\vec{\gamma} : I \rightarrow \vec{E}$ y sean x^r las coordenadas de γ relativas al sistema de referencia móvil $(O, \vec{u}_1, \vec{u}_2, \vec{u}_3)$. Esto significa que

$$\vec{\gamma} = \sum_i x_i^r \vec{u}_i.$$

Derivando,

$$\gamma' = \sum_i (x_i^r)' \vec{u}_i + \sum_i x_i^r \vec{u}_i' = \vec{\gamma}'_r + \sum_i x_i^r \vec{u}_i',$$

donde $\vec{\gamma}'_r$ es el vector que en el sistema de referencia móvil tiene coordenadas $(x_i^r)'$. Vamos a desarrollar el último término. Para ello tenemos en cuenta que todo vector \vec{u} cumple que $\vec{u} = \sum_j (\vec{u} \vec{u}_i) \vec{u}_i$, luego

$$\begin{aligned} \sum_i x_i^r \vec{u}_i' &= \sum_{i,j} x_i^r (\vec{u}_i' \vec{u}_j) \vec{u}_j \\ &= (-x_2^r \omega_3 + x_3^r \omega_2) \vec{u}_1 + (x_1^r \omega_3 - x_3^r \omega_1) \vec{u}_2 + (-x_1^r \omega_2 + x_2^r \omega_1) \vec{u}_3 = \vec{\omega} \times \vec{\gamma}. \end{aligned}$$

En resumen:

$$\vec{\gamma}' = \vec{\gamma}'_r + \vec{\omega} \times \vec{\gamma}. \quad (5.7)$$

Consideremos ahora una curva $\gamma : I \rightarrow E$, que podemos interpretar como la trayectoria de un móvil. Si llamamos igualmente x^r a sus coordenadas en el sistema móvil, ahora la relación es

$$\gamma = O + \sum_i x_i^r \vec{u}_i,$$

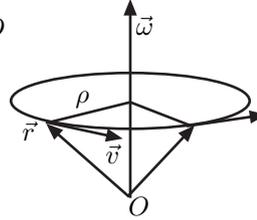
y al derivar obtenemos

$$\vec{v} = \vec{v}_r + \vec{v}_O + \vec{\omega} \times \vec{r}, \quad (5.8)$$

donde $\vec{v} = \gamma'$ es la velocidad del móvil, $\vec{v}_O = O'$ es la velocidad del origen del sistema móvil, \vec{v}_r es la velocidad del móvil medida desde el sistema de referencia móvil y $\vec{r} = \overrightarrow{O\gamma}$ es el vector de coordenadas x_r respecto el sistema móvil.

Esta relación nos da la interpretación del vector $\vec{\omega}$. Si γ es la trayectoria de un punto que desde el sistema de referencia móvil se ve en reposo, entonces \vec{r}' es constante y $\vec{v}_r = \vec{0}$, luego tenemos que $\vec{v} - \vec{v}_O = \vec{\omega} \times \vec{r}$.

Esto significa que si “descontamos” el movimiento de O (o, más precisamente, si sustituimos el sistema de referencia fijo por un segundo sistema de referencia móvil cuyo origen sea la misma curva O , pero cuya base ortonormal sea constante, lo cual equivale a su vez a suponer que O es constante y $\vec{v}_O = \vec{0}$) los objetos fijos respecto al sistema móvil se mueven con velocidad como indica la figura: tangente a una circunferencia de radio ρ con centro en el eje determinado por O y $\vec{\omega}$ y con módulo $v = \omega\rho$, de modo que si $\vec{\omega}$ permaneciera constante, los puntos en reposo respecto del sistema móvil estarían girando alrededor del eje determinado por O y $\vec{\omega}$ con velocidad angular ω .



Podemos expresar esto diciendo que, en cada instante, los ejes del sistema de referencia móvil están girando alrededor del eje determinado por O y $\vec{\omega}$ con velocidad angular ω .

Ahora volvemos a derivar usando la fórmula general (5.7) para derivar \vec{v}_r :

$$\vec{a} = \vec{a}_r + \vec{a}_O + \vec{\omega} \times \vec{v}_r + \vec{\omega}' \times \vec{r} + \vec{\omega} \times \vec{r}'.$$

Llamamos $\vec{\alpha} = \vec{\omega}'$ y aplicamos de nuevo (5.7) para derivar \vec{r}' :

$$\vec{a} = \vec{a}_r + \vec{a}_O + \vec{\omega} \times \vec{v}_r + \vec{\omega} \times (\vec{v}_r + \vec{\omega} \times \vec{r}) + \vec{\alpha} \times \vec{r}.$$

En definitiva:

$$\vec{a}_r = \vec{a} - \vec{a}_O - \vec{\omega} \times (\vec{\omega} \times \vec{r}) - 2\vec{\omega} \times \vec{v}_r - \vec{\alpha} \times \vec{r}. \quad (5.9)$$

Esta fórmula nos da la relación entre las aceleraciones observadas \vec{a}_r y \vec{a} en un mismo móvil respecto de un sistema de referencia móvil y otro fijo. Todo lo dicho tiene pleno sentido desde un punto de vista puramente matemático, pero para entender correctamente la interpretación física de esta fórmula debemos tener en consideración los principios básicos de la mecánica clásica. Nos ocupamos de ello en el apartado siguiente.

Sistemas de referencia inerciales y no inerciales Más arriba hemos dicho que la primera ley de Newton afirma que si un cuerpo no sufre ninguna influencia externa, se mueve con velocidad constante (en particular, en línea recta). Las “influencias externas” se llaman *fuerzas*. En un momento dado, un cuerpo puede estar sometido a distintas fuerzas. Por ejemplo, una persona parada sobre el suelo está sometida a una fuerza de atracción gravitatoria ejercida por la Tierra que la empuja hacia abajo y una fuerza de repulsión eléctrica ejercida por los electrones del suelo que la empuja hacia arriba, de tal modo que ambas fuerzas se compensan y se anulan mutuamente.

Cada fuerza se cuantifica mediante un vector, y la segunda ley de Newton dice que si \vec{F} es la suma de todas las fuerzas que actúan sobre un cuerpo en un instante dado, entonces $\vec{F} = m\vec{a}$, donde \vec{a} es la aceleración que experimenta en ese instante y m es la *masa* del cuerpo.

Así, la segunda ley de Newton generaliza a la primera: cuando la suma total de fuerzas que actúan sobre un cuerpo es nula en todo momento, entonces su aceleración es nula, luego su velocidad permanece constante.

Ahora bien, es fácil ver que estas leyes no son ciertas si no se precisan adecuadamente. En el apartado anterior hemos hablado de “un sistema de referencia fijo” en un espacio euclídeo E y de otro “móvil”, pero estas ideas no tienen una traducción física directa. En física es costumbre llamar *observadores* a los sistemas de referencia. Así, si pensamos, por ejemplo, en un tren en marcha, un observador que viaja en el tren es un sistema de referencia respecto al cual el tren está en reposo, por ejemplo, uno que tenga por ejes coordenados las tres rectas que concurren en una esquina de uno de sus vagones, mientras que un observador situado en una estación es un sistema de referencia respecto al que la estación está en reposo.

Supongamos ahora que A es un observador situado en una estación, que B es otro observador situado en un tren en marcha, y que en el suelo del tren hay una pelota en reposo. Si el tren frena en un momento dado, el observador B verá que la pelota estaba en reposo y, de repente, ha empezado a moverse sin que haya sufrido ninguna influencia externa. Para B no se cumple la segunda ley de Newton. En cambio, para A la situación es muy distinta: la pelota se estaba moviendo con el tren, pero cuando éste ha frenado, precisamente por la primera ley de Newton, ha mantenido su velocidad constante, y por eso ha empezado a moverse respecto del tren.

Tal vez el lector concluya que las leyes de la física son válidas respecto de observadores “en reposo”, y que en otro caso pueden fallar debido a que parte del movimiento observado en un objeto puede ser un reflejo del movimiento del observador. Desgraciadamente, no podemos afirmar tal cosa porque no tiene sentido físico hablar de objetos en reposo de forma absoluta. Podemos considerar que el tren se mueve y la estación está parada, pero la estación está sobre la Tierra, que se mueve alrededor del Sol, y el Sol también se mueve respecto del centro de la Galaxia. En general, si pensamos en un universo con muchos objetos en movimiento, no hay ningún criterio que nos permita decir que uno dado está en reposo “de verdad” y los que se mueven respecto de él se mueven “de verdad”. Pero esto no es realmente un problema. Las leyes de Newton se cumplen para todos los observadores *inerciales*, que son los sistemas de referencia determinados por objetos sobre los que la suma de fuerzas a las que están sometidos es nula (y cuyos ejes no giran).

Ahora podemos interpretar la fórmula a la que hemos llegado en el apartado anterior entendiendo que el sistema de referencia “fijo” en \mathbb{R}^3 es un observador inercial A , y que el sistema de referencia “móvil” es otro observador B , no necesariamente inercial. Si entendemos que la curva γ que hemos analizado es la trayectoria de un objeto de masa m , multiplicando por m la relación entre las aceleraciones resulta:

$$m\vec{a}_r = \vec{F} + \vec{f}_O + \vec{f}_{\text{cent}} + \vec{f}_{\text{cor}} + \vec{f}_{\text{Eu}}, \quad (5.10)$$

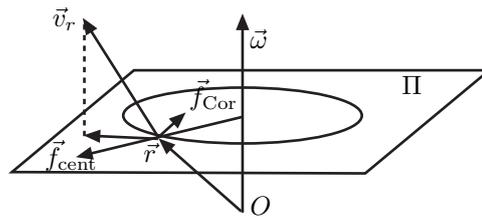
donde

- \vec{a}_r es la aceleración del objeto observada por B .

- $\vec{F} = m\vec{a}$ es la suma de todas las fuerzas que se están ejerciendo sobre el móvil (pues suponemos que el observador A es inercial, luego para él se cumple la segunda ley de Newton).¹
- $\vec{f}_O = -m\vec{a}_O$, $\vec{f}_{\text{cent}} = -m\vec{\omega} \times (\vec{\omega} \times \vec{r})$, $\vec{f}_{\text{Cor}} = -2m\vec{\omega} \times \vec{v}_r$, $\vec{f}_{\text{Eu}} = -m\vec{\alpha} \times \vec{r}$.

Si sobre el observador B es inercial, los cuatro últimos términos de la ecuación son nulos, y (5.10) se reduce a la segunda ley de Newton. En general, lo que afirma (5.10) es que las aceleraciones observadas en un objeto por el observador B son las que observaría un observador inercial bajo el supuesto de que, además de las “fuerzas reales” debidas a la interacción con otros objetos (gravedad, electricidad, etc.), sobre cada objeto actúan cuatro “fuerzas ficticias”, llamadas también *fuerzas de inercia*, que no son sino el reflejo del movimiento observador (de su origen y de sus ejes) respecto de un observador inercial. Vamos a describirlas:

- $\vec{f}_O = -m\vec{a}_O$ tiene el sentido opuesto a la fuerza total que está actuando sobre O (pero no la misma intensidad: si el tren tiene masa M y frena por efecto de una fuerza $\vec{F}_{\text{fren}} = M\vec{a}_O$, entonces sobre la pelota actúa una fuerza ficticia $\vec{f}_O = -m\vec{a}_O$). Ésta es la fuerza que “empuja” la pelota cuando el tren frena. Si la fuerza de frenado es en sentido opuesto a la marcha del tren, la pelota se moverá en el sentido de la marcha del tren.
- $\vec{f}_{\text{cent}} = -m\vec{\omega} \times (\vec{\omega} \times \vec{r})$ está en el plano Π perpendicular a $\vec{\omega}$ por $\gamma(t)$ y empuja al objeto radialmente en sentido opuesto al eje. Se llama *fuerza centrífuga*.
- $\vec{f}_{\text{Cor}} = -2m\vec{\omega} \times \vec{v}_r$ está también en Π y además es perpendicular a la velocidad relativa al observador B . Se llama *fuerza de Coriolis*.
- $\vec{f}_{\text{Eu}} = -m\vec{\alpha} \times \vec{r}$ no tiene ningún nombre en particular, aunque en algunos libros la llaman *fuerza de Euler*. Depende de la variación de $\vec{\omega}$.



¹Esto presupone que las fuerzas que se ejercen sobre un objeto son vectores que no dependen del observador. Por ejemplo, la fuerza gravitatoria que un objeto ejerce sobre otro depende de las masas de ambos y del vector que va de la posición de uno al del otro, y nada de esto depende del observador. En realidad no se puede decir lo mismo de las fuerzas electromagnéticas, pero ésta es precisamente la razón por la que la mecánica newtoniana debe reemplazarse por la mecánica relativista. Salvo en situaciones en las que intervienen objetos que se mueven a velocidades cercanas a las de la luz, las diferencias entre una y otra son inapreciables, por lo que podemos suponer que vector \vec{F} no depende del observador.

Por ejemplo, tomemos un sistema de referencia A con origen en el centro del Sol, con el eje Z paralelo al eje de rotación de la Tierra (que podemos suponer fijo) y con los otros dos ejes fijos con respecto a alguna estrella lejana, que también podemos suponer fija en la práctica. Por otra parte, consideremos un observador B que se mueve con la Tierra, cuyo origen de coordenadas O sea el centro de la Tierra, el eje Z sea el eje de rotación de la Tierra y cuyos otros dos ejes estén fijos respecto de dos puntos de la superficie terrestre. De este modo, su base asociada tendrá coordenadas

$$i = (\cos \omega t, \sin \omega t, 0), \quad j = (-\sin \omega t, \cos \omega t, 0), \quad k = (0, 0, 1),$$

donde ω es la velocidad angular de la Tierra, es decir, $\pi/12$ radianes por hora. Usando (5.6) se comprueba inmediatamente que $\vec{\omega} = (0, 0, \omega)$, como era de esperar. Ahora podemos estimar la intensidad de las fuerzas de inercia a las que está sometido todo objeto situado sobre la superficie terrestre. Supongamos que el objeto tiene una masa $m = 1$ kg.

La fuerza \vec{f}_O depende de las fuerzas externas que actúan sobre la Tierra. Las principales son la atracción del Sol y de la Luna. También podríamos considerar las de los demás planetas del sistema solar. Según la *ley de gravitación universal* de Newton, su módulo viene dado por la fórmula

$$F = \frac{GMm}{r^2},$$

donde G es la *constante gravitatoria*, M y m son las masas de los objetos considerados (por ejemplo la del Sol y la de la Tierra) y r es la distancia que los separa. Si consideramos únicamente la gravedad solar, entonces

$$M_T a_O = \frac{GM_S M_T}{r_S^2} \Rightarrow a_O = \frac{GM_S}{r_S^2} \Rightarrow f_O = \frac{GM_S m}{r_S^2}.$$

La tabla siguiente contiene toda la información necesaria para los cálculos, las masas de la Tierra, el Sol y la Luna, el radio de la Tierra, las distancias a la Tierra del Sol y la Luna, la constante de gravitación y la velocidad angular de rotación de la Tierra:

M_T	$5.9736 \cdot 10^{24}$ kg	r_T	$6.371 \cdot 10^6$ m
M_S	$1.9891 \cdot 10^{30}$ kg	R_S	$1.496 \cdot 10^{11}$ m
M_L	$7.349 \cdot 10^{22}$ kg	R_L	$3.844 \cdot 10^8$ m
G	$6.6738 \cdot 10^{-11}$ Nm ² /kg ²	ω	$7.2722 \cdot 10^{-5}$ s ⁻¹

Es fácil ver entonces que la fuerza de inercia debida a la atracción del Sol sobre la Tierra es de² 0.0059 N, lo que supone el 0.06% del peso del objeto, que es de 9.8 N. Similarmente, la fuerza de inercia debida a la Luna es de 0.00003 N.

Vemos, pues, que estas fuerzas son en general despreciables a la hora de estudiar el movimiento del objeto, como también lo es la atracción gravitatoria que sufre el cuerpo por parte del Sol o de la Luna, del mismo orden de magnitud.

²La unidad de fuerza es el Newton, definido como la fuerza necesaria para que un cuerpo de 1 kg de masa experimente una aceleración de 1 m/s².

La fuerza centrífuga es perpendicular al eje de rotación de la Tierra. Es nula en los polos, donde $\vec{\omega} \times \vec{r} = \vec{0}$, y es máxima en el ecuador, donde su intensidad es de $mR_T\omega^2 = 0.03369\text{ N}$, un 0.34% del peso del objeto. Si la Tierra girara mucho más rápidamente, la fuerza centrífuga haría que los cuerpos “cayeran” hacia el cielo, verticalmente en el ecuador y oblicuamente según la latitud.

La fuerza de Coriolis depende de la velocidad del objeto respecto de la Tierra. Si suponemos que se mueve a 1 m/s, será máxima si es perpendicular al eje de rotación de la Tierra, y entonces su módulo es $2\omega v = 0.0001454\text{ N}$, el 0.0015% del peso. Sobre un cuerpo que cae, su velocidad \vec{v}_r apunta al centro de la Tierra y $-\vec{\omega} \times \vec{v}_r$ apunta hacia el Este, por lo que los cuerpos que caen sufren una desviación hacia el Este, nula en los polos y máxima en el ecuador.

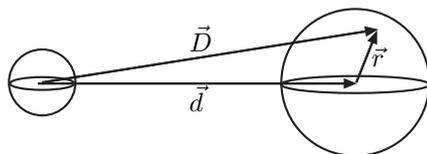
Si el movimiento se realiza sobre la superficie de la Tierra observamos que $\vec{\omega}$ apunta hacia el exterior de la misma en el hemisferio Norte y hacia el interior en el hemisferio Sur, por lo que $-\vec{\omega} \times \vec{v}_r$ apunta hacia la derecha de \vec{v}_r en el hemisferio norte y hacia la izquierda en el hemisferio sur. Si descomponemos esta fuerza en dos vectores, uno en la dirección del centro de la Tierra y otro tangente a la misma, la gravedad hace inadmisible la primera componente, pero la segunda es apreciable. Ésta es nula en el ecuador y máxima en los polos. En la subsección 7.3.3 veremos cómo la fuerza de Coriolis produce el movimiento del péndulo de Foucault.

Como hemos supuesto $\vec{\omega}$ constante, no hay fuerza de Euler.

Ejercicio: En realidad $\vec{\omega}$ describe un cono cuya generatriz forma un ángulo de 23.5° con su altura y tarda 25 776 años en dar una vuelta completa. Calcular la fuerza de Euler con estos datos.

Estos cálculos muestran que la Tierra puede considerarse un sistema de referencia inercial en la mayor parte de situaciones prácticas (y lo mismo es aplicable al Sol, y así lo hemos hecho en todo el razonamiento precedente). No obstante, hay situaciones en las que las fuerzas de inercia son relevantes. Ya hemos citado el caso del péndulo de Foucault, que estudiaremos más adelante, y a continuación vamos a analizar otro caso.

Mareas En las mismas condiciones del apartado precedente, consideremos la ecuación (5.10) para un objeto de masa m sobre el que vamos a considerar el efecto de la gravitación terrestre y la gravitación lunar. Esto significa que vamos a sustituir \vec{F} por $\vec{F} + \vec{F}_{\text{gr}} + \vec{F}_{\text{grL}}$, donde \vec{F}_{gr} es la fuerza de la gravedad terrestre, \vec{F}_{grL} es la fuerza de la gravedad Lunar y \vec{F} es ahora la suma de las demás fuerzas que puedan estar actuando sobre el cuerpo. Igualmente, calcularemos \vec{f}_O teniendo en cuenta la acción de la Luna sobre la Tierra.



Llamamos \vec{d} al vector que une el centro de la Luna con el centro de la Tierra, y \vec{D} al que une el centro de la Luna con la posición del objeto. Así, según la ley de gravitación universal:

$$\vec{F}_{\text{gr}} = -\frac{GmM_T}{r^3} \vec{r}, \quad \vec{F}_{\text{grL}} = -\frac{GmM_L}{D^3} \vec{D}.$$

Por otra parte, según la segunda ley de Newton tenemos que

$$M_T \vec{a}_O = -\frac{GM_T M_L}{d^3} \vec{d} \Rightarrow \vec{a}_O = -\frac{GM_L}{d^3} \vec{d} \Rightarrow \vec{f}_O = \frac{GM_L m}{d^3} \vec{d}.$$

De este modo, (5.10) tiene la forma

$$m\vec{a}_r = \vec{F} - \frac{GmM_T}{r^3} \vec{r} - \frac{GmM_L}{D^3} \vec{D} + \frac{GmM_L}{d^3} \vec{d} + \vec{f}_{\text{cent}} + \vec{f}_{\text{cor}}.$$

(Como $\vec{\omega}$ es constante, no hay fuerza de Euler.) Equivalentemente,

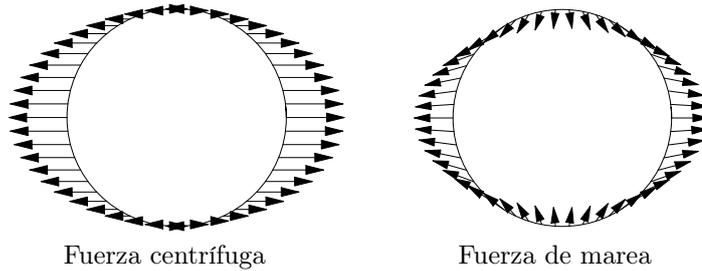
$$m\vec{a}_r = \vec{F} + \vec{F}_{\text{gr}} + \vec{f}_{\text{mar}} + \vec{f}_{\text{cor}} + \vec{f}_{\text{cent}},$$

donde

$$\vec{f}_{\text{mar}} = -GmM_L \left(\frac{1}{D^3} \vec{D} - \frac{1}{d^3} \vec{d} \right)$$

es la *fuerza de marea*, o la fuerza debida a la gravedad de la Luna, que es la diferencia entre la gravitación lunar propiamente dicha y la fuerza ficticia que refleja la aceleración que la Luna provoca al centro de la Tierra y que formalmente es la fuerza con la que la Luna atraería al objeto si estuviera situado en el centro de la Tierra.

La figura muestra la fuerza centrífuga y la fuerza de marea en diferentes puntos de la superficie terrestre, calculadas con las fórmulas y los datos que hemos dado más arriba, aunque están multiplicadas por un factor $k = 10^8$ en el caso de la fuerza centrífuga y $k = 2 \cdot 10^{12}$ en el caso de la fuerza de marea, para que sean apreciables (se supone que la Luna está a la izquierda):



La fuerza centrífuga es mucho más intensa que la fuerza de marea. Por ejemplo, en el punto del ecuador más cercano a la Luna, las intensidades son:

$$F_{\text{gr}} = 9.8\text{N/kg}, \quad F_{\text{cent}} = 0.3369\text{N/kg}, \quad F_{\text{mar}} = 1.1274 \cdot 10^{-6}\text{N/kg},$$

con lo que la fuerza centrífuga que experimenta un cuerpo es del orden del 0.34% de su peso y la fuerza de marea del orden del 0.000011% de su peso. Son demasiado pequeñas para que puedan apreciarse en un objeto en concreto, pero su efecto se nota en la masa de agua que forma los océanos. La fuerza centrífuga hace que el nivel del mar sea más elevado en el ecuador que en los polos, mientras que la fuerza de marea provoca un efecto similar en los puntos de longitud próxima a la del eje Tierra-Luna.

Sin embargo, una diferencia esencial entre ambas es que la fuerza centrífuga es simétrica respecto del eje de rotación de la Tierra (el eje vertical en la figura), mientras que la fuerza de marea es simétrica respecto del eje Tierra-Luna (el eje horizontal en la figura, de modo que la figura de la derecha es igualmente válida si consideramos que el eje de rotación de la Tierra es el vertical o el perpendicular al papel). Como consecuencia, aunque la fuerza de marea es mucho más débil que la centrífuga, su efecto es más llamativo, ya que, para un observador situado en la Tierra, la Luna gira a su alrededor una vez al día, por lo que los abultamientos que provoca en el nivel del mar giran también una vez al día y, como hay dos abultamientos opuestos, sucede que el nivel del mar sube y baja en cada sitio dos veces al día. Ésta es la explicación del fenómeno de las mareas. El Sol también produce mareas, aunque de menor intensidad, pero las mareas lunares se intensifican sensiblemente cuando la Luna está alineada con el Sol.

Puede parecer sorprendente que la presencia de la Luna “a la izquierda” de la Tierra provoque una fuerza “hacia la derecha” en la cara opuesta de la Tierra, pero no es tan extraño si nos damos cuenta de que se trata de una fuerza ficticia: cuando un coche acelera “hacia la izquierda”, sus ocupantes sienten una fuerza ficticia “hacia la derecha”.

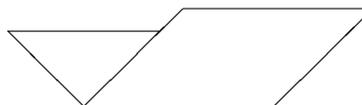
Capítulo VI

Introducción a las variedades diferenciables

En este capítulo aplicaremos el cálculo diferencial al estudio de las superficies. Si bien todos los ejemplos que consideraremos serán bidimensionales, la mayor parte de la teoría la desarrollaremos sobre un concepto general de “superficie de n dimensiones”. La idea básica es que una superficie es un espacio topológico que localmente se parece a un plano. El ejemplo típico es la superficie terrestre: tenemos que alejarnos mucho de ella para darnos cuenta de que no es plana, sino esférica. Una definición topológica que recoja estas ideas sería la siguiente:

Un subconjunto S de \mathbb{R}^n es una superficie si para cada punto $p \in S$ existe un entorno V de p , un abierto U en \mathbb{R}^2 y un homeomorfismo $X : U \rightarrow V \cap S$.

Es decir, S es una superficie si alrededor de cada punto es homeomorfa a un abierto de \mathbb{R}^2 . Notar que no pedimos que S sea homeomorfa a un abierto de \mathbb{R}^2 , sino sólo que lo sea alrededor de cada punto. Basta pensar en una esfera para comprender la importancia de este hecho. Una esfera no es homeomorfa a un abierto de \mathbb{R}^2 , pero un pequeño trozo de esfera es como un trozo de plano abombado, homeomorfo a un trozo de plano “llano”. Sin embargo nosotros estamos interesados en superficies diferenciables, en el sentido de que se parezcan a planos afines alrededor de cada punto. Podría pensarse que para conseguir esto bastaría exigir que el homeomorfismo X sea diferenciable, pero no es así. Por ejemplo, pensemos en $X(u, v) = (u^3, v, |u^3|)$. La aplicación X es diferenciable, y es un homeomorfismo entre \mathbb{R}^2 y un conjunto $S \subset \mathbb{R}^3$ cuya forma es la de una hoja de papel doblada por la mitad. Alrededor de los puntos de la forma $(0, v, 0)$ no se parece a ningún plano, sino que tiene un “pico”. Si la Tierra tuviera esta forma no necesitaríamos alejarnos de ella para darnos cuenta de que estaría “doblada”.



La razón es que $dX(0, b) = (0, dv(0, b), 0)$, de modo que alrededor de un

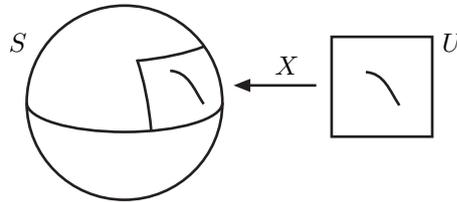
punto $(0, b)$ la función X se parece a la aplicación afín $f(u, v) = (0, v, 0)$, cuya imagen es la recta $x = z = 0$. Así pues, aunque topológicamente la imagen de X es homeomorfa a un plano, desde el punto de vista del cálculo diferencial la imagen de X alrededor de un punto $(0, y, 0)$ se parece a la recta $x = z = 0$, y no a un plano. Para evitar esto hemos de exigir que la imagen de $dX(u, v)$ sea un plano y no una recta. Esto es tanto como decir que la matriz jacobiana tenga rango 2.

6.1 Variedades

Definición 6.1 Un conjunto $S \subset \mathbb{R}^m$ es una *variedad diferenciable* de dimensión $n \leq m$ y de clase C^q si para cada punto $p \in S$ existe un entorno V de p , un abierto U en \mathbb{R}^n y una función $X : U \rightarrow \mathbb{R}^m$ de clase C^q de modo que el rango de la matriz JX sea igual a n en todo punto y $X : U \rightarrow S \cap V$ sea un homeomorfismo. Una aplicación X en estas condiciones se llama *carta* de S alrededor de p .

En lo sucesivo supondremos que las variedades con las que trabajamos son de clase C^q para un q suficientemente grande como para que existan las derivadas que consideremos (y sean continuas). Rara vez nos hará falta suponer $q > 3$, aunque de hecho todos los ejemplos que consideraremos serán de clase C^∞ .

La palabra “carta” hay que entenderla en el sentido de “mapa”. En efecto, podemos pensar en U como un mapa “plano” de una región de S , y la aplicación X es la que hace corresponder cada punto del mapa con el punto real que representa.



Alternativamente, podemos pensar en X^{-1} como una aplicación que asigna a cada punto $p \in S \cap V$ unas *coordenadas* $x = (x_1, \dots, x_n) \in U \subset \mathbb{R}^n$, de forma análoga a los sistemas de coordenadas en un espacio afín.¹ Dentro de poco será equivalente trabajar con cartas o con sistemas de coordenadas, pero por el momento podemos decir que las cartas son diferenciables y en cambio no tiene sentido decir que las funciones coordenadas lo sean, pues no están definidas sobre abiertos de \mathbb{R}^m .

¹Etimológicamente, una “variedad” no es más que un conjunto cuyos elementos vienen determinados por “varias” coordenadas. En los resultados generales llamaremos x_1, \dots, x_n a las coordenadas para marcar la analogía con \mathbb{R}^n , aunque en el caso de curvas seguiremos usando la variable t (o s si la parametrización es la natural) y en el caso de superficies $S \subset \mathbb{R}^3$ usaremos x, y, z para las coordenadas en \mathbb{R}^3 y u, v para las coordenadas en S .

Ejemplo Todo abierto U en \mathbb{R}^n es una variedad diferenciable de dimensión n y clase C^∞ . Basta tomar como carta la identidad en U . De este modo, todos los resultados sobre variedades valen en particular para \mathbb{R}^n y sus abiertos. ■

Ejercicio: Refinar el argumento del teorema 3.24 para concluir que dos puntos cualesquiera de una variedad conexa S de clase C^q pueden ser unidos por un arco de clase C^q contenido en S .

El teorema siguiente proporciona una clase importante de variedades diferenciables, pues a continuación vemos que toda variedad es localmente de este tipo.

Teorema 6.2 Sea $f : U \subset \mathbb{R}^n \rightarrow \mathbb{R}^k$ una aplicación de clase C^q sobre un abierto U y $X : U \rightarrow \mathbb{R}^{n+k}$ la aplicación dada por $X(x) = (x, f(x))$. Entonces $X[U]$ es una variedad diferenciable de dimensión n y clase C^q .

DEMOSTRACIÓN: Basta observar que X es obviamente un homeomorfismo en su imagen (su inversa es una proyección) y $JX(x)$ contiene una submatriz de orden n igual a la identidad, luego su rango es n . La definición se satisface tomando $V = \mathbb{R}^{n+k}$. ■

Notemos que $X[U]$ es la gráfica de f , luego el teorema anterior afirma que la gráfica de una función diferenciable es siempre una variedad diferenciable. Ahora veamos que todo punto de una variedad diferenciable tiene un entorno en el que la variedad es la gráfica de una función.

Teorema 6.3 Sea $S \subset \mathbb{R}^{n+k}$ una variedad de clase C^q y de dimensión n . Sea $p \in S$. Entonces existe un entorno V de p , un abierto U en \mathbb{R}^n y una función $f : U \rightarrow \mathbb{R}^k$ de clase C^q de modo que la aplicación $X : U \rightarrow \mathbb{R}^{n+k}$ dada por $X(x) = (x, f(x))$ es una carta alrededor de p .

En realidad hemos de entender que las coordenadas de x y $f(x)$ se intercalan en un cierto orden que no podemos elegir, tal y como muestra la prueba.

DEMOSTRACIÓN: Sea $Y : W \rightarrow \mathbb{R}^{n+k}$ una carta alrededor de p . Sea V un entorno de p tal que $Y : W \rightarrow S \cap V$ sea un homeomorfismo. Sea $t_0 \in W$ el vector de coordenadas de p , es decir, $Y(t_0) = p$. Puesto que $JY(t_0)$ tiene rango n , reordenando las funciones coordenadas de Y podemos suponer que el determinante formado por las derivadas parciales de las n primeras es no nulo. Digamos que $Y(t) = (Y_1(t), Y_2(t))$, donde $|JY_1(t_0)| \neq 0$. Sea $p_1 = Y_1(t_0)$.

Por el teorema de inyectividad local y el teorema de la función inversa, existe un entorno abierto $G \subset W$ de t_0 tal que Y_1 es inyectiva en G , $Y_1[G] = U$ es abierto en \mathbb{R}^n y la función $Y_1^{-1} : U \rightarrow G$ es de clase C^q .

El conjunto $Y[G]$ es un entorno abierto de p en $S \cap V$, luego existe un entorno abierto V' de p en \mathbb{R}^{n+k} tal que $Y[G] = S \cap V \cap V'$. Cambiando V' por $V \cap V'$ podemos suponer que $V' \subset V$ y que $Y[G] = S \cap V'$.

De este modo, cada punto $p \in S \cap V'$ está determinado por sus coordenadas $t \in G$, las cuales a su vez están determinadas por $x = Y_1(t) \in U$, con la particularidad de que x es el vector de las primeras componentes de p . Concretamente

p está formado por x y $f(x) = Y_2(Y_1^{-1}(x))$. La función f es de clase C^q en U . De este modo, si $x \in U$ y $t = Y_1^{-1}(x) \in G$, tenemos que

$$X(x) = (x, f(x)) = (Y_1(t), Y_2(t)) = Y(t) \in S \cap V'.$$

Recíprocamente, si $(x, y) \in S \cap V'$, entonces $x = Y_1(t)$, $y = Y_2(t)$ para un cierto $t \in G$, luego $(x, y) = (x, f(x)) = X(x)$. En definitiva tenemos que $X : U \rightarrow S \cap V'$ es biyectiva. Su inversa es la proyección en las primeras componentes, luego X es un homeomorfismo. ■

En las condiciones de la prueba anterior, sea $\pi : \mathbb{R}^m \rightarrow \mathbb{R}^n$ la proyección en las n primeras componentes y $g : V' \rightarrow G$ la aplicación dada por $g(x) = Y_1^{-1}(\pi(x))$. Notemos que si $x \in V'$ entonces $\pi(x) \in U$, luego g está bien definida y es de clase C^q . Si $t \in G$ entonces

$$g(Y(t)) = g(Y_1(t), Y_2(t)) = Y_1^{-1}(Y_1(t)) = t,$$

luego $(Y|_G)^{-1}$ es la restricción a $V' \cap S$ de g . Con esto hemos probado:

Teorema 6.4 *Sea $Y : U \rightarrow S \subset \mathbb{R}^m$ una carta de una variedad diferenciable de dimensión n y clase C^q . Para cada punto $t \in U$ existe un entorno $G \subset U$ de t , un entorno V de $Y(t)$ y una aplicación $g : V \rightarrow G$ de clase C^q tal que $(Y|_G)^{-1} = g|_{V \cap S}$.*

De aquí se sigue una propiedad fundamental de las cartas:

Teorema 6.5 *Sea $S \subset \mathbb{R}^m$ una variedad diferenciable de dimensión n y de clase C^q . Sea $p \in S$ y $X : U \rightarrow S \cap V$, $Y : U' \rightarrow S \cap V'$ dos cartas alrededor de p . Sean $V_0 = V \cap V'$, $U_0 = X^{-1}[V_0]$, $U'_0 = Y^{-1}[V_0]$. Entonces la aplicación $X \circ Y^{-1} : U'_0 \rightarrow U_0$ es biyectiva, de clase C^q y con determinante jacobiano no nulo, con lo que su inversa es también de clase C^q .*

DEMOSTRACIÓN: Si $t \in U_0$, por el teorema anterior existe una función g de clase C^q definida en un entorno de $X(t)$ de modo que $X \circ Y^{-1} = X \circ g$ (en un entorno de t), luego $X \circ Y^{-1}$ es de clase C^q en un entorno de t , luego en todo U_0 . Lo mismo vale para su inversa $Y \circ X^{-1}$, luego la regla de la cadena nos da que sus diferenciales son mutuamente inversas, luego los determinantes jacobianos son no nulos. ■

Veremos ahora otro ejemplo importante de variedades diferenciables. Primeramente consideraremos el caso lineal al cual generaliza.

Ejemplo Una variedad afín de dimensión n en \mathbb{R}^m es también una variedad diferenciable de la misma dimensión y de clase C^∞ . En efecto, una tal variedad está formada por los puntos que satisfacen un sistema de $m - n$ ecuaciones lineales linealmente independientes. Esto implica que la matriz de coeficientes del sistema tiene un determinante de orden $m - n$ no nulo, luego agrupando

adecuadamente las variables podemos expresar el sistema como $xA + yB = c$, donde $x \in \mathbb{R}^n$, $y \in \mathbb{R}^{m-n}$, $|B| \neq 0$, luego podemos despejar

$$y = f(x) = (c - xA)B^{-1},$$

donde la función f es obviamente de clase C^∞ . Esto significa que la variedad lineal está formada por los puntos (x, y) tales que $y = f(x)$, luego es la gráfica de f y por consiguiente es una variedad de clase C^∞ . ■

Ahora probamos que las soluciones de un sistema de $k = m - n$ ecuaciones diferenciables con m incógnitas constituyen una variedad de dimensión n supuesto que se cumpla una condición de independencia similar a la independencia lineal que exigíamos en el ejemplo anterior.

Definición 6.6 Si $f : A \subset \mathbb{R}^{n+k} \rightarrow \mathbb{R}^k$ es diferenciable en $(x, y) \in A$, donde $x \in \mathbb{R}^n$, $y \in \mathbb{R}^k$, definimos

$$\frac{\partial(f_1 \cdots f_k)}{\partial(y_1 \cdots y_k)}(x, y) = \begin{vmatrix} D_{n+1}f_1(x, y) & \cdots & D_{n+1}f_k(x, y) \\ \vdots & & \vdots \\ D_{n+k}f_1(x, y) & \cdots & D_{n+k}f_k(x, y) \end{vmatrix}.$$

Teorema 6.7 (Teorema de la función implícita) Consideremos una aplicación $f : A \subset \mathbb{R}^{n+k} \rightarrow \mathbb{R}^k$ de clase C^q en el abierto A , con $q \geq 1$. Sea $(x^0, y^0) \in A$ tal que $f(x^0, y^0) = 0$ y supongamos que

$$\frac{\partial(f_1 \cdots f_k)}{\partial(y_1 \cdots y_k)}(x^0, y^0) \neq 0.$$

Entonces existen abiertos $V \subset A$, $U \subset \mathbb{R}^n$ de modo que $(x^0, y^0) \in A$, $x^0 \in U$ y una función $g : U \rightarrow \mathbb{R}^k$ de clase C^q tal que

$$\{(x, y) \in V \mid f(x, y) = 0\} = \{(x, y) \in \mathbb{R}^{n+k} \mid x \in U, y = g(x)\}.$$

DEMOSTRACIÓN: Sea $F : A \rightarrow \mathbb{R}^{n+k}$ la función $F(x, y) = (x, f(x, y))$. Sus funciones coordenadas son las proyecciones en las componentes de \mathbb{R}^n más las funciones coordenadas de f , luego F es de clase C^q . Su determinante jacobiano es

$$\begin{vmatrix} 1 & \cdots & 0 & D_1f_1(x, y) & \cdots & D_1f_k(x, y) \\ \vdots & & \vdots & \vdots & & \vdots \\ 0 & \cdots & 1 & D_nf_1(x, y) & \cdots & D_nf_k(x, y) \\ 0 & \cdots & 0 & D_{n+1}f_1(x, y) & \cdots & D_{n+1}f_k(x, y) \\ \vdots & & \vdots & \vdots & & \vdots \\ 0 & \cdots & 0 & D_{n+k}f_1(x, y) & \cdots & D_{n+k}f_k(x, y) \end{vmatrix},$$

que claramente coincide (salvo signo) con

$$\frac{\partial(f_1 \cdots f_k)}{\partial(y_1 \cdots y_k)}(x, y).$$

Así $|J_F(x^0, y^0)| \neq 0$. Por el teorema de inyectividad local existe un entorno V de (x^0, y^0) donde F es inyectiva y su determinante jacobiano no se anula. Por el teorema de la función inversa tenemos que $W = F[V]$ es abierto en \mathbb{R}^{n+k} y la función $G = F^{-1} : W \rightarrow V$ es de clase C^q .

Podemos expresar $G(x, y) = (G_1(x, y), G_2(x, y))$. Claramente G_1 y G_2 son ambas de clase C^q . Si $(x, y) \in W$, entonces

$$(x, y) = F(G(x, y)) = F(G_1(x, y), G_2(x, y)) = (G_1(x, y), f(G(x, y))),$$

luego $G_1(x, y) = x$, con lo que en general $G(x, y) = (x, G_2(x, y))$.

Definimos $U = \{x \in \mathbb{R}^n \mid (x, 0) \in W\}$. Es claro que se trata de un abierto. Además, $F(x^0, y^0) = (x^0, 0) \in W$, luego $x^0 \in U$. Definimos $g : U \rightarrow \mathbb{R}^k$ mediante $g(x) = G_2(x, 0)$. Claramente g es de clase C^q .

Tomemos ahora $x \in U$ e $y = g(x)$. Hemos de probar que $(x, y) \in V$ y $f(x, y) = 0$. En efecto, por definición de U es $(x, 0) \in W$, luego $G(x, 0) \in V$, pero

$$G(x, 0) = (x, G_2(x, 0)) = (x, g(x)) = (x, y).$$

Además $(x, 0) = F(G(x, 0)) = F(x, y) = (x, f(x, y))$, luego $f(x, y) = 0$.

Recíprocamente, si $(x, y) \in V$ y $f(x, y) = 0$ entonces

$$F(x, y) = (x, f(x, y)) = (x, 0) \in W,$$

luego $x \in U$ y $(x, y) = G(F(x, y)) = G(x, 0) = (x, G_2(x, 0)) = (x, g(x))$, con lo que $g(x) = y$. ■

Lo que afirma este teorema es que si $S = \{x \in \mathbb{R}^{n+k} \mid f(x) = 0\}$ es un conjunto determinado por un sistema de k ecuaciones de clase C^q y $p \in S$ cumple la hipótesis entonces $V \cap S = X[U]$, donde $X(x) = (x, g(x))$, de donde se sigue que X cumple las condiciones para ser una carta de S alrededor de p . Si la hipótesis se cumple en todo punto entonces S es una variedad diferenciable de dimensión n .

Por ejemplo, si $f(x, y, z) = x^2 + y^2 + z^2 - r^2$, entonces el conjunto S es una esfera. Para comprobar que se trata de una superficie de clase C^∞ basta comprobar que en cada punto al menos una de las derivadas

$$\frac{\partial f}{\partial x} = 2x, \quad \frac{\partial f}{\partial y} = 2y, \quad \frac{\partial f}{\partial z} = 2z,$$

es no nula, pero las tres sólo se anulan simultáneamente en $(0, 0, 0)$, que no es un punto de S , luego, efectivamente, la esfera es una superficie diferenciable.

Es importante observar que la derivada que no se anula no siempre es la misma. Por ejemplo, en el polo norte $(0, 0, r)$ la única derivada que no se anula es la de z , luego en un entorno podemos expresar z como función $z(x, y)$. Concretamente, $z = \sqrt{r^2 - x^2 - y^2}$. Similarmente, la porción de esfera alrededor del polo sur es la gráfica de la función $z = -\sqrt{r^2 - x^2 - y^2}$. En cambio, alrededor de $(r, 0, 0)$ la esfera no es la gráfica de ninguna función $z(x, y)$. Es fácil ver

que dado cualquier entorno U de $(r, 0, 0)$ y cualquier entorno V de $(r, 0)$ siempre hay puntos (x, y) en U para los cuales hay dos puntos distintos $(x, y, \pm z)$ en U (con lo que (x, y) debería tener dos imágenes) y puntos (x, y) con $x^2 + y^2 > r^2$ para los que no existe ningún z tal que $(x, y, z) \in U$. Sin embargo, alrededor de este punto la esfera es la gráfica de la función $x = \sqrt{r^2 - y^2 - z^2}$.

El mismo argumento prueba en general que la *esfera* de dimensión n

$$S^n = \{x \in \mathbb{R}^{n+1} \mid \|x\|_2^2 = 1\}$$

es una variedad diferenciable.

Ejemplo: superficies de revolución Sea C una variedad diferenciable de dimensión 1 en \mathbb{R}^2 . Supongamos que todos sus puntos (x, z) cumplen $x > 0$. Llamaremos *superficie de revolución* generada por C al conjunto

$$S = \{(x, y, z) \in \mathbb{R}^3 \mid (\sqrt{x^2 + y^2}, z) \in C\}.$$

El conjunto S está formado por todos los puntos que resultan de girar alrededor del eje Z los puntos de C . Vamos a ver que se trata de una variedad diferenciable de dimensión 2.

Tomemos $(x_0, y_0, z_0) \in S$ y $\bar{x}_0 = \sqrt{x_0^2 + y_0^2}$. Entonces $(\bar{x}_0, z_0) \in C$. Sea $\alpha(u) = (r(u), z(u))$ una carta de C alrededor de este punto, digamos $r(u_0) = \bar{x}_0$, $z(u_0) = z_0$. Por definición existe un entorno V_0 de u_0 y un entorno U_0 de (\bar{x}_0, z_0) de modo que $C \cap U_0 = \alpha[V_0]$. Sea

$$X(u, v) = (r(u) \cos v, r(u) \sin v, z(u)), \quad (u, v) \in V_0 \times \mathbb{R}.$$

Claramente X es diferenciable (de la misma clase que α) y su matriz jacobiana es

$$JX(u, v) = \begin{pmatrix} r'(u) \cos v & r'(u) \sin v & z'(u) \\ -r(u) \sin v & r(u) \cos v & 0 \end{pmatrix}.$$

El menor formado por las dos primeras columnas es $r(u)r'(u)$. Por hipótesis r no se anula y, por ser α una carta, su matriz jacobiana (r', z') no puede ser nula tampoco, luego si $r'(u) = 0$, entonces $z'(u) \neq 0$, luego uno de los menores $r(u)z'(u) \sin v$ o $-r(u)z'(u) \cos v$ es no nulo. En cualquier caso el rango de JX es 2.

Es claro que existe un $v_0 \in \mathbb{R}$ tal que $X(u_0, v_0) = (x_0, y_0, z_0)$. La aplicación X no es inyectiva, pero sí lo es su restricción a $V = V_0 \times]v_0 - \pi, v_0 + \pi[$. Veamos que es una carta para el punto dado. Sea

$$U = \{(x \cos v, x \sin v, z) \mid (x, z) \in U_0, |v - v_0| < \pi\}.$$

Sin la restricción sobre v , el conjunto U sería la antiimagen de U_0 por la aplicación continua $(x, y, z) \mapsto (\sqrt{x^2 + y^2}, z)$. En realidad U es la intersección de este abierto con el complementario del semiplano formado por los puntos $(x \cos v_0, x \sin v_0, z)$, con $x \geq 0$, que es un cerrado, luego U es abierto. Es fácil ver que $X[V] = U \cap S$. Falta probar que X^{-1} es continua, ahora bien, dado

$(x, y, z) \in U \cap S$ podemos obtener su coordenada u como $u = \alpha^{-1}(\sqrt{x^2 + y^2}, z)$, que es una aplicación continua, y su coordenada v se obtiene aplicando a $(x/r(u), y/r(u))$ la inversa del homeomorfismo $v \mapsto (\cos v, \sen v)$ definido para $|v - v_0| < \pi$. Por lo tanto $X|_V$ es una carta alrededor de (x_0, y_0, z_0) .

Si la variedad C es cubrible por una única carta $(r(u), z(u))$, lo que se traduce en que C es una curva regular, entonces tenemos una única función X , de modo que todo punto de S admite como carta a una restricción de X . Expresaremos esto diciendo simplemente que X es una carta de S .

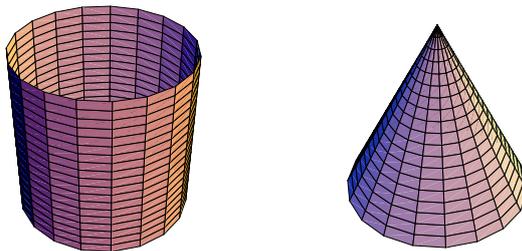
Las líneas de la forma $X(u, v_0)$ y $X(u_0, v)$, donde u_0 y v_0 son constantes, se llaman *meridianos* y *paralelos* de la superficie S . Los paralelos son siempre circunferencias paralelas entre sí, los meridianos son giros de la curva C .

Los ejemplos más simples de superficies de revolución se obtienen al girar una recta. Si ésta es paralela al eje de giro obtenemos un *cilindro*, y en caso contrario un *cono*. En el caso del cono hemos de considerar en realidad una semirrecta abierta $(r(u), z(u)) = (mu, u)$, para $u > 0$, pues para $u = 0$ tenemos el vértice del cono, donde S no es diferenciable. Una carta del cilindro es

$$X(u, v) = (r \cos v, r \sen v, u),$$

y para el cono tenemos

$$X(u, v) = (mu \cos v, mu \sen v, u).$$



Las figuras muestran algunos de los meridianos y paralelos del cilindro y el cono. Los meridianos son rectas y los paralelos circunferencias.

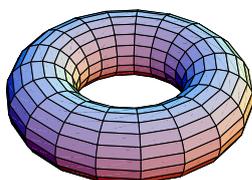
Un caso más sofisticado aparece al girar una circunferencia de radio r cuyo centro esté a una distancia R del eje Z , es decir, tomando

$$(r(u), z(u)) = (R + r \cos u, r \sen u), \quad \text{con } 0 < r < R.$$

Todo punto de la circunferencia admite como carta a una restricción de esta curva. Así obtenemos un tubo de sección circular cerrado sobre sí mismo. Recibe el nombre de *toro*.² En este caso

$$X(u, v) = (R \cos v + r \cos u \cos v, R \sen v + r \cos u \sen v, r \sen u).$$

²Del latín *torus*, que es el nombre dado en arquitectura a los salientes tubulares de las columnas. Obviamente no tiene nada que ver con *taurus*, el animal del mismo nombre en castellano.



Obviamente X es de clase C^∞ . Su restricción a $]0, 2\pi[\times]0, 2\pi[$ es inyectiva y cubre todos los puntos del toro excepto los de las circunferencias $u = 0$ y $v = 0$. Si llamamos U al complementario de la unión de estas dos circunferencias tenemos un abierto en \mathbb{R}^3 , y es claro que con él se cumple la definición de variedad. Igualmente se prueba que la restricción a $] -\pi, \pi[\times] -\pi, \pi[$ constituye una carta para los puntos exceptuados. Así pues, el toro es una superficie diferenciable de clase C^∞ . Sus meridianos son circunferencias de radio r .

La esfera menos dos puntos antípodas puede considerarse como la superficie de revolución generada por la semicircunferencia $(r \operatorname{sen} \phi, r \operatorname{cos} \phi)$, para $\phi \in]0, \pi[$. La carta correspondiente es

$$X(\phi, \theta) = (r \operatorname{sen} \phi \cos \theta, r \operatorname{sen} \phi \operatorname{sen} \theta, r \operatorname{cos} \phi), \quad \phi \in]0, \pi[, \theta \in]0, 2\pi[.$$

Si $p = X(\phi, \theta)$ entonces θ es la longitud de p en el sentido geográfico y ϕ es la “colatitud”, es decir, el ángulo respecto al polo Norte. Los meridianos y paralelos coinciden con los geográficos. La carta no cubre los polos, aunque girando la esfera obtenemos otra carta similar que los cubra. ■

Ejemplo: La cinta de Möbius Consideremos la aplicación

$$X : \mathbb{R} \times]-1/2, 1/2[\longrightarrow \mathbb{R}^3$$

dada por

$$X(u, v) = \left(\left(1 + \frac{v}{2} \cos \pi u\right) \cos 2\pi u, \left(1 + \frac{v}{2} \cos \pi u\right) \operatorname{sen} 2\pi u, \frac{v}{2} \operatorname{sen} \pi u \right).$$

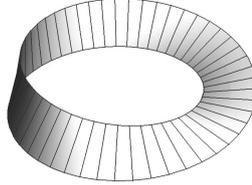
Para interpretarla observamos que la curva $X(u, 0)$ recorre la circunferencia de centro $(0, 0, 0)$ y radio 1 contenida en el plano XY , y da una vuelta completa cada vez que u aumenta 1 unidad, mientras que, para cada u fijo, la curva $X(u, v)$ es un segmento de longitud 1 perpendicular a la circunferencia, pues sus vectores tangentes respectivos son

$$\begin{aligned} X_u(u, 0) &= (-2\pi \operatorname{sen} 2\pi u, 2\pi \operatorname{cos} 2\pi u, 0), \\ X_v(u, v) &= \left(\frac{1}{2} \cos \pi u \cos 2\pi u, \frac{1}{2} \cos \pi u \operatorname{sen} 2\pi u, \frac{1}{2} \operatorname{sen} \pi u \right), \end{aligned}$$

y ciertamente son vectores ortogonales. Más en general, un cálculo rutinario muestra que la matriz jacobiana $JX(u, v)$ tiene rango 2 en cada punto, con lo que no es difícil probar que la situación es la misma que en el caso de las superficies de revolución: la imagen M de X es una variedad diferencial de clase C^∞ que admite como cartas las restricciones de X a cualquier rectángulo

$$]u_0 - 1/2, u_0 + 1/2[\times]-1/2, 1/2[.$$

La variedad M se conoce como *cinta de Möbius* y su aspecto es el que muestra la figura:



Podemos pensar que M se obtiene pegando los extremos opuestos de un rectángulo, pero no de la forma natural que da lugar a un cilindro, sino dando media vuelta a uno de los extremos antes de unirlo al otro. ■

Ejemplo: Producto de variedades Si $S_1 \subset \mathbb{R}^{m_1}$ y $S_2 \subset \mathbb{R}^{m_2}$ son variedades entonces $S_1 \times S_2 \subset \mathbb{R}^{m_1+m_2}$ es también una variedad. Si $X_1 : U_1 \rightarrow V_1 \cap S_1$ es una carta alrededor de un punto $p_1 \in S_1$ y $X_2 : U_2 \rightarrow V_2 \cap S_2$ es una carta alrededor de $p_2 \in S_2$, entonces $X_1 \times X_2 : U_1 \times U_2 \rightarrow (V_1 \times V_2) \cap (S_1 \times S_2)$ dada por $(X_1 \times X_2)(u_1, u_2) = (X_1(u_1), X_2(u_2))$ es una carta alrededor de (p_1, p_2) .

Sean $\pi_i : S_1 \times S_2 \rightarrow S_i$ las proyecciones. Si la carta X_1 tiene coordenadas x_1, \dots, x_{n_1} y la carta X_2 tiene coordenadas y_1, \dots, y_{n_2} , entonces las coordenadas de $X_1 \times X_2$ son las funciones $\pi_1 \circ x_i$ y $\pi_2 \circ y_i$, a las que podemos seguir llamando x_i e y_i sin riesgo de confusión. ■

6.2 Espacios tangentes, diferenciales

Al principio de la sección anterior anticipábamos que los sistemas de coordenadas en una variedad son un análogo a los sistemas de coordenadas en un espacio afín. La diferencia principal es que en el caso afín las coordenadas están definidas sobre todo el espacio, mientras que en una variedad las tenemos definidas sólo en un entorno de cada punto. En esta sección desarrollaremos esta analogía mostrando que toda variedad diferenciable se confunde en un entorno de cada punto con una variedad afín. Para empezar, si p es un punto de una variedad S , X es una carta alrededor de p y x es su sistema de coordenadas asociado, sabemos que en un entorno de $x(p)$ el punto $X(x)$ se confunde con $p + dX(x(p))(x - x(p))$, con lo que los puntos de S se confunden con los de $p + dX(x(p))[\mathbb{R}^n]$.

Definición 6.8 Sea $S \subset \mathbb{R}^m$ una variedad diferenciable de dimensión n y sea $X : U \rightarrow S$ una carta alrededor de un punto $p \in S$. Sea $x \in U$ tal que $X(x) = p$. Llamaremos *espacio tangente* a S en p a la variedad lineal $T_p(S) = dX(x)[\mathbb{R}^n]$. Llamaremos *variedad tangente* a S por p a la variedad afín $p + T_p(S)$.

Puesto que $JX(x)$ tiene rango n , es claro que las variedades tangentes tienen dimensión n . El teorema 6.5 prueba que el espacio tangente no depende de la

carta con la que se construye, pues si X e Y son dos cartas alrededor de p , digamos $X(x) = Y(y) = p$, sabemos que $g = X \circ Y^{-1}$ es diferenciable en un entorno de x y $X = g \circ Y$, luego $dX(x) = dg(x) \circ dY(y)$, luego $dX(x)$ y $dY(y)$ tienen la misma imagen (pues $dg(x)$ es un isomorfismo). El teorema siguiente muestra más explícitamente que $T_p(S)$ sólo depende de S .

Teorema 6.9 *Sea $S \subset \mathbb{R}^m$ una variedad diferenciable de dimensión n . Entonces $T_p(S)$ está formado por el vector nulo más los vectores tangentes en p a todas las curvas regulares que pasan por p contenidas en S .*

DEMOSTRACIÓN: Sea $X : U \rightarrow S$ una carta alrededor de p . Podemos suponer que es de la forma $X(x) = (x, f(x))$, para una cierta función diferenciable f . Sea $X(p_1) = p$. Sea $v \in \mathbb{R}^n$ no nulo. Consideremos la curva $x(t) = p_1 + tv$. Para valores suficientemente pequeños de t se cumple que $x(t) \in U$. Consideremos la curva $\alpha(t) = X(x(t))$. Claramente α está contenida en S y cumple $\alpha(0) = p$. Su vector tangente en p es

$$\alpha'(0) = dX(x(0))(x'(0)) = dX(p_1)(v).$$

Esto prueba que todo vector de $T_p(S)$ es de la forma indicada. Recíprocamente, si $\alpha(t)$ es una curva regular contenida en S que pasa por p , digamos $\alpha(t_0) = p$, sea $x(t) = X^{-1}(\alpha(t))$, definida en un entorno de t_0 . Se cumple que $x(t)$ es derivable, pues X^{-1} no es más que la restricción de la proyección $\pi : \mathbb{R}^m \rightarrow \mathbb{R}^n$, que es diferenciable, luego $x = \alpha \circ \pi$. Tenemos $\alpha = x \circ X$, luego $\alpha'(t) = dX(x(t))(x'(t))$. Esta relación prueba que $x'(t) \neq 0$ o de lo contrario también se anularía $\alpha'(t)$. Por lo tanto x es regular. Además la tangente de α en p es $\alpha'(t_0) = dX(p_1)(x'(t_0)) \in T_p(S)$. ■

En la prueba de este teorema hemos visto un hecho importante: si α es una curva contenida en una variedad S y pasa por un punto p , dada una carta $X : U \rightarrow S$ alrededor de p , podemos trasladar a la carta el arco de curva alrededor de p , es decir, existe otra curva x en U de modo que $\alpha = x \circ X$ (en un entorno de las coordenadas de p). En otras palabras, x es la representación de α en el mapa de S determinado por X .

Ejercicio: Probar que el plano tangente a una gráfica vista como variedad diferenciable coincide con el que ya teníamos definido.

Precisemos la interpretación geométrica de la variedad tangente. Ya hemos justificado que los puntos de S se confunden con los de la variedad tangente $T_p(S)$ en un entorno de p , pero más exactamente, si X es una carta alrededor de p y x es su sistema de coordenadas, hemos visto que cada punto $q \in S$ suficientemente próximo a p se confunde con el punto

$$p + dX(x(p))(x(q) - x(p)) \in p + T_p(S).$$

Definición 6.10 Sea $S \subset \mathbb{R}^m$, $p \in S$, sea $X : U \rightarrow V \cap S$ una carta alrededor de p y sea x su sistema de coordenadas. Llamaremos *proyección* asociada a X a la aplicación $\pi_p : S \cap V \rightarrow T_p(S)$ dada por $\pi_p(q) = dX(x(p))(x(q) - x(p))$.

Según hemos visto, la interpretación geométrica de estas proyecciones consiste en que el paso $q \mapsto \pi_p(q)$ es imperceptible si tomamos puntos q suficientemente próximos a p . Ahora veamos que las coordenadas de q en la carta coinciden con las coordenadas de $\pi_p(q)$ asociadas a un cierto sistema de referencia afín en $T_p(S)$.

Sea $X : U \rightarrow S$ una carta de una variedad S . Sea $X(x) = p$. Entonces $dX(x) : \mathbb{R}^n \rightarrow T_p(S)$ es un isomorfismo. Por consiguiente, si e_1, \dots, e_n son los vectores de la base canónica en \mathbb{R}^n , sus imágenes $dX(x)(e_i) = D_i X(x)$ forman una base de $T_p(S)$. El espacio tangente no tiene una base canónica pero, según acabamos de ver, cada carta alrededor de p determina una base en $T_p(S)$. Es claro que si $q \in S$ está en el entorno de p cubierto por la carta, las coordenadas de $\pi_p(q)$ en la base asociada en $T_p(S)$ son $x(q) - x(p)$, luego si con dicha base formamos un sistema de referencia afín en $p + T_p(S)$ cuyo origen sea el punto $O = p - dX(x(p))(x(p))$, tenemos que las coordenadas de $\pi_p(q)$ en este sistema son precisamente $x(q)$. Cuando hablemos del sistema de referencia afín asociado a la carta nos referiremos a éste. En conclusión, cada punto q de un entorno de p en S se confunde con el punto $\pi_p(q)$ de idénticas coordenadas afines en la variedad tangente $p + T_p(S)$.

Ejercicio: Probar que si S_1 y S_2 son variedades diferenciables y $(p, q) \in S_1 \times S_2$ entonces $T_{(p,q)}(S_1 \times S_2) = T_p(S_1) \times T_q(S_2)$.

Seguidamente generalizamos la noción de diferenciabilidad al caso de aplicaciones entre variedades cualesquiera (no necesariamente abiertos de \mathbb{R}^n).

Definición 6.11 Diremos que una aplicación continua $f : S \rightarrow T$ entre dos variedades es *diferenciable* (de clase C^q) en un punto $p \in S$ si existen cartas X e Y alrededor de p y $f(p)$ respectivamente de modo que $X \circ f \circ Y^{-1}$ sea diferenciable (de clase C^q) en $X^{-1}(p)$.

El teorema 6.5 implica que la diferenciabilidad de f en p no depende de la elección de las cartas X e Y , en el sentido de que si unas cartas prueban que f es diferenciable, otras cualesquiera lo prueban igualmente.

Es fácil ver que la composición de aplicaciones diferenciables es diferenciable. Una aplicación $f : U \rightarrow \mathbb{R}^m$ definida en un abierto U de \mathbb{R}^n es diferenciable en el sentido que ya teníamos definido si y sólo si lo es considerando a U y a \mathbb{R}^m como variedades diferenciables (con la identidad como carta).

Por el teorema 6.4, si $S \subset T \subset \mathbb{R}^m$ son variedades diferenciables, la inclusión $i : S \rightarrow T$ es diferenciable, lo que se traduce en que las restricciones a S de las funciones diferenciables en T son funciones diferenciables en S .

Es claro que todas estas propiedades valen también si sustituimos la diferenciabilidad por la propiedad de ser de clase C^q .

Una aplicación $f : S \rightarrow T$ entre dos variedades es un *difeomorfismo* si es biyectiva, diferenciable y su inversa es diferenciable. Dos variedades son *difeomorfas* si existe un difeomorfismo entre ellas.

Es obvio que las cartas de una variedad son difeomorfismos en su imagen. Más aún:

Teorema 6.12 *Todo difeomorfismo entre un abierto de \mathbb{R}^n y un abierto de una variedad S en \mathbb{R}^m es una carta para S .*

DEMOSTRACIÓN: Sea $f : U \rightarrow W$ un difeomorfismo, donde $W \subset S$ es abierto en S . Entonces existe un abierto V en \mathbb{R}^m tal que $f[U] = W = V \cap S$. Obviamente df tiene rango máximo en cada punto, con lo que se cumple la definición de carta. ■

En particular tenemos que las coordenadas $x_i : S \cap V \rightarrow \mathbb{R}$ asociadas a una carta X son funciones diferenciables (son la composición de X^{-1} con las proyecciones $\pi_i : \mathbb{R}^n \rightarrow \mathbb{R}$). Ahora definimos la diferencial de una función diferenciable.

Supongamos que $f : S \rightarrow T$ es una aplicación entre dos variedades diferenciable en un punto p . Sean $X : U \rightarrow S$ e $Y : W \rightarrow T$ cartas alrededor de p y $f(p)$. Digamos que $U(x) = p$, $Y(y) = f(p)$. Entonces $j = X \circ f \circ Y^{-1}$ es diferenciable en x y tenemos las aplicaciones lineales siguientes:

$$\begin{array}{ccc} T_p(S) & & T_{f(p)}(T) \\ dX(x) \uparrow & & \uparrow dY(y) \\ \mathbb{R}^n & \xrightarrow{dj(x)} & \mathbb{R}^m \end{array}$$

Las flechas verticales representan isomorfismos, luego podemos definir la *diferencial* de f en p como la aplicación lineal $df(p) : T_p(S) \rightarrow T_{f(p)}(T)$ dada por $df(p) = dX(x)^{-1} \circ dj(x) \circ dY(y)$. Teniendo en cuenta que las diferenciales aproximan localmente a las funciones correspondientes no es difícil convencerse de que $df(p)$ se confunde con f cuando los puntos de $T_p(S)$ se confunden con los de S . El teorema siguiente prueba que $df(p)$ no depende de la elección de las cartas X e Y .

Teorema 6.13 *Sea $f : S \rightarrow T$ una aplicación diferenciable en un punto $p \in S$. Sea $v \in T_p(S)$. Si α es cualquier curva contenida en S que pase por p con tangente v , entonces $\alpha \circ f$ es una curva contenida en T que pasa por $f(p)$ con tangente $df(p)(v)$.*

DEMOSTRACIÓN: Sean X e Y cartas alrededor de p y $f(p)$ respectivamente. Digamos que $X(x) = p$ e $Y(y) = f(p)$. Sea β la representación de α en la carta X , es decir, $\alpha = \beta \circ X$. Entonces $v = \alpha'(t_0) = dX(x)(\beta'(t_0))$.

Podemos descomponer $\alpha \circ f = \alpha \circ X^{-1} \circ X \circ f \circ Y^{-1} \circ Y$. Con la notación que hemos empleado en la definición de $df(p)$ tenemos $\alpha \circ f = \beta \circ j \circ Y$. Esto prueba que $\alpha \circ f$ es derivable en t_0 y además

$$(\alpha \circ f)'(t_0) = dY(y)(dj(x)(\beta'(t_0))) = dY(y)(dj(x)(dX(x)^{-1}(v))) = df(p)(v).$$

■

Es inmediato comprobar que la regla de la cadena sigue siendo válida para aplicaciones diferenciables entre variedades, es decir,

$$d(f \circ g)(p) = df(p) \circ dg(f(p)).$$

De aquí se sigue en particular que si f es un difeomorfismo, entonces $df(p)$ es un isomorfismo y $df^{-1}(f(p)) = df(p)^{-1}$.

Si $S \subset T \subset \mathbb{R}^m$ son variedades diferenciables entonces el teorema anterior prueba que la diferencial de la inclusión $i : S \rightarrow T$ en cada punto $p \in S$ es simplemente la inclusión de $T_p(S)$ en $T_p(T)$. De aquí se sigue que la diferencial en un punto p de la restricción a S de una función f diferenciable en T es simplemente la restricción de $df(p)$ a $T_p(S)$, pues la restricción no es más que la composición con la inclusión.

Si $f, g : S \rightarrow \mathbb{R}$ son funciones diferenciables, es fácil ver que

$$d(f + g) = df + dg, \quad d(fg) = gdf + f dg,$$

donde las operaciones hay que entenderlas como definidas puntualmente, por ejemplo, $gdf + f dg$ representa la aplicación que a cada $p \in S$ le asigna la aplicación $g(p)df(p) + f(p)dg(p) : T_pS \rightarrow \mathbb{R}$ dada por

$$v \mapsto g(p)df(p)(v) + f(p)dg(p)(v).$$

Probamos la fórmula para el producto y dejamos al lector la de la suma, que es más simple: Dado $p \in S$, tomamos una carta $X : U \rightarrow S$ alrededor de p y $x \in U$ tal que $X(x) = p$. Sean $\bar{f} = X \circ f$, $\bar{g} = X \circ g$, que son funciones diferenciables $\bar{f}, \bar{g} : U \rightarrow \mathbb{R}$. Además $X \circ (fg) = \bar{f}\bar{g}$. Por lo tanto,

$$\begin{aligned} d(fg)(p)(v) &= d(\bar{f}\bar{g})(x)(dX(x)^{-1}(v)) = (\bar{g}(x)d\bar{f}(x) + \bar{f}(x)d\bar{g}(x))(dX(x)^{-1}(v)) \\ &= g(p)d\bar{f}(x)(dX(x)^{-1}(v)) + f(p)d\bar{g}(x)(dX(x)^{-1}(v)) \\ &= g(p)df(p)(v) + f(p)dg(p)(v) = (g(p)df(p) + f(p)dg(p))(v). \end{aligned}$$

Si $X : U \rightarrow S$ es una carta de una variedad S alrededor de un punto p , entonces sus coordenadas asociadas x_i son ciertamente diferenciables. Más concretamente, si $\pi_i : \mathbb{R}^n \rightarrow \mathbb{R}$ es la proyección en la i -ésima coordenada, tenemos que $x_i = X^{-1} \circ \pi_i$, luego $dx_i(p) = dX(p)^{-1} \circ d\pi_i(x)$ y en particular

$$dx_i(p)(D_j X(x)) = d\pi_i(x)(e_j) = \begin{cases} 1 & \text{si } i = j, \\ 0 & \text{si } i \neq j, \end{cases}$$

luego las aplicaciones $dx_i(p)$ forman la base dual de $D_1 X(x), \dots, D_n X(x)$. Por consiguiente, para cada $v \in T_p(S)$ se cumple que $dx_i(p)(v)$ es la coordenada correspondiente a $D_i X(x)$ en la expresión de v como combinación lineal de las derivadas de X .

Ejemplo Consideremos el plano tangente a \mathbb{R}^2 en el punto $p = (1, 1)$ (que es el propio \mathbb{R}^2). La base asociada a la carta identidad es simplemente la base canónica (e_1, e_2) , y su base dual es la dada por las proyecciones $dx(p)$, $dy(p)$. También podemos considerar también la carta determinada por las coordenadas

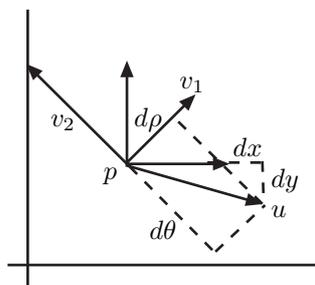
polares (ρ, θ) , es decir, $(x, y) = (\rho \cos \theta, \rho \sen \theta)$. Su base asociada es la formada por las derivadas parciales:

$$v_1(\rho, \theta) = (\cos \theta, \sen \theta), \quad v_2(\rho, \theta) = (-\rho \sen \theta, \rho \cos \theta).$$

En particular, en el punto $(1, 1)$ queda

$$v_1 = \left(\frac{1}{\sqrt{2}}, \frac{1}{\sqrt{2}} \right), \quad v_2 = (-1, 1).$$

Dado un vector $u \in \mathbb{R}^2$, sus coordenadas (en la base canónica) son $(dx(p)(u), dy(p)(u))$, mientras que $(d\rho(p)(u), d\theta(p)(u))$ son sus coordenadas en la base (v_1, v_2) .



Conocemos la relación entre las diferenciales:

$$dx = \cos \theta d\rho - \rho \sen \theta d\theta, \quad dy = \sen \theta d\rho + \rho \cos \theta d\theta.$$

Concretamente, en el punto $(1, 1)$ se cumple

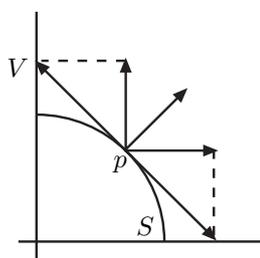
$$dx = \frac{\sqrt{2}}{2} d\rho - d\theta, \quad dy = \frac{\sqrt{2}}{2} d\rho + d\theta. \tag{6.1}$$

Sea ahora S la circunferencia de radio $\sqrt{2}$. Tres posibles cartas de S alrededor de $(1, 1)$ son

$$g_1(x) = (x, \sqrt{2-x^2}), \quad g_2(y) = (\sqrt{2-y^2}, y), \quad g_3(\theta) = \sqrt{2}(\cos \theta, \sen \theta).$$

Sus funciones coordenadas son respectivamente (las restricciones de) las funciones x, y, θ , luego sus diferenciales asociadas son las restricciones de las diferenciales correspondientes, que seguiremos llamando $dx(p), dy(p), d\theta(p)$. Es fácil ver que las bases asociadas a las tres cartas son respectivamente

$$v_x = (1, -1), \quad v_y = (-1, 1), \quad v_\theta = (-1, 1).$$



Obviamente no podemos tomar a ρ como coordenada, pues ρ es constante en S . Esto se traduce en que $d\rho(p) = 0$ (sobre $T_p(S)$). Alternativamente, vemos que los vectores de $T_p(S)$ tienen nula la primera coordenada de su expresión en la base (v_1, v_2) . Como consecuencia, de (6.1) se sigue ahora que $dy = d\theta = -dx$. ■

Ejemplo Consideremos el toro T de carta

$$X(u, v) = (R \cos v + r \cos u \cos v, R \sen v + r \cos u \sen v, r \sen u).$$

Ya hemos comentado que X no es exactamente una carta de T , sino que las cartas de T son restricciones de X a dominios adecuados. Consideremos

la circunferencia unidad $S^1 = \{x \in \mathbb{R}^2 \mid \|x\| = 1\}$. Entonces la aplicación $f : S^1 \times S^1 \rightarrow T$ dada por

$$f(x, y) = (Ry_1 + rx_1y_1, Ry_2 + rx_1y_2, rx_2)$$

es un difeomorfismo. Notemos que si $x = (\cos u, \sin u)$, $y = (\cos v, \sin v)$, entonces $f(x, y) = X(u, v)$. Teniendo esto en cuenta es fácil ver que f es biyectiva. Además es diferenciable porque sus funciones coordenadas son polinómicas (es la restricción de una función diferenciable en \mathbb{R}^4). En un entorno de cada punto de T , la función f^{-1} puede expresarse como $(\cos u, \sin u, \cos v, \sin v)$, donde u, v son las funciones coordenadas de la carta de T alrededor de punto obtenida por restricción de X . Por consiguiente f es un difeomorfismo. ■

Ejercicio: Probar que un cilindro es difeomorfo al producto de un segmento por una circunferencia y que una bola abierta menos su centro es difeomorfa al producto de un segmento por una esfera.

Definición 6.14 Sea $f : S \rightarrow \mathbb{R}$ una función definida sobre una variedad y sea $p \in S$ un punto donde f sea diferenciable. Sea X una carta de S alrededor de p y sean x_1, \dots, x_n sus coordenadas asociadas. Definimos la *derivada parcial* de f respecto a x_i en p como

$$\frac{\partial f}{\partial x_i}(p) = df(p)(D_i X(x)),$$

donde x es el vector de coordenadas de p en la carta dada.

Es claro que esta noción de derivada parcial generaliza a la que ya teníamos para el caso de funciones definidas en abiertos de \mathbb{R}^n . En el caso general sea $j = X \circ f$. Según la definición de $df(p)$ resulta que

$$\frac{\partial f}{\partial x_i}(p) = dj(x)(e_i) = \frac{\partial j}{\partial x_i}(x),$$

donde e_i es el i -ésimo vector de la base canónica de \mathbb{R}^n . Si f es diferenciable en un entorno de p tenemos

$$\frac{\partial f}{\partial x_i} = X^{-1} \circ \frac{\partial j}{\partial x_i}.$$

Ahora es claro que una función f es de clase C^q en S si y sólo si tiene derivadas parciales continuas de orden q .

Puesto que $dx_1(p), \dots, dx_n(p)$ es la base dual de $D_1 X(x), \dots, D_n X(x)$, de la propia definición de derivada parcial se sigue que

$$df = \frac{\partial f}{\partial x_1} dx_1 + \dots + \frac{\partial f}{\partial x_n} dx_n.$$

También es fácil ver que las reglas usuales de derivación de sumas y productos siguen siendo válidas, así como el teorema de Schwarz. Además

$$\frac{\partial x_j}{\partial x_i} = \begin{cases} 1 & \text{si } i = j \\ 0 & \text{si } i \neq j. \end{cases}$$

Es importante observar que la derivada de una función f respecto a una coordenada x_i no depende sólo de f y x_i , sino de la carta de la cual forma parte x_i . Por ejemplo, si en la esfera de centro 0 y radio 1 consideramos un punto cuyas tres coordenadas (x, y, z) sean no nulas, en un entorno podemos considerar la carta de coordenadas (x, y) , respecto a la cual

$$\frac{\partial z}{\partial x} = -\frac{x}{\sqrt{1-x^2-y^2}}.$$

Sin embargo, también podemos considerar la carta de coordenadas (x, z) y entonces resulta que

$$\frac{\partial z}{\partial x} = 0.$$

6.3 La métrica de una variedad

Todas las propiedades métricas de \mathbb{R}^n se derivan de su producto escalar, que es una forma bilineal $\mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$. En una variedad $S \subset \mathbb{R}^m$ no tenemos definido un producto escalar, pero sí tenemos uno en cada uno de sus espacios tangentes: la restricción del producto escalar en \mathbb{R}^m .

Conviene introducir ciertos hechos básicos sobre formas bilineales. Puesto que son puramente algebraicas las enunciaremos para un espacio vectorial arbitrario E , pero en la práctica E será siempre el espacio tangente $T_p(S)$ de una variedad S en un punto p . Fijada una base (v_1, \dots, v_n) de E , representaremos su base dual por (dx_1, \dots, dx_n) . Esta notación —puramente formal en un principio— se ajusta al único ejemplo que nos interesa, pues si $E = T_p(S)$ y (v_1, \dots, v_n) es la base asociada a una carta X , entonces la base dual que en general hemos llamado (dx_1, \dots, dx_n) es concretamente la formada por las diferenciales $dx_1(p), \dots, dx_n(p)$, donde x_1, \dots, x_n son las funciones en S que a cada punto le asignan sus coordenadas respecto a X .

Definición 6.15 Sea E un espacio vectorial de dimensión n . Llamaremos $B(E)$ al conjunto de todas las formas bilineales $F : E \times E \rightarrow \mathbb{R}$, que es claramente un espacio vectorial con la suma y el producto definidos puntualmente.³

Si $f, g : E \rightarrow \mathbb{R}$ son aplicaciones lineales, definimos su *producto tensorial* como la forma bilineal $f \otimes g \in B(E)$ dada por $(f \otimes g)(u, v) = f(u)g(v)$.

Las propiedades siguientes son inmediatas:

a) $f \otimes (g + h) = f \otimes g + f \otimes h, \quad (f + g) \otimes h = f \otimes h + g \otimes h.$

b) $(\alpha f) \otimes g = f \otimes (\alpha g) = \alpha(f \otimes g), \quad \text{para } \alpha \in \mathbb{R}.$

Teorema 6.16 *Todo elemento de $B(E)$ se expresa de forma única como*

$$F = \sum_{i,j=1}^n \alpha_{ij} dx_i \otimes dx_j, \quad \text{con } \alpha_{ij} \in \mathbb{R}.$$

³Los elementos de $B(E)$ se llaman *tensores dos veces covariantes*, pero aquí no vamos a entrar en el cálculo tensorial.

Concretamente $\alpha_{ij} = F(v_i, v_j)$.

DEMOSTRACIÓN: Basta observar que

$$(dx_i \otimes dx_j)(v_r, v_s) = \begin{cases} 1 & \text{si } i = r, j = s \\ 0 & \text{en caso contrario.} \end{cases}$$

De aquí se sigue que F y el miembro derecho de la igualdad actúan igual sobre todos los pares de vectores básicos. La unicidad es clara. ■

Por ejemplo, en estos términos el producto escalar en \mathbb{R}^n viene dado por

$$dx_1 \otimes dx_1 + \cdots + dx_n \otimes dx_n.$$

Definición 6.17 Un *campo tensorial* (dos veces covariante) en una variedad $S \subset \mathbb{R}^m$ es una aplicación que a cada $p \in S$ le hace corresponder una forma bilineal en $T_p(S)$. El *tensor métrico* de S es el campo g que a cada punto p le asigna la restricción a $T_p(S)$ del producto escalar en \mathbb{R}^m .

Si llamamos $T(S)$ al conjunto de los campos tensoriales en S según la definición anterior, es claro que se trata de un espacio vectorial con las operaciones definidas puntualmente. Más aún, podemos definir el producto de una función $f : S \rightarrow \mathbb{R}$ por un campo $F \in T(S)$ como el campo $fF \in T(S)$ dado por $(fF)(p) = f(p)F(p)$.

Claramente, si $x_i : S \rightarrow \mathbb{R}$ son las restricciones de las funciones coordenadas de \mathbb{R}^m , una expresión para el tensor métrico es

$$dx_1 \otimes dx_1 + \cdots + dx_m \otimes dx_m,$$

pero es más conveniente expresar el tensor métrico en términos de las coordenadas de una carta de S y no de una carta de \mathbb{R}^m . Veamos cuál es esta expresión:

Sea $X : U \rightarrow S$ una carta de S . Representaremos por x_1, \dots, x_n las funciones coordenadas respecto a X . Si $x \in U$ y $p = X(x)$, sabemos que $D_1X(x), \dots, D_nX(x)$ es una base de $T_p(S)$ y $dx_1(p), \dots, dx_n(p)$ es su base dual. Por consiguiente, todo $w \in T_p(S)$ se expresa como

$$w = dx_1(p)(w)D_1X(x(p)) + \cdots + dx_n(p)(w)D_nX(x(p)).$$

Así pues, si $w_1, w_2 \in T_p(S)$, su producto escalar es

$$g_p(w_1, w_2) = \sum_{i,j=1}^n D_iX(x(p))D_jX(x(p))dx_i(p)(w_1)dx_j(p)(w_2),$$

luego

$$g_p = \sum_{i,j=1}^n g_{ij}(p)dx_i(p) \otimes dx_j(p), \quad \text{con } g_{ij}(p) = D_iX(x(p))D_jX(x(p)),$$

o, más brevemente, como igualdad de campos:

$$g = \sum_{i,j=1}^n g_{ij} dx_i \otimes dx_j, \quad (6.2)$$

Esta expresión recibe el nombre de *expresión en coordenadas* del tensor métrico de S en la carta X . Las funciones g_{ij} se llaman *coeficientes* del tensor métrico en la carta dada. Claramente son funciones diferenciables. Notemos que la expresión coordenada no está definida en toda la variedad S , sino sólo sobre los puntos del rango V de la carta X .

La matriz $(g_{ij}(p))$ es la matriz del producto escalar de $T_p(S)$ en una cierta base. Es claro entonces que su determinante es no nulo. Este hecho será relevante en varias ocasiones.

A través del difeomorfismo $X : U \rightarrow V$ juntamente con los isomorfismos $dX(x) : \mathbb{R}^n \rightarrow T_p(S)$ podemos transportar la restricción a V del tensor métrico de S hasta un campo tensorial en U , concretamente el dado por

$$\begin{aligned} h_X(w_1, w_2) &= g_p(dX(x)(w_1), dX(x)(w_2)) \\ &= \sum_{i,j=1}^n g_{ij}(X(x)) dx_i(X(x))(dX(x)(w_1)) dx_j(X(x))(dX(x)(w_2)) \\ &= \sum_{i,j=1}^n g_{ij}(X(x)) d(X \circ x_i)(x)(w_1) d(X \circ x_j)(x)(w_2) \\ &= \sum_{i,j=1}^n g_{ij}(x) dx_i(x)(w_1) dx_j(x)(w_2), \end{aligned}$$

donde en el último término x_i es simplemente la proyección en la i -ésima coordenada de U y $g_{ij}(x) = (X \circ g_{ij})(x)$. Por lo tanto h_X tiene la misma expresión (6.2) interpretando convenientemente las funciones.

Al transportar a la carta el tensor métrico, podemos calcular el producto de dos vectores tangentes a dos curvas α y β que se cortan en p a partir de sus representaciones en X . Digamos que $\alpha(t) = X(x(t))$ y $\beta(t) = X(\bar{x}(t))$ y supongamos que en t_0 pasan por p . Entonces

$$g_p(\alpha'(t_0), \beta'(t_0)) = h_X(x'(t_0), \bar{x}'(t_0)).$$

Del mismo modo que el tensor métrico de una variedad S asigna a cada punto p el producto escalar de $T_p(S)$, también podemos considerar la aplicación que a cada punto p le asigna la norma en $T_p(S)$. Ésta recibe el nombre de *elemento de longitud* de S y se representa por ds . Así pues,

$$ds(p)(v) = \|v\| = \sqrt{g_p(v, v)}.$$

El tensor métrico y el elemento de longitud se determinan mutuamente por la relación

$$ds^2(p)(u + v) = ds^2(p)(u) + ds^2(p)(v) + 2g_p(u, v),$$

luego en la práctica es equivalente trabajar con uno o con otro y ds suele dar lugar a expresiones más simples. Por ejemplo, la expresión de ds^2 en una carta es

$$ds^2 = \sum_{i,j=1}^n g_{ij} du_i du_j. \quad (6.3)$$

La misma expresión es válida para el campo que resulta de transportarlo al dominio de la carta interpretando adecuadamente las funciones.

El nombre de elemento de longitud se debe a que si $\alpha : [a, b] \rightarrow S$ es una curva regular cuya imagen está contenida en el rango de una carta X y $\alpha(t) = X(x(t))$, entonces $ds^2(x'(t)) = \|\alpha'(t)\|^2$, luego la longitud de α es

$$\begin{aligned} L &= \int_a^b \|\alpha'(t)\| dt = \int_a^b \sqrt{\sum_{i,j=1}^n g_{ij}(x(t)) x'_i(t) x'_j(t)} dt \\ &= \int_a^b \sqrt{\sum_{i,j=1}^n g_{ij}(x(t)) x'_i(t) dt x'_j(t) dt} = \int_a^b ds, \end{aligned}$$

entendiendo ahora que en (6.3) $x = x(t)$ y $dx_i = x'_i(t) dt$.

En el caso de una superficie $S \subset \mathbb{R}^3$ es costumbre representar las derivadas parciales de una carta $X(u, v)$ mediante X_u , X_v y los coeficientes del tensor métrico como $E = X_u X_u$, $F = X_u X_v$, $G = X_v X_v$, de modo que la expresión en coordenadas del tensor métrico es

$$E du \otimes du + F(du \otimes dv + dv \otimes du) + G dv \otimes dv. \quad (6.4)$$

El elemento de longitud es

$$ds^2 = E du^2 + 2F dudv + G dv^2.$$

Ejemplo Vamos a calcular los coeficientes del tensor métrico de la superficie de revolución dada por

$$X = (r(u) \cos v, r(u) \sin v, z(u)).$$

Tenemos

$$\begin{aligned} X_u &= (r'(u) \cos v, r'(u) \sin v, z'(u)) \\ X_v &= (-r(u) \sin v, r(u) \cos v, 0), \end{aligned}$$

luego

$$E = r'(u)^2 + z'(u)^2, \quad F = 0, \quad G = r(u)^2.$$

Observemos que E es el módulo al cuadrado de la curva que genera la superficie, luego si su parametrización es la natural tenemos simplemente $E = 1$.

En el caso del toro tenemos $(r(u), z(u)) = (R + r \cos u, r \sin u)$, luego

$$E = r^2, \quad F = 0, \quad G = (R + r \cos u)^2.$$

Por lo tanto la longitud de una curva que sobre la carta venga dada por $(u(t), v(t))$ se calcula integrando

$$ds^2 = r^2 du^2 + (R + r \cos u)^2 dv^2.$$

Por ejemplo, la longitud de un arco de paralelo (u_0, t) , donde $t \in [0, k]$ es

$$\int_0^k ds = \int_0^k (R + r \cos u_0) dt = (R + r \cos u_0)k,$$

como era de esperar, dado que el paralelo es un arco de circunferencia de radio $R + r \cos u_0$. ■

Ejercicio: Calcular los coeficientes del tensor métrico del cilindro, el cono y la esfera.

Definición 6.18 Diremos que un difeomorfismo $f : S \rightarrow T$ entre dos variedades es una *isometría* si para todo arco α contenido en S se cumple que $\alpha \circ f$ tiene la misma longitud.⁴

Explícitamente, si f es una isometría y $\alpha : [a, b] \rightarrow S$ es un arco y $\alpha(t_0) = p$, entonces

$$\int_a^t \|\alpha'(x)\| dx = \int_a^t \|(\alpha \circ f)'(x)\| dx,$$

y derivando resulta

$$\|\alpha'(t_0)\| = \|(\alpha \circ f)'(t_0)\| = \|df(p)(\alpha'(t_0))\|.$$

Ahora bien, todo vector no nulo de $T_p(S)$ es de la forma $\alpha'(t_0)$ para un cierto arco α , luego tenemos que $df(p) : T_p(S) \rightarrow T_{f(p)}(T)$ es una isometría para todo punto p . Igualmente se prueba el recíproco.

Ejercicio: Probar que las isometrías de \mathbb{R}^n en \mathbb{R}^n en el sentido que acabamos de definir coinciden con las isometrías en el sentido del álgebra lineal.

Si f es una isometría, X es una carta alrededor de p con $X(x) = p$ y llamamos $Y = X \circ f$, es claro que Y es una carta alrededor de $f(p)$. Además tenemos que $D_i Y(x) = dY(x)(e_i) = df(p)(dX(x)(e_i)) = df(p)(D_i X(x))$, de donde se sigue que los coeficientes del tensor métrico son iguales en ambas cartas, es decir,

$$g_{ij}(x) = D_i X(x) D_j X(x) = D_i Y(x) D_j Y(x).$$

Similarmente se concluye que si dos variedades tienen cartas con un mismo dominio y con los mismos coeficientes g_{ij} del tensor métrico entonces los fragmentos de superficie cubiertos por las cartas son superficies isométricas.

⁴Para que esta definición tenga sentido tanto α como $\alpha \circ f$ tienen que ser rectificables. En la práctica podemos suponer que f es de clase C^1 y restringir la definición a curvas de clase C^1 , pues entonces tanto α como $\alpha \circ f$ son de clase C^1 , por lo que $\|\alpha'\|$ y $\|(\alpha \circ f)'\|$ son funciones continuas en $[a, b]$, y demostraremos (teorema 8.59) que toda función continua en un intervalo cerrado tiene una primitiva.

Ejemplo Consideremos la carta del cilindro dada por

$$X(u, v) = \left(r \cos \frac{v}{r}, r \sin \frac{v}{r}, u \right).$$

El elemento de longitud del cilindro es, en esta carta, $ds^2 = du^2 + dv^2$, que es exactamente la misma que la del plano con la identidad como carta. La aplicación X no es una isometría porque no es biyectiva, pero sí es una isometría local, en el sentido de que todo punto del plano tiene un entorno V de modo que la restricción de X es una isometría entre U y $X[U]$. Así pues, un cilindro es localmente isométrico a un plano. ■

6.4 Geodésicas

Imaginemos la superficie S de un planeta cuyos habitantes creen que es plano. Cuando éstos creen caminar en línea recta en realidad sus trayectorias son curvas, sin embargo su distinción entre rectas y curvas tiene un significado objetivo. Tratemos de explicitarlo. Sea $N_p(S)$ el *espacio normal* a S en p , es decir, el complemento ortogonal de $T_p(S)$. Consideremos una curva α contenida en S . Entonces $\alpha'(t) \in T_{\alpha(t)}(S)$. Podemos descomponer $\alpha''(t) = v_t(t) + v_n(t)$, donde $v_t(t) \in T_{\alpha(t)}(S)$ y $v_n(t) \in N_{\alpha(t)}(S)$. La descomposición es única. El vector v_n contiene la parte de la aceleración que mantiene a los habitantes del planeta pegados a su superficie (la gravedad) y es “invisible” para ellos, pues si el planeta fuera realmente plano la gravedad no curvaría sus trayectorias. El vector v_t contiene la variación de la velocidad que ellos detectan: determina si la trayectoria se curva a la izquierda o a la derecha. Ellos llaman rectas a las curvas que cumplen $v_t = 0$. A continuación desarrollamos estas ideas en un contexto más general:

Definición 6.19 Sea $S \subset \mathbb{R}^m$ una variedad de dimensión n , sea $\alpha : I \rightarrow S$ una curva regular y $V : I \rightarrow \mathbb{R}^m$ una función de clase C^1 tal que para todo $t \in I$ se cumpla $V(t) \in T_{\alpha(t)}(S)$. En estas condiciones diremos que V es un *campo de vectores* sobre α . Llamaremos *derivada covariante* de V en cada punto t a la proyección ortogonal de $V'(t)$ sobre $T_{\alpha(t)}(S)$. La representaremos por $DV(t)$.

En la situación que describíamos antes, el vector v_t es la derivada covariante del campo dado por $V(t) = \alpha'(t)$. Para ilustrar el caso general podemos pensar en un habitante del planeta S que camina rumbo norte con su brazo derecho apuntando hacia el noreste. Si interpretamos el brazo como un campo de vectores sobre su trayectoria, desde el punto de vista del caminante éste apunta siempre en la misma dirección, pues él camina “recto”, es decir, sin desviarse ni hacia el este ni hacia el oeste, y su brazo forma un ángulo fijo con su dirección de avance. En otras palabras, considera que el campo vectorial es constante y su derivada es nula. Esto es falso, pues en realidad su trayectoria no es recta, sino una circunferencia y su brazo sí cambia de dirección (el único caso en que la dirección no variaría sería si apuntara al este o al oeste, con lo que siempre marcaría la dirección perpendicular al plano de la circunferencia en que se

mueve). La que en realidad es nula es la derivada covariante del campo, que los habitantes confunden con la derivada total al desconocer la curvatura de su planeta.

En las condiciones de la definición anterior, sea $X : U \rightarrow S$ una carta de S y expresemos la curva (localmente) como $\alpha(t) = X(x(t))$. Entonces una base de $T_{\alpha(t)}(S)$ en cada punto es

$$D_1X(x(t)), \dots, D_nX(x(t)),$$

luego podremos expresar

$$V(t) = a_1(t)D_1X(x(t)) + \dots + a_n(t)D_nX(x(t)), \quad (6.5)$$

para ciertas funciones $a_i(t)$. Multiplicando la igualdad por $D_iX(x(t))$ se obtiene un sistema de ecuaciones lineales con coeficientes $g_{ij}(x(t))$. Como el determinante es no nulo, resolviendo el sistema concluimos que las funciones $a_i(t)$ son derivables. Entonces

$$V'(t) = \sum_{i=1}^n a'_i(t)D_iX(x(t)) + \sum_{i,j=1}^n a_i(t)D_{ij}X(x(t))x'_j(t). \quad (6.6)$$

El primer término es tangente a S , luego no se altera al tomar la proyección ortogonal. Para calcular la proyección del segundo conviene introducir un nuevo concepto:

Definición 6.20 Sea $X : U \rightarrow S$ una carta de una variedad S . Llamaremos *símbolos de Christoffel* de S en la carta X a las funciones $\Gamma_{ij}^k : U \rightarrow \mathbb{R}$ que cumplen

$$D_{ij}X = \sum_{k=1}^n \Gamma_{ij}^k D_kX + N_{ij}, \quad (6.7)$$

donde $N_{ij}(x) \in N_p(S)$ (con $p = X(x)$). Observemos que $\Gamma_{ij}^k = \Gamma_{ji}^k$.

Las proyecciones de las segundas parciales $D_{ij}X$ se obtienen eliminando la componente N_{ij} , con lo que al calcular la proyección de (6.6) llegamos a que la derivada covariante de V viene dada por

$$DV = \sum_{k=1}^n \left(a'_k + \sum_{i,j=1}^n a_i \Gamma_{ij}^k x'_j \right) D_kX. \quad (6.8)$$

Un hecho muy importante es que los símbolos de Christoffel, y por consiguiente la derivada covariante, dependen únicamente de los coeficientes g_{ij} de la primera forma fundamental de S . En efecto, multiplicando las ecuaciones (6.7) por D_lX obtenemos

$$D_{ij}X D_lX = \sum_{k=1}^n g_{kl} \Gamma_{ij}^k.$$

Una simple comprobación nos da que

$$D_{ij}X D_lX = \frac{1}{2}(D_i g_{jl} + D_j g_{il} - D_l g_{ij}),$$

luego en total resulta

$$\sum_{k=1}^n g_{kl} \Gamma_{ij}^k = \frac{1}{2}(D_i g_{jl} + D_j g_{il} - D_l g_{ij}). \quad (6.9)$$

Fijando i, j y variando l obtenemos un sistema de n ecuaciones lineales con n incógnitas y coeficientes (g_{kl}) , que nos permite despejar los símbolos Γ_{ij}^k en términos de los coeficientes g_{ij} y sus derivadas, como queríamos probar. Ahora nos ocupamos con detalle del caso particular que describíamos al principio de la sección:

Definición 6.21 Sea $\alpha(t)$ una curva contenida en una variedad S . Llamaremos *aceleración geodésica*⁵ de α a la derivada covariante del campo vectorial α' .

Supongamos que α está parametrizada por el arco. Entonces $\|\alpha'(s)\| = 1$, luego derivando resulta $\alpha''(s)\alpha'(s) = 0$, y esta ortogonalidad se conserva al proyectar sobre $T_p(S)$, de modo que $D\alpha'(s)$ es perpendicular al vector tangente de α . Llamaremos *curvatura geodésica* de α a $\kappa_g = \|D\alpha'\|$. Si $\kappa_g \neq 0$ definimos el *vector normal geodésico* de α como el vector $\kappa_g^{-1}D\alpha'$, de modo que $D\alpha' = \kappa_g n_g$.

En el caso de que α no esté parametrizada por el arco el vector normal geodésico y la curvatura geodésica se definen a través de su parametrización natural. Explícitamente, si $\alpha(t)$ es una curva contenida en S y $s(t)$ es su longitud de arco, usando la notación $v = s'(t) = \|\alpha'(t)\|$, $a = v'(t)$ para la velocidad y aceleración sobre la trayectoria y $T = \alpha'(s)$ para el vector tangente, tenemos

$$\alpha'(t) = v\alpha'(s), \quad \alpha''(t) = aT + v^2\alpha''(s).$$

Al proyectar sobre el espacio tangente resulta

$$D\alpha'(t) = aT + v^2\kappa_g n_g.$$

De este modo, la aceleración geodésica de α se descompone en una aceleración tangencial, cuyo módulo a es la tasa de variación de la velocidad v , y una aceleración normal, cuyo módulo es $v^2\kappa_g$. Un habitante del planeta S que “crea” vivir en $T_p(S)$ confundirá la aceleración geodésica, el vector normal geodésico y la curvatura geodésica de α con la aceleración, el vector normal y la curvatura de α . Por lo tanto llamará rectas a las curvas sin aceleración normal geodésica:

Definición 6.22 Una curva α contenida en una variedad S es una *geodésica*⁶ si cumple $\kappa_g = 0$, o equivalentemente, si $D\alpha'$ es proporcional a α' en cada punto. En tal caso el módulo de $D\alpha'$ es $a = s''(t)$, donde s es la longitud de arco, por lo que si α está parametrizada por el arco entonces α es una geodésica si y sólo si $D\alpha' = 0$.

⁵La *geodesia* (gr. = división de la tierra) estudia la forma de la Tierra, deducida a partir de mediciones realizadas desde su superficie. La geometría diferencial ha adoptado este adjetivo para referirse en general a los conceptos que puede medir un “habitante” de una variedad arbitraria sin salir de ella.

⁶Deberíamos decir “recta geodésica”, es decir, el equivalente en S a una recta, pero es preferible contraer el término pues, al fin y al cabo, normalmente las geodésicas no son rectas.

Vamos a particularizar las ecuaciones que determinan la derivada covariante de un campo al caso de la aceleración geodésica de una curva. Si $\alpha(t) = X(x(t))$, entonces

$$\alpha'(t) = \sum_{i=1}^n D_i X(x(t)) x'_i(t),$$

luego si en (6.5) hacemos $V = \alpha'$ tenemos $a_i = x'_i$, luego la fórmula (6.8) se convierte en

$$DV = \sum_{k=1}^n \left(x''_k + \sum_{i,j=1}^n \Gamma_{ij}^k x'_i x'_j \right) D_k X.$$

El vector

$$\left(x''_k + \sum_{i,j=1}^n \Gamma_{ij}^k x'_i x'_j \right)_{k=1}^n$$

es la antiimagen por dX de $D\alpha'$, es decir, la representación en el mapa de la aceleración geodésica de α . Lo llamaremos *expresión en coordenadas* de dicha aceleración geodésica.

La condición necesaria y suficiente para que una curva parametrizada por el arco de coordenadas $x(s)$ sea una geodésica es

$$x''_k + \sum_{i,j=1}^n \Gamma_{ij}^k x'_i x'_j = 0, \quad k = 1, \dots, n. \quad (6.10)$$

Si la parametrización es arbitraria sólo hemos de exigir que el vector formado por los miembros izquierdos sea proporcional a x' .

Ejemplo Si una carta $X(u, v)$ de una superficie $S \subset \mathbb{R}^3$ cumple $F = 0$, las ecuaciones (6.9) se reducen a

$$\begin{aligned} \Gamma_{11}^1 &= \frac{E_u}{2E}, & \Gamma_{11}^2 &= -\frac{E_v}{2G}, & \Gamma_{12}^1 &= \frac{E_v}{2E}, \\ \Gamma_{12}^2 &= \frac{G_u}{2G}, & \Gamma_{22}^1 &= -\frac{G_u}{2E}, & \Gamma_{22}^2 &= \frac{G_v}{2G}. \end{aligned} \quad (6.11)$$

■

Ejemplo En un plano (tomando como carta la identidad) todos los símbolos de Christoffel son nulos, por lo que las geodésicas parametrizadas por el arco son las curvas que cumplen $(u'', v'') = (0, 0)$, es decir, las rectas. ■

Ejemplo En la superficie de revolución generada por la curva $(r(u), z(u))$, suponiendo a ésta parametrizada por el arco, los únicos símbolos de Christoffel no nulos son

$$\Gamma_{12}^2 = \frac{r'(u)}{r(u)}, \quad \Gamma_{22}^1 = -r(u)r'(u).$$

Por lo tanto las ecuaciones de las geodésicas parametrizadas por el arco son

$$u'' = v'^2 r(u) r'(u), \quad v'' = -2u' v' \frac{r'(u)}{r(u)}.$$

Es inmediato comprobar que los meridianos (t, v_0) cumplen estas ecuaciones, luego son geodésicas. Si se cumple $r'(u) = 0$, (por ejemplo en los extremos locales de r) entonces el paralelo (u_0, t) también cumple las ecuaciones, luego es una geodésica.

En el caso concreto de la esfera los meridianos son los arcos de circunferencia de radio máximo que unen los polos. Dada la simetría de la esfera, que permite tomar cualquier par de puntos antípodas como polos, podemos afirmar que todas las circunferencias máximas son geodésicas. Para una carta dada, el único paralelo (u_0, t) que cumple $r'(u_0) = 0$ es el ecuador de la esfera, que también es una circunferencia máxima, luego ya sabíamos que es una geodésica. ■

6.5 Superficies

Terminaremos el capítulo con algunos resultados específicos sobre superficies $S \subset \mathbb{R}^3$. Si X es una carta de una superficie S , entonces X_u, X_v son en cada punto (u, v) una base del plano tangente en $X(u, v)$, luego el vector $X_u \times X_v$ es no nulo y perpendicular a dicho plano. Si llamamos α al ángulo formado por X_u y X_v en un punto dado, entonces

$$\|X_u \times X_v\|^2 = \|X_u\|^2 \|X_v\|^2 (1 - \cos^2 \alpha) = X_u X_u X_v X_v - (X_u X_v)^2 = EG - F^2.$$

Así pues, $\|X_u \times X_v\| = \sqrt{EG - F^2}$.

Definición 6.23 La *aplicación de Gauss* asociada a una carta $X : U \rightarrow S$ de una superficie S es la aplicación $n : U \rightarrow \mathbb{R}^3$ dada por

$$n(u, v) = \frac{X_u \times X_v}{\|X_u \times X_v\|} = \frac{X_u \times X_v}{\sqrt{EG - F^2}}.$$

De este modo, $n(u, v)$ es en cada punto un vector unitario perpendicular a S en $X(u, v)$. Esto lo determina completamente salvo en su sentido. Si cambiamos de carta, el sentido de n puede cambiar.

Si X y \bar{X} son dos cartas que cubren una misma región conexa de una variedad, entonces $n(u, v) = \epsilon(u, v) \bar{n}(u, v)$, donde $\epsilon(u, v) = \pm 1$. Es claro que ϵ es una función continua en un conexo, luego ha de ser constante. En definitiva, $n(u, v) = \pm \bar{n}(u, v)$. Resulta, pues, que en un entorno de cada punto de S existen exactamente dos determinaciones opuestas del vector normal. A cualquiera de ellas la llamaremos también *aplicación de Gauss* de la superficie.

Si llamamos $G \subset S$ a la imagen de X , la aplicación n induce otra aplicación $n : G \rightarrow S^2$, donde S^2 es la esfera de centro $(0, 0, 0)$ y radio 1, que está unívocamente determinada en un entorno de cada punto excepto por su signo. Es importante notar que no siempre es posible extender esta aplicación n a toda la superficie S (sin perder la continuidad).

Ejemplo Consideremos de nuevo la cinta de Möbius M . Un simple cálculo muestra que, sobre la circunferencia $X(u, 0)$, se cumple que

$$X_u \times X_v = (\pi \cos 2\pi u \operatorname{sen} \pi u, \pi \operatorname{sen} 2\pi u \operatorname{sen} \pi u, -\pi \cos \pi u)$$

cuyo módulo es π , por lo que

$$n(u, 0) = (\cos 2\pi u \operatorname{sen} \pi u, \operatorname{sen} 2\pi u \operatorname{sen} \pi u, -\cos \pi u)$$

Si se pudiera definir una aplicación continua $N : M \rightarrow S^2$ que fuera perpendicular a M en cada punto, tendría que ser $n(u, 0) = \alpha(u)N(X(u, 0))$, donde $\alpha(u) = n(u, 0) \cdot N(X(u, 0))$ sería una función continua igual a ± 1 en cada $u \in \mathbb{R}$, luego sería constante. Sin embargo,

$$n(0, 0) = (0, 0, -1), \quad n(1, 0) = (0, 0, 1),$$

mientras que $X(0, 0) = X(1, 0) = (1, 0, 0)$, luego $N(1, 0, 0)$ tendría que ser a la vez dos vectores opuestos.

Esto significa que si desplazamos de forma continua un vector normal a M hasta dar una vuelta completa a la cinta, terminamos con el vector apuntando en sentido opuesto al vector de partida, por lo que no es posible asignar de forma continua un vector normal a cada punto.

No vamos a entrar en ello aquí, pero esto se expresa diciendo que la cinta de Möbius es una variedad diferencial *no orientable*. ■

La aplicación de Gauss aporta información importante sobre las superficies y simplifica algunos de los conceptos que hemos estudiado para variedades arbitrarias. Por ejemplo, en la sección anterior hemos estudiado la componente tangencial (o geodésica) de la curvatura de una curva contenida en una variedad. Del mismo modo podemos definir la curvatura normal como el módulo de la componente normal de la segunda derivada. En el caso de las superficies en \mathbb{R}^3 podemos apoyarnos en la aplicación de Gauss.

Definición 6.24 Sea S una superficie (al menos de clase C^2) y α una curva contenida en S parametrizada por el arco y que pase por un punto p . Fijada una determinación n del vector normal a S alrededor de p , llamaremos *curvatura normal* de α a $\kappa_n = \alpha'' \cdot n$. Definimos $N_n = \kappa_n n$ y $N_t = \alpha'' - N_n$.

Notemos que el signo de κ_n depende de la determinación que elijamos de la aplicación de Gauss. Supongamos que sobre una carta la curva es $(u(t), v(t))$. Entonces

$$\alpha' = X_u u' + X_v v', \quad \alpha'' = X_{uu} u'^2 + X_{uu} u'' + X_{uv} u' v' + X_{uv} u' v' + X_{vv} v'^2 + X_v v'',$$

luego

$$\kappa_n = \alpha'' \cdot n = (X_{uu} n) u'^2 + 2(X_{uv} n) u' v' + (X_{vv} n) v'^2.$$

Llamamos

$$e = X_{uu} n, \quad f = X_{uv} n, \quad g = X_{vv} n,$$

que son funciones de la carta X (salvo por el signo, que depende de la elección del sentido de n). Si la parametrización de la curva no es la natural y $s(t)$ es la longitud de arco, usamos la regla de la cadena:

$$\frac{du}{dt} = \frac{du}{ds} \frac{ds}{dt}, \quad \frac{dv}{dt} = \frac{dv}{ds} \frac{ds}{dt},$$

con la que la fórmula, llamando ahora u' , v' , s' a las derivadas respecto de t (hasta ahora eran las derivadas respecto de s), se convierte en

$$\kappa_n = \frac{e u'^2 + 2f u'v' + g v'^2}{s'^2}.$$

Observemos que esta expresión no depende de la curva (u, v) , sino sólo de su derivada (u', v') (recordemos que $s' = \|(u', v')\|$). De aquí deducimos:

Teorema 6.25 (Teorema de Meusnier) *Si S es una superficie, todas las curvas contenidas en S que pasan por un punto p con un mismo vector tangente tienen la misma curvatura normal. Ésta viene dada por*

$$\kappa_n = \frac{e du^2 + 2f dudv + g dv^2}{E du^2 + 2F dudv + G dv^2}.$$

DEMOSTRACIÓN: Es claro que X es una carta alrededor de p y α es una curva contenida en S que pasa por p con vector tangente w , entonces la representación de α en la carta X es $X^{-1} \circ \alpha$, luego el vector tangente de esta representación —el que en la discusión previa al teorema llamábamos (u', v') — es $(du(p)(w), dv(p)(w))$, donde ahora u y v son las funciones coordenadas de X .

Así pues, la fórmula que habíamos obtenido nos da que

$$\kappa_n(p)(w) = \frac{e(p) du(p)^2(w) + 2f(p) du(p)(w)dv(p)(w) + g(p) dv(p)^2(w)}{E(p) du(p)^2(w) + 2F(p) du(p)(w)dv(p)(w) + G dv(p)^2(w)},$$

entendiendo aquí a e, f, g como las composiciones con X^{-1} de las funciones del mismo nombre que teníamos definidas. ■

Definición 6.26 El elemento de longitud de una superficie S se conoce también con el nombre que le dio Gauss: *la primera forma fundamental* de S . Definimos *la segunda forma fundamental* de S como la aplicación que a $p \in S$ y cada vector $w \in T_p(S)$ le asigna la curvatura normal en p de las curvas contenidas en S que pasan por p con tangente w multiplicada por $\|w\|^2$.

El teorema anterior prueba que la segunda forma fundamental es en cada punto p una forma cuadrática definida sobre $T_p(S)$. Concretamente, si fijamos una carta tenemos

$$F^1 = E du^2 + 2F dudv + G dv^2, \quad F^2 = e du^2 + 2f dudv + g dv^2.$$

Ambas formas cuadráticas pueden considerarse definidas tanto sobre la superficie S como sobre el dominio de la carta (en cuyo caso du y dv representan simplemente las proyecciones de \mathbb{R}^2). Sin embargo, una diferencia importante es que, aunque las expresiones anteriores son válidas únicamente sobre el rango de una carta, la primera forma fundamental está definida sobre toda la superficie y está completamente determinada por la misma, mientras que la segunda sólo la tenemos definida en un entorno de cada punto y además salvo signo.

Para calcular explícitamente la segunda forma fundamental de una superficie notamos que

$$e = X_{uu}n = X_{uu} \frac{X_u \times X_v}{\|X_u \times X_v\|} = \frac{(X_{uu}, X_u, X_v)}{\sqrt{EG - F^2}},$$

e igualmente

$$f = \frac{(X_{uv}, X_u, X_v)}{\sqrt{EG - F^2}}, \quad g = \frac{(X_{vv}, X_u, X_v)}{\sqrt{EG - F^2}}.$$

Ejemplo Los coeficientes de la segunda forma fundamental de la superficie de revolución generada por la curva $(r(u), z(u))$ son

$$e = \frac{z''(u)r'(u) - z'(u)r''(u)}{\sqrt{r'(u)^2 + z'(u)^2}}, \quad f = 0, \quad g = \frac{z'(u)r(u)}{\sqrt{r'(u)^2 + z'(u)^2}}.$$

Para el caso del toro tenemos $(r(u), z(u)) = (R + r \cos v, r \sin v)$ luego queda

$$e = r, \quad f = 0, \quad g = R \cos u + r \cos^2 u.$$

■

Ejercicio: Comprobar que la curvatura normal en todo punto de la esfera de radio r y en toda dirección es igual a $\pm 1/r$, donde el signo es positivo si elegimos el vector normal que apunta hacia dentro de la esfera y negativo en caso contrario.

6.6 La curvatura de Gauss

Es un hecho conocido que si F es una forma bilineal simétrica en un espacio euclídeo existe una base ortonormal en la que la matriz de F es diagonal. Podemos aplicar esto a un plano tangente $T_p(S)$ de una superficie tomando el producto escalar determinado por la primera forma fundamental y como F la segunda forma fundamental. Entonces concluimos que existe una base (e_1, e_2) de $T_p(S)$ en la cual las expresiones en coordenadas de las formas fundamentales es

$$F^1(x, y) = x^2 + y^2 \quad \text{y} \quad F^2(x, y) = \lambda_1 x^2 + \lambda_2 y^2.$$

Los números λ_1 y λ_2 son los valores propios de cualquiera de las matrices de F^2 en cualquier base ortonormal de $T_p(S)$, luego están unívocamente determinados salvo por el hecho de que un cambio de carta puede cambiar sus signos. Podemos suponer $\lambda_1 \leq \lambda_2$. Entonces se llaman respectivamente *curvatura mínima* y *curvatura máxima* de S en p . En efecto, se trata del menor y el mayor valor que toma F^2 entre los vectores de norma 1, pues

$$\lambda_1 = \lambda_1(x^2 + y^2) \leq \lambda_1 x^2 + \lambda_2 y^2 = F^2(x, y) \leq \lambda_2(x^2 + y^2) = \lambda_2.$$

Si $w \in T_p(S)$ tiene norma arbitraria entonces aplicamos esto a $w/\|w\|$ y concluimos que

$$\lambda_1 \leq \frac{F^2(w)}{F^1(w)} \leq \lambda_2,$$

es decir, $\lambda_1 \leq \kappa_n \leq \lambda_2$. Así pues, λ_1 y λ_2 son la menor y la mayor curvatura normal que alcanzan las curvas que pasan por p . Además se alcanzan en direcciones perpendiculares e_1 y e_2 , llamadas *direcciones principales* en p . Notemos que puede ocurrir $\lambda_1 = \lambda_2$, en cuyo caso la curvatura normal es la misma en todas direcciones y no hay direcciones principales distinguidas. Los puntos de S donde $\lambda_1 = \lambda_2$ se llaman *puntos umbilicales*.

Veamos ahora cómo calcular las direcciones principales en una carta. Consideremos la fórmula de Meusnier como función (diferenciable) de dos variables. Si (du, dv) marca una dirección principal⁷ entonces κ_n es máximo o mínimo en este punto, luego el teorema 5.16 afirma que sus derivadas parciales han de anularse en él. Así pues, se ha de cumplir

$$\begin{aligned} \frac{\partial \kappa_n}{\partial du} &= \frac{2(e du + f dv)}{F^1(du, dv)} - \frac{2(E du + F dv)}{F^1(du, dv)^2} F^2(du, dv) = 0, \\ \frac{\partial \kappa_n}{\partial dv} &= \frac{2(f du + g dv)}{F^1(du, dv)} - \frac{2(F du + G dv)}{F^1(du, dv)^2} F^2(du, dv) = 0. \end{aligned}$$

Despejando obtenemos

$$\begin{aligned} \kappa_n &= \frac{F^2(du, dv)}{F^1(du, dv)} = \frac{e du + f dv}{E du + F dv}, \\ \kappa_n &= \frac{F^2(du, dv)}{F^1(du, dv)} = \frac{f du + g dv}{F du + G dv}. \end{aligned} \tag{6.12}$$

Al igualar ambas ecuaciones obtenemos una condición necesaria para que un vector indique una dirección principal. Es fácil ver que puede expresarse en la forma:

$$\begin{vmatrix} dv^2 & -dudv & du^2 \\ E & F & G \\ e & f & g \end{vmatrix} = 0.$$

Si (E, F, G) (en un punto) es múltiplo de (e, f, g) entonces la ecuación se cumple trivialmente, pero por otra parte es claro que la curvatura normal es constante y no hay direcciones principales. En caso contrario es claro tenemos una forma cuadrática con al menos dos coeficientes no nulos. Si suponemos, por ejemplo, que el coeficiente de dv^2 es no nulo, entonces $du \neq 0$, y al dividir entre du^2 la forma cuadrática se convierte en una ecuación de segundo grado

⁷Aquí podemos considerar $(du, dv) \in \mathbb{R}^2$. La notación diferencial está motivada por lo siguiente: Fijada una carta con coordenadas u, v , una curva regular en la superficie viene determinada por una representación coordenada $(u(t), v(t))$. El vector tangente a la curva en un punto dado marcará una dirección principal si y sólo si la fórmula de Meusnier evaluada en $(u'(t), v'(t))$ toma un valor máximo o mínimo, pero dicha fórmula depende sólo de las diferenciales $(du(t), dv(t))$, por lo que en realidad buscamos una relación entre du y dv .

en la razón dv/du . Esta ecuación tiene a lo sumo dos soluciones linealmente independientes, luego éstas han de ser necesariamente las direcciones principales. Por consiguiente la ecuación caracteriza dichas direcciones.

Definición 6.27 Se llama *curvatura media* y *curvatura total o de Gauss* de una superficie S en un punto p a los números

$$H = \frac{\lambda_1 + \lambda_2}{2}, \quad K = \lambda_1 \lambda_2.$$

Notemos que el signo de H depende de la carta, mientras que el de K es invariante. Si operamos en (6.12) obtenemos

$$\begin{aligned} (e - E\kappa_n)du + (f - F\kappa_n)dv &= 0, \\ (f - F\kappa_n)du + (g - G\kappa_n)dv &= 0. \end{aligned}$$

Puesto que el sistema tiene una solución no trivial en (du, dv) se ha de cumplir

$$\begin{vmatrix} e - E\kappa_n & f - F\kappa_n \\ f - F\kappa_n & g - G\kappa_n \end{vmatrix} = 0,$$

o equivalentemente

$$(EG - F^2)\kappa_n^2 - (eG - 2Ff + gE)\kappa_n + (eg - f^2) = 0.$$

Esta ecuación la cumplen las curvaturas principales $\kappa_n = \lambda_1, \lambda_2$ y por otro lado tiene sólo dos soluciones, luego

$$\begin{aligned} H &= \frac{eG - 2Ff + gE}{2(EG - F^2)}, \\ K &= \frac{eg - f^2}{EG - F^2}. \end{aligned} \tag{6.13}$$

En particular vemos que la curvatura de Gauss es el cociente de los determinantes de las dos formas fundamentales.

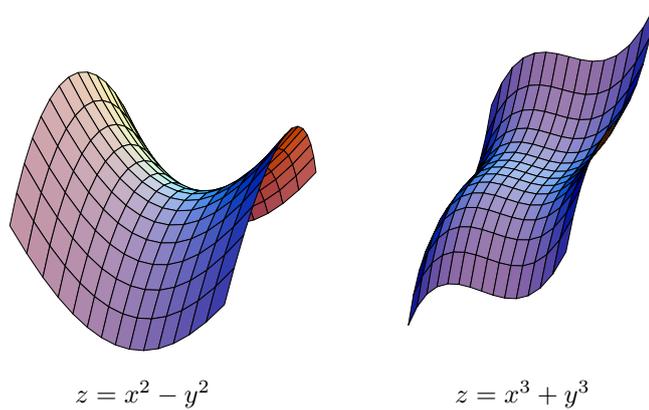
Ejercicio: Calcular la curvatura media y la curvatura de Gauss del cilindro, el cono, el toro y la esfera.

Definición 6.28 Un punto p de una superficie S es *elíptico* o *hiperbólico* según si $K(p) > 0$ o $K(p) < 0$. Si $K(p) = 0$ distinguiremos entre puntos *parabólicos*, cuando sólo una de las curvaturas extremas es nula y puntos *planos*, cuando las dos curvaturas extremas son nulas.

Si un punto es elíptico todas las curvas que pasan por él tienen la curvatura normal del mismo signo, por lo que la superficie se curva toda hacia el mismo lado del plano tangente, como es el caso de la esfera o del toro. Si un punto es hiperbólico entonces hay curvas (perpendiculares, de hecho) que pasan por él con curvaturas en sentidos opuestos, luego la superficie tiene puntos próximos a ambos lados del plano tangente. Es el caso del hiperboloide $z = x^2 - y^2$, cuya curvatura en la carta $(u, v, u^2 - v^2)$ viene dada por $K = -4/\sqrt{4u^2 + 4v^2 + 1}$ ³.

Los puntos de un cilindro son parabólicos. Las curvas $u = \text{cte.}$ y $v = \text{cte.}$ son circunferencias de radio r y rectas, respectivamente. Las primeras tienen curvatura normal $\lambda_2 = 1/r$ y las segundas $\lambda_1 = 0$. Es fácil ver que se trata de las curvaturas principales. Todos los cálculos son sencillos.

Todos los puntos de un plano son puntos planos. Otro ejemplo es el punto $(0, 0)$ en la gráfica de $x^3 + y^3$.



Probamos ahora una caracterización algebraica de la curvatura de Gauss que más adelante nos dará una interpretación geométrica de la misma. Observemos que si n es una determinación del vector normal alrededor de un punto p en una superficie S y llamamos S^2 a la esfera de centro $(0, 0, 0)$ y radio 1, entonces $dn(p) : T_p(S) \rightarrow T_{n(p)}(S^2)$, pero como $n(p)$ es perpendicular a $T_p(S)$, en realidad $T_{n(p)}(S^2) = T_p(S)$, luego podemos considerar a $dn(p)$ como un endomorfismo de $T_p(S)$.

Teorema 6.29 Sea S una superficie y n una determinación del vector normal alrededor de un punto p . Entonces $K(p) = |dn(p)|$.

DEMOSTRACIÓN: Sea X una carta alrededor de p . Entonces una base de $T_p(S)$ la forman los vectores X_u y X_v . Llamemos $n(u, v)$ a $X \circ n$. Entonces

$$\begin{aligned} dn(p)(X_u) &= dn(p)(dX(u, v)(1, 0)) = dn(u, v)(1, 0) = n_u, \\ dn(p)(X_v) &= dn(p)(dX(u, v)(0, 1)) = dn(u, v)(0, 1) = n_v. \end{aligned}$$

Si expresamos

$$\begin{aligned} n_u &= aX_u + bX_v \\ n_v &= cX_u + dX_v \end{aligned}$$

entonces el determinante de $dn(p)$ es el de la matriz formada por a, b, c, d . Notemos que derivando las igualdades $nX_u = nX_v = 0$ se deduce la relación $n_uX_u = -nX_{uu} = -e$ y similarmente $n_uX_v = n_vX_u = -f$, $n_vX_v = -g$. Por consiguiente al multiplicar las ecuaciones anteriores por X_u y X_v obtenemos

$$-e = aE + bF, \quad -f = aF + bG, \quad -f = cE + dF, \quad -g = cF + dG,$$

de donde

$$-\begin{pmatrix} e & f \\ f & g \end{pmatrix} = \begin{pmatrix} a & b \\ c & d \end{pmatrix} \begin{pmatrix} E & F \\ F & G \end{pmatrix}.$$

Tomando determinantes concluimos que

$$eg - f^2 = |dn(p)| (EG - F^2),$$

luego efectivamente $|dn(p)| = K(p)$. ■

Consideremos ahora las fórmulas (6.7) que definen los símbolos de Christoffel. Al particularizarlas al caso de una superficie se convierten en

$$\begin{aligned} X_{uu} &= \Gamma_{11}^1 X_u + \Gamma_{11}^2 X_v + en, \\ X_{uv} &= \Gamma_{12}^1 X_u + \Gamma_{12}^2 X_v + fn, \\ X_{vv} &= \Gamma_{22}^1 X_u + \Gamma_{22}^2 X_v + gn, \end{aligned}$$

(en principio la componente normal ha de ser de la forma αn para cierto α , y multiplicando la igualdad por n se sigue que $\alpha = e, f, g$ según el caso.)

De estas ecuaciones se sigue

$$\begin{aligned} X_{uu}X_{vv} - X_{uv}^2 &= eg - f^2 + (\Gamma_{11}^1\Gamma_{22}^1 - (\Gamma_{12}^1)^2)E \\ &+ (\Gamma_{11}^1\Gamma_{22}^2 + \Gamma_{11}^2\Gamma_{22}^1 - 2\Gamma_{12}^1\Gamma_{12}^2)F \\ &+ (\Gamma_{11}^2\Gamma_{22}^2 - (\Gamma_{12}^2)^2)G. \end{aligned}$$

Por otra parte, derivando respecto a v y u respectivamente las relaciones

$$X_{uu}X_v = F_u - \frac{1}{2}E_v, \quad X_{uv}X_v = \frac{1}{2}G_u$$

y restando los resultados obtenemos

$$X_{uu}X_{vv} - X_{uv}^2 = -\frac{1}{2}E_{vv} + F_{uv} - \frac{1}{2}G_{uu}.$$

En definitiva resulta la expresión

$$\begin{aligned} eg - f^2 &= -\frac{1}{2}E_{vv} + F_{uv} - \frac{1}{2}G_{uu} - (\Gamma_{11}^1\Gamma_{22}^1 - (\Gamma_{12}^1)^2)E \\ &- (\Gamma_{11}^1\Gamma_{22}^2 + \Gamma_{11}^2\Gamma_{22}^1 - 2\Gamma_{12}^1\Gamma_{12}^2)F \\ &- (\Gamma_{11}^2\Gamma_{22}^2 - (\Gamma_{12}^2)^2)G. \end{aligned}$$

La fórmula (6.13) muestra ahora que la curvatura de Gauss de un punto depende únicamente de los coeficientes E, F, G de la primera forma fundamental y sus derivadas. Puesto que dos superficies localmente isométricas tienen cartas con los mismos coeficientes E, F, G , hemos probado el resultado que Gauss, en sus *Dquisitiones generales circa superficies curvas*, presentó con el nombre de *theorema egregium*:

Teorema 6.30 (Gauss) *Las isometrías locales conservan la curvatura.*

Las ecuaciones (6.11) nos dan la siguiente expresión para la curvatura respecto a una carta con $F = 0$:

$$K = \frac{E_u G_u + E_v^2}{4E^2 G} + \frac{E_v G_v + G_u^2}{4EG^2} - \frac{E_{vv} + G_{uu}}{2EG}, \quad \text{si } F = 0.$$

Desde aquí es fácil deducir a su vez los siguientes casos particulares:

$$K = -\frac{1}{2A} \left(\frac{\partial^2 \log A}{\partial u^2} + \frac{\partial^2 \log A}{\partial v^2} \right), \quad \text{si } F = 0, \quad E = G = A,$$

$$K = -\frac{1}{\sqrt{G}} \frac{\partial^2 \sqrt{G}}{\partial u^2}, \quad \text{si } F = 0, \quad E = 1.$$

En particular, la curvatura de la superficie de revolución definida por la curva $(r(u), z(u))$ es

$$K = -\frac{r''(u)}{r(u)}.$$

Ejemplo Se llama *pseudoesfera* a la superficie de revolución P generada por la tractriz. Recordemos que la tractriz es

$$(r(u), z(u)) = \left(l \operatorname{sen} u, l \cos u + l \log \tan \frac{u}{2} \right), \quad \text{para } \pi/2 < u < \pi.$$

Por lo tanto la pseudoesfera está dada por

$$X(u, v) = \left(l \operatorname{sen} u \cos v, l \operatorname{sen} u \operatorname{sen} v, l \cos u + l \log \tan \frac{u}{2} \right).$$

Recordemos también que la longitud de arco de la tractriz es $w = -l \log \operatorname{sen} u$, luego $\operatorname{sen} u = e^{-w/l}$. Las cartas

$$Y :]0, +\infty[\times]v_0 - \pi, v_0 + \pi[\longrightarrow \mathbb{R}^3$$

que resultan de tomar la tractriz parametrizada por el arco tienen la primera forma fundamental determinada por

$$E = 1, \quad F = 0, \quad G = r(w)^2 = l^2 e^{-2w/l},$$

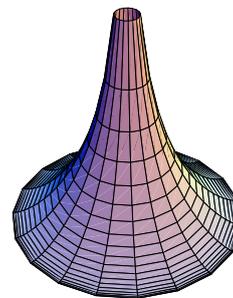
es decir,

$$ds^2 = dw^2 + e^{-w/l} dv^2.$$

De aquí se sigue fácilmente que $K = -1/l^2$.

Por lo tanto un ejemplo de superficie de curvatura constante igual a $K < 0$ es la pseudoesfera

$$\frac{1}{\sqrt{-K}} \left(\operatorname{sen} u \cos v, \operatorname{sen} u \operatorname{sen} v, \cos u + \log \tan \frac{u}{2} \right).$$



Hay una parametrización alternativa respecto a la cual la primera forma fundamental tiene una expresión especialmente simple, que consiste en tomar $x = v$, $y = e^{w/l}$. Así obtenemos una carta $Z :]x_0 - \pi, x_0 + \pi[\times]1, +\infty[\rightarrow \mathbb{R}^3$ respecto a la cual la primera forma fundamental es

$$ds^2 = \frac{l^2}{y^2}(dx^2 + dy^2).$$

Observemos que las líneas coordenadas $y = \text{cte.}$ siguen siendo circunferencias y las de la forma $x = \text{cte.}$ siguen siendo tractrices, pero con una parametrización distinta de la dada por el arco. ■

Capítulo VII

Ecuaciones diferenciales ordinarias

Una *ecuación diferencial ordinaria* es una relación de la forma

$$f(t, y(t), y'(t), \dots, y^{(n)}(t)) = 0,$$

donde $f : D \subset \mathbb{R}^{n+2} \rightarrow \mathbb{R}$ e $y : I \rightarrow \mathbb{R}$ es una función definida en un intervalo I derivable n veces. Normalmente la función y es desconocida, y entonces se plantea el problema de *integrar* la ecuación, es decir, encontrar todas las funciones y que la satisfacen. El adjetivo “ordinaria” se usa para indicar que la función incógnita tiene una sola variable. Las ecuaciones que relacionan las derivadas parciales de una función de varias variables se llaman *ecuaciones diferenciales en derivadas parciales*, pero no vamos a ocuparnos de ellas. Dentro de las ecuaciones ordinarias, nos vamos a ocupar únicamente de un caso más simple pero suficientemente general: aquel en que tenemos despejada la derivada de orden mayor, es decir, una ecuación de la forma

$$y^{(n)}(t) = f(t, y(t), y'(t), \dots, y^{(n-1)}(t)).$$

El número n se llama *orden* de la ecuación. El caso más simple es la ecuación de primer orden $y' = f(t)$. Sabemos que si la ecuación tiene solución de hecho hay infinitas de ellas pero, en un intervalo dado, cada una se diferencia de las demás en una constante, de modo que una solución queda completamente determinada cuando se especifica un valor $y(t_0) = y_0$. Veremos que esto sigue siendo válido para todas las ecuaciones de primer orden. Por ello se define un *problema de Cauchy* como

$$\left. \begin{array}{l} y' = f(t, y) \\ y(t_0) = y_0 \end{array} \right\}$$

Resolver el problema significa encontrar una función y definida alrededor de t_0 de modo que satisfaga la ecuación diferencial y cumpla la condición inicial $y(t_0) = y_0$. Probaremos que bajo condiciones muy generales los problemas de Cauchy tienen solución única.

Toda la teoría se aplica igualmente al caso de sistemas de ecuaciones diferenciales. De hecho un sistema de ecuaciones puede verse como una única ecuación vectorial. Basta considerar que $y : I \rightarrow \mathbb{R}^n$ y $f : D \subset \mathbb{R}^{n+1} \rightarrow \mathbb{R}^n$.

Son muchas las ocasiones en las que el único conocimiento que tenemos de una o varias funciones es el hecho de que satisfacen un sistema de ecuaciones diferenciales. Por ejemplo, las ecuaciones (6.10) del capítulo anterior son un sistema de ecuaciones de segundo orden que determinan cuándo una curva $x(s)$ representa a una geodésica de una variedad en una carta dada. Los resultados que probaremos en este capítulo nos asegurarán en particular la existencia de geodésicas.

Antes de abordar la teoría general vamos a ver un ejemplo concreto. Una forma de asegurar que una ecuación diferencial tiene solución es encontrarla explícitamente:

Ejemplo *Un depósito de 1m^2 de base contiene agua hasta una altura de 1m . En la base abrimos un agujero de 1cm^2 de superficie. ¿Cuánto tardará en descender 0.5m el nivel del agua?*

SOLUCIÓN: Nos basamos en la *ley de Torricelli*, que afirma que la velocidad de salida del agua en cada instante es $\sqrt{2gh}$, donde $g = 9.8\text{m/s}^2$ es la intensidad del campo gravitatorio en la superficie de la Tierra y $h = h(t)$ es el nivel del agua.¹

Es fácil convencerse entonces de que el volumen de agua que sale por el agujero por unidad de tiempo es

$$\frac{dV}{dt} = A\sqrt{2gh(t)},$$

donde A es el área del agujero. Por otro lado, el volumen de agua en el depósito es $V = Bh$, donde B es el área de la base (y suponiendo que el depósito tiene paredes horizontales),² luego la variación de dicho volumen por unidad de tiempo será

$$\frac{dV}{dt} = B\frac{dh}{dt}.$$

Al igualar ambas expresiones tenemos una ecuación diferencial de primer orden:

$$B\frac{dh}{dt} = -A\sqrt{2gh},$$

a la que hay que añadir la condición inicial $h(0) = h_0$, que indica que en $t = 0$ la altura del agua era la dada por el enunciado. Notemos que hay que poner un signo negativo, pues la derivada de h es negativa, luego el miembro derecho también tiene que serlo.

¹La ley de Torricelli sólo es válida bajo el supuesto de que el depósito es lo suficientemente grande como para que el descenso en la altura sea inapreciable en intervalos cortos de tiempo.

²Esto es un hecho elemental bien conocido que será evidente a partir de los resultados del cálculo integral que presentaremos en los capítulos siguientes.

La técnica más elemental para resolver una ecuación diferencial es la *separación de variables*, que consiste en dejar la variable dependiente h en un miembro y la variable independiente t en el otro (cosa que no siempre se puede hacer):

$$\frac{dh}{\sqrt{h}} = -\frac{A\sqrt{2g}}{B} dt.$$

Ahora integramos ambos miembros:³

$$2\sqrt{h} = -\frac{A\sqrt{2g}}{B}t + k.$$

La constante k podemos calcularla mediante la condición inicial, que se traduce en que $k = 2\sqrt{h_0}$. Ahora despejamos h :

$$h(t) = \left(\sqrt{h_0} - \sqrt{\frac{g}{2}} \frac{A}{B} t \right)^2.$$

Un simple cálculo con los datos del problema muestra que el nivel del agua tardará 22 minutos en reducirse a la mitad. En teoría, el depósito tardará 75 minutos en vaciarse, aunque el cálculo no es fiable porque la ley de Torricelli no es aplicable para niveles bajos. ■

Volviendo a la situación general, observemos que de momento no tenemos garantizada la existencia de solución ni en el caso más simple: $y' = f(x)$. Que esta ecuación diferencial tenga solución equivale a que la función $f(x)$ tenga primitiva. En el ejemplo precedente hemos justificado la existencia de primitivas de los dos miembros de la ecuación diferencial encontrándolas explícitamente. En este capítulo aceptaremos sin prueba que toda función continua en un intervalo abierto tiene una primitiva en dicho intervalo, hecho que demostraremos en el capítulo siguiente (teorema 8.58), y a partir de ahí obtendremos resultados generales sobre existencia de solución de sistemas de ecuaciones diferenciales.

De este modo, si $f : I \rightarrow \mathbb{R}$ es una función continua en un intervalo abierto I y $a \in I$, una primitiva de f en I es la función

$$F(x) = \int_a^x f(x) dx$$

(lo cual no es más que una notación para la primitiva, no una prueba de su existencia). Para una función continua $f : [a, b] \rightarrow \mathbb{R}^n$, definimos su integral como

$$\int_a^b f(x) dx = \left(\int_a^b f_1(x) dx, \dots, \int_a^b f_n(x) dx \right).$$

³Una justificación de estas operaciones que no requiera separar las diferenciales es como sigue: tenemos

$$\frac{1}{\sqrt{h}} h'(t) = -\frac{A\sqrt{2g}}{B},$$

y una primitiva del primer miembro es $2\sqrt{h(t)}$, mientras que una del segundo miembro es $-A\sqrt{2g}t/B$, y dos primitivas de una misma función deben diferenciarse a lo sumo en una constante k .

7.1 Ecuaciones diferenciales de primer orden

Nos ocupamos ahora de asegurar la existencia y unicidad de la solución de los problemas de Cauchy. Nos basaremos en un resultado general sobre espacios métricos completos:

Teorema 7.1 (Teorema de punto fijo de Banach) *Sea M un espacio métrico completo y $T : M \rightarrow M$ una aplicación tal que existe un número real $0 < \alpha < 1$ de modo que*

$$d(T(x), T(y)) < \alpha d(x, y), \quad \text{para todo } x, y \in M.$$

Entonces existe un único $x \in M$ tal que $T(x) = x$.

Las aplicaciones T que cumplen la propiedad indicada se llaman *contractivas*. Los puntos x que cumplen $T(x) = x$ se llaman *puntos fijos* de T . El teorema afirma, pues, que toda aplicación contractiva en un espacio métrico completo tiene un único punto fijo.

DEMOSTRACIÓN: Tomamos un punto arbitrario $x_0 \in M$ y consideramos la sucesión dada por $x_{n+1} = T(x_n)$. Por la propiedad de T , tenemos que

$$\begin{aligned} d(x_1, x_2) &= d(T(x_0), T(x_1)) < \alpha d(x_0, x_1), \\ d(x_2, x_3) &= d(T(x_1), T(x_2)) < \alpha d(x_1, x_2) = \alpha^2 d(x_0, x_1), \end{aligned}$$

y en general concluimos $d(x_n, x_{n+1}) < \alpha^n d(x_0, x_1)$. Aplicando la desigualdad triangular resulta, para $n < m$,

$$d(x_n, x_m) < \left(\sum_{i=n}^{m-1} \alpha^i \right) d(x_0, x_1) < \left(\sum_{i=n}^{\infty} \alpha^i \right) d(x_0, x_1) = \frac{\alpha^n}{1 - \alpha} d(x_0, x_1).$$

El término de la derecha tiende a 0, lo que significa que la sucesión x_n es de Cauchy. Como el espacio M es completo existe $x = \lim_n x_n \in M$. Veamos que x es un punto fijo de T . Para ello observamos que

$$\begin{aligned} d(x, T(x)) &\leq d(x, x_n) + d(x_n, x_{n+1}) + d(x_{n+1}, T(x)) \\ &< (1 + \alpha)d(x, x_n) + \alpha^n d(x_0, x_1). \end{aligned}$$

El último término tiende a 0, luego ha de ser $d(x, T(x)) = 0$, es decir, $T(x) = x$. Si y es otro punto fijo de T , entonces $d(T(x), T(y)) = d(x, y)$, en contradicción con la propiedad contractiva, luego el punto fijo es único. ■

Es frecuente que una ecuación diferencial dependa de uno más parámetros. Por ejemplo, la fuerza $F(t, x)$ que afecta a un móvil de masa m es, por lo general, función del tiempo t y de la posición x . La segunda ley de Newton afirma que su trayectoria $x(t)$ obedece la ecuación diferencial de segundo orden

$$F(t, x) = m x''(t),$$

donde la masa m es un parámetro. En casos como éste podemos considerar la solución como función de los parámetros, es decir, $x(t, m)$ es la posición en el instante t de un cuerpo de masa m sometido a la fuerza $F(t, x)$ (y en unas condiciones iniciales dadas). En el teorema de existencia y unicidad que damos a continuación contemplamos la existencia de estos parámetros y probamos que la solución depende continuamente de ellos.

Teorema 7.2 Sean $t_0, a, b_1, \dots, b_n, y_1^0, \dots, y_n^0$ números reales. Consideremos una aplicación continua

$$f : [t_0 - a, t_0 + a] \times \prod_{i=1}^n [y_i^0 - b_i, y_i^0 + b_i] \times K \longrightarrow \mathbb{R}^n,$$

donde K es un espacio métrico compacto. Sea M una cota de f respecto a la norma $\| \cdot \|_\infty$ en \mathbb{R}^n . Supongamos que existe una constante N tal que

$$\|f(t, y, \mu) - f(t, z, \mu)\|_\infty \leq N \|y - z\|_\infty.$$

Entonces el problema de Cauchy

$$\left. \begin{aligned} y'(t, \mu) &= f(t, y, \mu) \\ y(t_0, \mu) &= y_0 \end{aligned} \right\}$$

tiene solución única $y : [t_0 - h_0, t_0 + h_0] \times K \longrightarrow \mathbb{R}^n$, continua en su dominio, donde h_0 es cualquier número real tal que

$$0 < h_0 \leq \min \left\{ a, \frac{b_1}{M}, \dots, \frac{b_n}{M} \right\}, \quad h_0 < \frac{1}{N}.$$

Se entiende que derivada de y que aparece en el problema de Cauchy es respecto de la variable t . Teóricamente deberíamos usar la notación de derivadas parciales, pero es costumbre usar la notación del análisis de una variable para evitar que el problema parezca una ecuación diferencial en derivadas parciales, cuando en realidad no lo es.

DEMOSTRACIÓN: Sea h_0 en las condiciones indicadas, sea $I = [t - h_0, t + h_0]$, sea $D = \prod_{i=1}^n [y_i^0 - b_i, y_i^0 + b_i]$ y sea $M = C(I \times K, D)$, que es un espacio de Banach con la norma supremo. Definimos el operador $T : M \longrightarrow M$ mediante

$$T(y)(t, \mu) = y^0 + \int_{t_0}^t f(t, y(t, \mu), \mu) dt.$$

Hemos de probar que $T(y)(t, \mu) \in D$ y que $T(y)$ es una aplicación continua. En primer lugar,

$$\begin{aligned} |T(y)_i(t, \mu) - y_i^0| &= \left| \int_{t_0}^t f_i(t, y, \mu) dt \right| \leq \int_{t_0}^t |f_i(t, y, \mu)| dt \leq M \left| \int_{t_0}^t dt \right| \\ &= M|t - t_0| \leq Mh_0 \leq b_i. \end{aligned}$$

Esto prueba que $T(y)(t, \mu) \in D$. La continuidad es consecuencia de un cálculo rutinario:

$$\begin{aligned} \|T(y)(t_1, \mu_1) - T(y)(t_2, \mu_2)\|_\infty &= \max_i \left| \int_{t_0}^{t_1} f_i(t, y, \mu_1) dt - \int_{t_0}^{t_2} f_i(t, y, \mu_2) dt \right| \\ &\leq \max_i \left(\left| \int_{t_0}^{t_1} f_i(t, y, \mu_1) dt - \int_{t_0}^{t_1} f_i(t, y, \mu_2) dt \right| + \right. \\ &\quad \left. \left| \int_{t_0}^{t_1} f_i(t, y, \mu_2) dt - \int_{t_0}^{t_2} f_i(t, y, \mu_2) dt \right| \right) \\ &\leq \max_i \left(\left| \int_{t_0}^{t_1} (f_i(t, y, \mu_1) - f_i(t, y, \mu_2)) dt \right| + \left| \int_{t_2}^{t_1} f_i(t, y, \mu_2) dt \right| \right) \\ &\leq \max_i \left| \int_{t_0}^{t_1} |f_i(t, y, \mu_1) - f_i(t, y, \mu_2)| dt \right| + M|t_1 - t_2|. \end{aligned}$$

Sea $\epsilon > 0$. La función $f_i(t, y(t, \mu), \mu)$ es uniformemente continua en el compacto $I \times K$, luego existe un $\delta > 0$ tal que si $d(\mu_1, \mu_2) < \delta$, entonces $|f_i(t, y, \mu_1) - f_i(t, y, \mu_2)| < \epsilon/2h_0$. Podemos suponer que esto vale para todo $i = 1, \dots, n$, y si suponemos también que $|t_1 - t_2| < \epsilon/2M$ concluimos que

$$\|T(y)(t_1, \mu_1) - T(y)(t_2, \mu_2)\|_\infty < \epsilon.$$

Esto prueba la continuidad de $T(y)$ en el punto (t_1, μ_1) .

Ahora probamos que T es contractivo, con constante $\alpha = Nh_0 < 1$. En efecto,

$$\begin{aligned} \|T(y)(t, \mu) - T(z)(t, \mu)\|_\infty &= \max_i \left| \int_{t_0}^t (f_i(t, y(t, \mu), \mu) - f_i(t, z(t, \mu), \mu)) dt \right| \\ &\leq \max_i \left| \int_{t_0}^t N \|y(t, \mu) - z(t, \mu)\|_\infty dt \right| \leq N \max_i \left| \int_{t_0}^t \|y - z\| dt \right| \leq Nh_0 \|y - z\|. \end{aligned}$$

Por definición de norma supremo resulta $\|T(y) - T(z)\| \leq \alpha \|y - z\|$. El teorema anterior implica ahora la existencia de una única función $y \in M$ tal que

$$y(t, \mu) = y^0 + \int_{t_0}^t f(t, y, \mu) dt,$$

pero es claro que esto equivale a ser solución del problema de Cauchy, luego éste tiene solución única. ■

En la práctica, hay una hipótesis más fuerte que la condición de Lipschitz que hemos exigido en el teorema anterior pero que es más fácil de comprobar. Se trata de exigir simplemente que la función f sea de clase C^1 .

Teorema 7.3 *Sea $f : D \subset \mathbb{R}^n \rightarrow \mathbb{R}^m$ una función de clase C^1 en un abierto D . Para todo subconjunto compacto convexo $C \subset D$ existe una constante N tal que si $y, z \in C$ entonces $\|f(y) - f(z)\|_\infty \leq N \|y - z\|_\infty$.*

DEMOSTRACIÓN: Si llamamos f_1, \dots, f_m a las funciones coordenadas de f , basta probar que $|f_i(y) - f_i(z)| \leq N_i \|y - z\|_\infty$ para todo $y, z \in C$, pues tomando como N la mayor de las constantes N_i se cumple la desigualdad buscada. Equivalentemente, podemos suponer que $m = 1$.

Dados $y, z \in C$, consideramos la función $g(h) = f(y + h(z - y))$, definida en $[0, 1]$, pues C es convexo. Se cumple $g(0) = f(y)$, $g(1) = f(z)$. Por el teorema del valor medio existe $0 < h_0 < 1$ tal que

$$f(z) - f(y) = g'(h_0) = df(\xi)(z - y) = \nabla f(\xi)(z - y),$$

donde $\xi = y + h_0(z - y) \in C$. Sea N_0 una cota del módulo de las derivadas parciales de f (que por hipótesis son continuas) en el compacto C . Entonces

$$|f(z) - f(y)| = \left| \sum_{i=1}^n D_i f(\xi)(z_i - y_i) \right| \leq \sum_{i=1}^n N_0 \|z - y\|_\infty = nN_0 \|z - y\|_\infty.$$

■

En realidad, si tomamos como hipótesis que la ecuación diferencial sea de clase C^1 obtenemos no sólo la continuidad de la solución respecto de los parámetros, sino también la derivabilidad.

Teorema 7.4 Sea $f : D \subset \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ una función de clase C^1 en el abierto D . Sea $(t_0, y_0, \mu) \in D$. Entonces el problema de Cauchy

$$\left. \begin{aligned} y'(t, \mu) &= f(t, y, \mu) \\ y(t_0, \mu) &= y_0 \end{aligned} \right\}$$

tiene solución única definida y de clase C^1 en un entorno de (t_0, μ) .

DEMOSTRACIÓN: Tomamos un entorno C de (t_0, y_0, μ) que esté contenido en D y sea producto de intervalos cerrados de centro cada una de las componentes del punto. En particular es convexo y compacto. El teorema anterior garantiza que se cumplen las hipótesis del teorema de existencia y unicidad. Falta probar que la función $y(t, \mu)$ es de clase C^1 en su dominio.

Obviamente y es derivable respecto de t y la derivada es continua. Veamos que lo mismo sucede con las demás variables. Sea e_i un vector de la base canónica de \mathbb{R}^m . Consideramos un punto (t_1, μ_1) del dominio de y . La función

$$Q(t, \mu, h) = \frac{y(t, \mu + he_i) - y(t, \mu)}{h}$$

está definida en los puntos de un entorno de $(t_1, \mu_1, 0)$ para los que $h \neq 0$. Hemos de probar que tiene límite cuando h tiende a 0. Claramente

$$\frac{\partial Q}{\partial t} = \frac{f(t, y(t, \mu + he_i), \mu + he_i) - f(t, y(t, \mu), \mu)}{h}.$$

Llamemos

$$E(p, x) = \begin{cases} \frac{f(p+x) - f(p) - df(p)(x)}{\|x\|} & \text{si } x \neq 0, \\ 0 & \text{si } x = 0. \end{cases}$$

Se comprueba que es continua en un entorno de $(t_1, y(t_1, \mu_1), \mu_1, 0)$. Sólo hay que ver la continuidad en los puntos de la forma $(q, 0)$. Se demuestra para cada función coordenada independientemente, y a su vez para ello se aplica el teorema del valor medio la la función $f_j(p + tx)$. El resultado es que

$$E_j(p, x) = (\nabla f_j(p') - \nabla f_j(p)) \frac{x}{\|x\|},$$

donde p' es un punto entre p y $p + x$, y ahora basta aplicar la continuidad de las derivadas parciales de f .

En términos de E tenemos

$$\begin{aligned} \frac{\partial Q}{\partial t} &= df(t, y(t, \mu), \mu)(0, Q(t, \mu, h), e_i) \\ &+ \|(0, Q(t, \mu, h), 1)\| \frac{|h|}{h} E(0, y(t, \mu + he_i) - y(t, \mu), h). \end{aligned}$$

Más brevemente

$$\frac{\partial Q}{\partial t} = df(t, y(t, \mu), \mu)(0, Q(t, \mu, h), e_i) + \|(0, Q(t, \mu, h), 1)\| E^*(t, \mu, h),$$

entendiendo que $E^*(t, \mu, h)$ es continua en un entorno de $(t_1, \mu_1, 0)$ y se anula en los puntos donde $h = 0$.

Esto significa que Q es la solución de una ecuación diferencial determinada por la función continua

$$g(t, Q, \mu, h) = df(t, y(t, \mu), \mu)(0, Q, e_i) + \|(0, Q, 1)\| E^*(t, \mu, h),$$

donde μ y h son parámetros. Esta función no es diferenciable, pero cumple claramente la hipótesis del teorema 7.2. Concretamente la consideramos definida en un producto de intervalos de centros t_1 , $Q(t_1, \mu_1, h)$ (para un h fijo) por un entorno compacto K de (μ_1, h) que contenga a $(\mu_1, 0)$. Si tomamos como condición inicial en el punto (t_1, μ_1, h) la determinada por la función Q que ya tenemos definida, el teorema 7.2 nos garantiza la existencia de solución continua en un conjunto de la forma $[t_1 - r, t_1 + r] \times K$. Por la unicidad la solución debe coincidir con la función Q que ya teníamos. En particular coincidirá con ella en los puntos de un entorno de $(t_1, \mu_1, 0)$ tales que $h \neq 0$. De aquí se sigue que existe

$$\lim_{h \rightarrow 0} Q(t_1, \mu_1, h) = \frac{\partial y}{\partial \mu_i}(t_1, \mu_1).$$

Más aún, esta derivada satisface la ecuación diferencial

$$\frac{\partial}{\partial t} \frac{\partial y}{\partial \mu_i} = df(t, y(t, \mu), \mu) \left(0, \frac{\partial y}{\partial \mu_i}, e_i \right).$$

Esto implica que es continua, luego y es de clase C^1 . ■

Hemos probado que las derivadas respecto a los parámetros de una ecuación diferencial determinada por una función de clase C^k satisfacen una ecuación

diferencial de clase C^{k-1} . Una simple inducción prueba entonces que la solución y de una ecuación de clase C^k es una función de clase C^k .

Notemos que la solución y de un problema de Cauchy puede considerarse también como función de las condiciones iniciales, es decir, $y(t, \mu, t_0, y_0)$. Del teorema anterior se deduce que y es continua respecto a todas las variables, es decir, como función definida en un entorno de (t_0, μ, t_0, y_0) en $\mathbb{R} \times \mathbb{R}^m \times \mathbb{R} \times \mathbb{R}^n$. Para ello basta ver que si hacemos $z = y(t, \mu, t_0, y_0) - y_0$ y $r = t - t_0$, el problema

$$\left. \begin{aligned} y'(t) &= f(t, y, \mu) \\ y(t_0) &= y_0 \end{aligned} \right\}$$

es equivalente a

$$\left. \begin{aligned} z'(r) &= f(r + t_0, z + y_0, \mu) \\ z(0) &= 0 \end{aligned} \right\}$$

en el sentido de que una solución de uno da una del otro mediante los cambios de variable indicados. El segundo miembro del segundo problema es una función $g(r, z, \nu)$, donde $\nu = (t_0, y_0, \mu) \in D$. Concretamente, el dominio de g es el abierto

$$\{(r, z, t_0, y_0, \mu) \mid (r + t_0, z + y_0, \mu) \in D, (t_0, y_0, \mu) \in D\}.$$

Una solución z del segundo problema definida en un entorno de $(0, t_0, y_0, \mu)$ se traduce en una solución $y(t, t_0, y_0, \mu)$ del primer problema definida en un entorno de (t_0, t_0, y_0, μ) . En definitiva tenemos el siguiente enunciado, más completo, del teorema de existencia y unicidad:

Teorema 7.5 Sea $f : D \subset \mathbb{R} \times \mathbb{R}^n \times \mathbb{R}^m \rightarrow \mathbb{R}^n$ una función de clase C^k ($k \geq 1$) en el abierto D . Sea $(t_0, y_0, \mu) \in D$. Entonces el problema de Cauchy

$$\left. \begin{aligned} y'(t, t_0, y_0, \mu) &= f(t, y, \mu) \\ y(t_0, t_0, y_0, \mu) &= y_0 \end{aligned} \right\}$$

tiene solución única de clase C^k en un entorno de (t_0, t_0, y_0, μ) .

Con más precisión, fijados (t_0, y_0, μ) , observamos que si tenemos dos intervalos abiertos I_1, I_2 que contengan a t_0 de modo que existan soluciones $y_i : I_i \rightarrow \mathbb{R}^n$ que sean soluciones del problema de Cauchy con $y_i(t_0) = y_0$, entonces la unicidad de la solución garantiza que $y_1|_{I_1 \cap I_2} = y_2|_{I_1 \cap I_2}$. Por lo tanto, podemos considerar la unión de todos los intervalos en estas condiciones, y así obtenemos un intervalo abierto $I_{t_0, y_0, \mu}$ en el que está definida la *solución máxima* $y_{t_0, y_0, \mu} : I_{t_0, y_0, \mu} \rightarrow \mathbb{R}^n$, que tiene la propiedad adicional de no ser prolongable a ningún intervalo mayor. El teorema anterior garantiza que

$$G = \bigcup_{(t_0, y_0, \mu) \in D} I_{t_0, y_0, \mu} \times \{t_0\} \times \{y_0\} \times \{\mu\}$$

es un abierto en $\mathbb{R} \times D$ donde está definida la *solución máxima* $y(t, t_0, y_0, \mu)$, la cual no es prolongable a ningún abierto mayor.

Resulta ilustrativo adoptar una notación que sugiera una interpretación física “típica” de los problemas de Cauchy en la que, en particular, la variable t se concibe como el tiempo:

Definición 7.6 Un *campo de velocidades* es una aplicación $\vec{v} : \mathbb{R} \times \Omega \rightarrow \mathbb{R}^n$ de clase C^k (con $k \geq 1$), donde Ω es un abierto en \mathbb{R}^n . La solución máxima $\vec{r} : G \rightarrow \Omega$ (donde G es un abierto en $\mathbb{R}^2 \times \Omega$) del problema de Cauchy

$$\left. \begin{aligned} \frac{\partial \vec{r}}{\partial t}(t, t_0, \vec{r}_0) &= \vec{v}(t, \vec{r}(t, t_0, \vec{r}_0)) \\ \vec{r}(t_0, t_0, \vec{r}_0) &= \vec{r}_0 \end{aligned} \right\}$$

recibe el nombre de *flujo* asociado al campo de velocidades.

Observemos que, aparte del cambio de notación, sólo hemos introducido dos pequeñas simplificaciones respecto de la situación general planteada en el teorema anterior, y es que, por una parte hemos omitido el parámetro μ (pero esto es meramente por aligerar la notación, y no hay ningún inconveniente en mantenerlo en todo lo que sigue) y que por otra parte hemos supuesto que el campo \vec{v} está definido en todo punto de Ω para todo instante t . Esto sucede en particular si \vec{v} no depende de t , en cuyo caso se dice que es un *campo de velocidades estacionario*.

Las observaciones precedentes a la definición justifican que el abierto G contiene a todas las ternas $(t_0, t_0, \vec{r}_0) \in \mathbb{R}^2 \times \Omega$, y que, para cada $(t_0, \vec{r}_0) \in \mathbb{R} \times \Omega$, el conjunto $I_{t_0, \vec{r}_0} = \{t \in \mathbb{R} \mid (t, t_0, \vec{r}_0) \in G\}$ es un intervalo abierto que contiene a t_0 , sobre el que está definida la función $\vec{r}_{t_0, \vec{r}_0} : I_{t_0, \vec{r}_0} \rightarrow \mathbb{R}^n$, que podemos ver como la trayectoria de un móvil que pasa por el punto \vec{r}_0 en el instante t_0 y cuya velocidad en cada instante es la que el campo de velocidades prescribe para su posición en ese instante, es decir: $\vec{v}_{t_0, \vec{r}_0}(t) = \vec{v}(t, \vec{r}_{t_0, \vec{r}_0}(t))$.

Las curvas \vec{r}_{t_0, \vec{r}_0} se llaman *curvas integrales* del campo de velocidades. La unicidad de la solución hace que dos curvas integrales que pasen por el mismo punto \vec{r}_1 en el mismo instante t_1 son iguales, pues ambas tienen que ser \vec{r}_{t_1, \vec{r}_1} .

Como ejemplo de aplicación del teorema 7.5 demostramos lo siguiente:

Teorema 7.7 Dadas dos funciones $\kappa, \tau : I \rightarrow \mathbb{R}$ de clase C^2 de modo que $\kappa \geq 0$ y un punto $t_0 \in I$, existe una curva regular $x :]t_0 - \epsilon, t_0 + \epsilon[\rightarrow \mathbb{R}^3$ parametrizada por el arco tal que κ y τ son respectivamente su curvatura y su torsión. La curva es única salvo isometrías.

DEMOSTRACIÓN: La unicidad nos la da el teorema 5.25. Para probar la existencia consideramos el sistema de ecuaciones diferenciales determinado por las fórmulas de Frenet:

$$T' = \kappa N, \quad N' = -\kappa T - \tau B, \quad B' = \tau N.$$

Se trata de un sistema de nueve ecuaciones diferenciales con incógnitas las nueve funciones coordenadas de T , N y B . Tomamos unas condiciones iniciales cualesquiera T_0, N_0, B_0 tales que formen una base ortonormal positivamente orientada. Por el teorema anterior existen unas únicas funciones (T, B, N) que satisfacen las ecuaciones. Un simple cálculo nos da

$$\begin{aligned} (TN)' &= \kappa NN - \kappa TT - \tau TB, & (TB)' &= \kappa NB + \tau TN, \\ (NB)' &= -\kappa TB - \tau BB + \tau NN & (TT)' &= 2\kappa TN, \\ (NN)' &= -2\kappa TN - 2\tau NB & (BB)' &= 2\tau NB. \end{aligned}$$

Vemos que las seis funciones TN , TB , NB , TT , NN , BB satisfacen un sistema de ecuaciones diferenciales con la condición inicial $(0, 0, 0, 1, 1, 1)$ que por otra parte es claro que tiene por solución a la función constante $(0, 0, 0, 1, 1, 1)$. La unicidad implica que (T, N, B) es una base ortonormal de \mathbb{R}^3 . Definimos

$$x(s) = \int_{t_0}^s T(s) ds.$$

Entonces es claro que $x'(s) = T(s)$, luego en particular x está parametrizada por el arco. Además $x''(s) = \kappa N$, luego κ es la curvatura de x . Un simple cálculo nos da que la torsión es τ . ■

Veamos ahora una aplicación a las variedades diferenciables. Sea $S \subset \mathbb{R}^m$ una variedad de dimensión n y $\alpha : I \rightarrow S$ una curva parametrizada por el arco. Sea $\alpha(s_0) = p$. Sea X una carta de S alrededor de p . Entonces α tiene asociada una representación en la carta $x(s)$ de modo que $\alpha(s) = X(x(s))$. Un vector arbitrario $w_0 \in T_p(S)$ es de la forma $w_0 = dX(x(s_0))(a_0)$. Teniendo en cuenta las ecuaciones (6.8) del capítulo anterior es claro que la solución $a(s)$ del problema

$$a'_k + \sum_{i,j=1}^n a_i \Gamma_{ij}^k x'_j = 0$$

con la condición inicial $a(s_0) = a_0$ determina un campo de vectores

$$w(s) = a_1(s) D_1 X(x(s)) + \cdots + a_n(s) D_n X(x(s))$$

con derivada covariante nula. Con esto casi tenemos demostrado el teorema siguiente:

Teorema 7.8 Sea S una variedad y $\alpha : I \rightarrow S$ una curva parametrizada por el arco. Sea $\alpha(s_0) = p$ y sea $w_0 \in T_p(S)$. Entonces existe un único campo $w : I \rightarrow \mathbb{R}^m$ tal que $w(s) \in T_{\alpha(s)}(S)$, $w(s_0) = w_0$ y $Dw = 0$. Lo llamaremos *transporte paralelo* de w_0 a través de α .

En realidad hemos probado la existencia de w en un entorno de s_0 . Vamos a justificar que la solución puede prolongarse a todo I . Para ello conviene observar que al ser $Dw = 0$ tenemos que $w'(s)$ es perpendicular a $T_{\alpha(s)}(S)$, luego $(ww)' = 2ww' = 0$, es decir, $\|w\|$ es constante. De aquí se sigue que el transporte paralelo es único: si w y w' son transportes paralelos de un mismo vector, entonces $w - w'$ es un transporte paralelo del vector nulo (porque su derivada covariante es la resta de las de los dos campos, luego es nula), luego es la aplicación constantemente nula.

Acabamos de usar la linealidad de la derivada covariante. Más en general, si tenemos dos vectores $w_0, w'_0 \in T_p(S)$ y ambos tienen transporte paralelo w y w' , entonces $\alpha w + \beta w'$ es un transporte paralelo de $\alpha w_0 + \beta w'_0$. Según lo que sabemos, existe un entorno de s_0 a lo largo del cual todos los vectores de una base de $T_p(S)$ tienen transporte paralelo, con lo que de hecho todos los vectores de $T_p(S)$ lo tienen.

Dado cualquier $s \in I$, es claro que el intervalo $[s_0, s]$ (o $[s, s_0]$ si $s < s_0$) puede cubrirse con un número finito de estos entornos donde existe transporte paralelo, de donde se sigue inmediatamente la existencia de transporte paralelo desde s_0 hasta s , luego el transporte paralelo existe sobre todo I . ■

Definición 7.9 Sea S una variedad y $\alpha : [s_0, s_1] \rightarrow S$ una curva parametrizada por el arco de extremos p y q . Según el teorema anterior, para cada vector $w_0 \in T_p(S)$ tenemos definido el transporte paralelo $w(s)$ a lo largo de α de modo que $w(s_0) = w_0$. Al vector $\text{tp}_{pq}^\alpha(w_0) = w(s_1) \in T_q(S)$ lo llamaremos *trasladado* de w_0 a lo largo de α . Esto nos define una aplicación $\text{tp}_{pq}^\alpha : T_p(S) \rightarrow T_q(S)$ a la que llamaremos *transporte paralelo* de $T_p(S)$ a $T_q(S)$ a lo largo de α .

Ejercicio: Probar que el transporte paralelo $\text{tp}_{pq}^\alpha : T_p(S) \rightarrow T_q(S)$ es una isometría.

7.2 Ecuaciones diferenciales de orden superior

Las ecuaciones diferenciales que aparecen con mayor frecuencia en física y en geometría son de orden 2. Afortunadamente, toda la teoría sobre existencia y unicidad que vamos a necesitar para ecuaciones diferenciales de orden superior se deduce inmediatamente del caso de orden 1. En efecto:

Teorema 7.10 Sea $f : D \subset \mathbb{R} \times \mathbb{R}^{nm} \times \mathbb{R}^k \rightarrow \mathbb{R}^n$ una función de clase C^k con $k \geq 1$ en un abierto D . Entonces la ecuación diferencial

$$\left. \begin{aligned} y^{(m)}(t) &= f(t, y, y', \dots, y^{(m-1)}, \mu) \\ y(t_0) &= y_0 \\ y'(t_0) &= y'_0 \\ \dots\dots\dots &\dots \\ y^{(m-1)}(t_0) &= y_0^{(m-1)} \end{aligned} \right\}$$

tiene solución única $y(t, \mu, t_0, y_0, y'_0, \dots, y_0^{(m-1)})$ de clase C^k en un entorno de cada punto $(t_0, \mu, t_0, y_0, y'_0, \dots, y_0^{(m-1)})$.

DEMOSTRACIÓN: Basta observar que el problema equivale al sistema de ecuaciones de primer orden

$$\left. \begin{aligned} y'(t) &= y_1 \\ y'_1(t) &= y_2 \\ \dots\dots\dots &\dots \\ y'_{m-2}(t) &= y_{m-1} \\ y'_{m-1}(t) &= f(t, y, y_1, \dots, y_{m-1}, \mu) \\ (y, y_1, \dots, y_{m-1})(t_0) &= (y_0, y'_0, \dots, y_0^{(m-1)}) \end{aligned} \right\}$$

donde hemos introducido las variables auxiliares y_i , que representan funciones con valores en \mathbb{R}^n . Todo este sistema se puede expresar como una única ecuación vectorial en las condiciones de la sección anterior. ■

Ejemplo Sea S una variedad diferenciable, sea $p \in S$ y $w \in T_p(S)$ un vector unitario. Sea X una carta alrededor de p , $p = X(x_0)$ y $w = dX(x'_0)$. Entonces existe una única curva $x(t)$ que verifica las ecuaciones (6.10) del capítulo anterior con las condiciones iniciales $x(0) = x_0$, $x'(0) = x'_0$. La curva $g(t) = X(x(t))$ es una geodésica de S parametrizada por el arco tal que $g(0) = p$ y $g'(0) = w$. Recíprocamente, la representación en la carta de cualquier geodésica que cumpla esto ha de ser solución de las ecuaciones (6.10), luego g es única (salvo cambio de parámetro). En resumen:

En una variedad, por cada punto pasa una única geodésica en cada dirección.

Por ejemplo, en el capítulo anterior probamos que los círculos máximos son geodésicas de la esfera. Puesto que por cada punto y en cada dirección pasa un círculo máximo, concluimos que los círculos máximos son las únicas geodésicas de la esfera. ■

Ejemplo 1 Consideremos la ecuación diferencial $x''(t) = -kx$, con $k > 0$.

Es fácil comprobar que $x(t) = a \cos(\sqrt{k}t + \phi)$ es una solución, para $a, \phi \in \mathbb{R}$ cualesquiera. Para comprobar que son las únicas soluciones posibles basta ver que, fijadas unas condiciones iniciales $x(t_0) = x_0$, $x'(t_0) = x'_0$ existe una solución de este tipo que las cumple. Ahora bien, esto equivale a que

$$(x_0, -x'_0/\sqrt{k}) = a(\cos(\sqrt{k}t_0 + \phi), \sin(\sqrt{k}t_0 + \phi)),$$

por lo que sólo hay que elegir a y ϕ de modo que a y $\sqrt{k}t_0 + \phi$ sean el módulo y un argumento del par del miembro izquierdo. ■

Ejemplo 2 Vamos a calcular todas las soluciones de la ecuación

$$y''(t) = \frac{k}{t} y'(t), \quad k \in \mathbb{R}, t > 0.$$

Obviamente las funciones constantes son soluciones de la ecuación. Si y no es constante existe un punto $t_0 > 0$ tal que $y'(t_0) \neq 0$. Llamemos $y_1 = y'(t)$. En un entorno de t_0 tenemos

$$\frac{y'_1(t)}{y_1(t)} = \frac{k}{t},$$

luego integrando entre t_0 y t queda

$$\log y_1(t) - \log y_1(t_0) = \log x^k,$$

de donde $y_1(t) = y_1(t_0)x^k$, es decir, $y'(t) = y'(t_0)x^k$. Integrando de nuevo concluimos que

$$y(t) = \begin{cases} \frac{y'(t_0)}{k+1} x^{k+1} + y(t_0) & \text{si } k \neq -1, \\ y'(t_0) \log t + y(t_0) & \text{si } k = -1. \end{cases}$$

Ahora es claro que las soluciones de la ecuación dada están todas definidas en $]0, +\infty[$ y vienen dadas por

$$y(t) = \begin{cases} Ax^{k+1} + B & \text{si } k \neq -1, \\ A \log t + B & \text{si } k = -1. \end{cases}$$

Estas expresiones incluyen las funciones constantes, que habíamos dejado aparte. ■

Ejemplo 3 Vamos a calcular todas las soluciones de la ecuación

$$z''(t) + 2a z'(t) + b z(t) = 0.$$

(Obviamente el 2 en el coeficiente de z' es irrelevante, y lo hemos incluido porque así la solución se expresa de forma ligeramente más simple.)

Para resolverla vamos a usar un “truco”, y es que vamos a suponer que los coeficientes $2a = u + vi$, $b = r + si$ son números complejos, al igual que la función $z(t) = x(t) + y(t)i$. Esto tiene sentido y no se sale de la teoría que hemos desarrollado, pues si sustituimos en la ecuación y separamos la parte real de la imaginaria lo que tenemos es el sistema

$$x''(t) = -ux'(t) + vy'(t) - rx(t) + sy(t),$$

$$y''(t) = -vx'(t) - uy'(t) - sx(t) - ry(t),$$

y sabemos que tiene que tener una única solución para cada conjunto de condiciones iniciales $z_0 = x_0 + y_0i$, $z'_0 = x'_0 + y'_0i$ (es decir, para las condiciones $z(0) = z_0$, $z'(0) = z'_0$, donde suponemos $t_0 = 0$ sin pérdida de generalidad).

Sólo necesitamos comprobar que cuando consideramos a una función vectorial $z(t) = (x(t), y(t))$ como función compleja $z(t) = x(t) + y(t)i$ se siguen cumpliendo las propiedades siguientes:

- $(z_1 + z_2)'(t) = z'_1(t) + z'_2(t)$.
- $(az)'(t) = az'(t)$ (para $a \in \mathbb{C}$).
- $(e^{at})' = ae^{at}$ (para $a \in \mathbb{C}$).

Por ejemplo, para demostrar la última propiedad expresamos $a = u + vi$, de modo que

$$e^{at} = (e^{ut} \cos vt, e^{ut} \sen vt),$$

y al derivar queda

$$\begin{aligned} (e^{at})' &= ue^{ut}(\cos vt, \sen vt) + e^{ut}(-v \sen vt, v \cos vt) \\ &= ue^{ut}(\cos vt + i \sen vt) + e^{ut}v(-\sen vt + i \cos vt) \\ &= ue^{ut}(\cos vt + i \sen vt) + e^{ut}iv(\cos vt + i \sen vt) = ue^{ut}e^{vit} + e^{ut}ive^{vt} = ae^{at}. \end{aligned}$$

En el fondo, estas propiedades se cumplen porque es posible definir una derivación de funciones de variable compleja totalmente análoga a la derivación de funciones de variable real, pero no vamos a entrar en ello, y para este ejemplo basta comprobar las propiedades anteriores.

Consideremos una posible solución de la forma $z(t) = e^{kt}$, con $k \in \mathbb{C}$. Se comprueba inmediatamente que satisface la ecuación diferencial si y sólo si k cumple la ecuación

$$k^2 + 2ak + b = 0,$$

es decir, si $k = -a \pm c$, donde $c = \sqrt{a^2 - b}$. (Observemos que c puede ser un número complejo aunque partamos de coeficientes reales, y ésta es la razón por la que hemos pasado al caso complejo.) Supongamos en primer lugar que $c \neq 0$, con lo cual hemos encontrado dos soluciones: $z = e^{-at}e^{\pm ct}$. Ahora bien, es fácil ver que cualquier combinación lineal de ellas

$$z = e^{-at}(Ae^{ct} + Be^{-ct}), \quad A, B \in \mathbb{C}$$

es también solución. Vamos a probar que éstas son todas las soluciones de la ecuación, lo cual equivale a probar que, para cada par de condiciones iniciales $z_0, z'_0 \in \mathbb{C}$ es posible encontrar A y B tales que la solución correspondiente las satisfaga. No perdemos generalidad si suponemos que $t_0 = 0$, en cuyo caso las condiciones iniciales son

$$A + B = z_0, \quad -a(A + B) + c(A - B) = z'_0.$$

Es claro que este sistema de ecuaciones tiene solución para en A y B , pero en lugar de calcularla es más práctico expresar la solución en términos de $A + B$ y $A - B$, así:

$$z(t) = e^{-at} \left((A + B) \frac{e^{ct} + e^{-ct}}{2} + (A - B) \frac{e^{ct} - e^{-ct}}{2} \right),$$

con lo que

$$z(t) = e^{-at} \left(z_0 \frac{e^{ct} + e^{-ct}}{2} + \left(\frac{az_0}{c} + \frac{z'_0}{c} \right) \frac{e^{ct} - e^{-ct}}{2} \right) \quad (7.1)$$

es una solución que cumple $z(0) = z_0$, $z'(0) = z'_0$. Una expresión alternativa para la solución se obtiene llamando $q = \sqrt{b - a^2}$, con lo que $c = iq$, y entonces las relaciones 4.42 nos dan que

$$z(t) = e^{-at} \left(z_0 \cos qt + \left(\frac{az_0}{q} + \frac{z'_0}{q} \right) \operatorname{sen} qt \right). \quad (7.2)$$

De este modo, si los coeficientes de la ecuación de partida son reales, al igual que las condiciones iniciales, tenemos que o bien (7.1) o bien (7.2) muestra que la solución que hemos obtenido es real, según que sea $a^2 > b$ (en cuyo caso $c \in \mathbb{R}$) o bien $a^2 < b$ (en cuyo caso $q \in \mathbb{R}$).

Nos falta considerar el caso en que $c = 0$, en el que sólo tenemos una solución $z = e^{-at}$, pero en este caso la ecuación se reduce a

$$z''(t) + 2a z'(t) + a^2 z(t) = 0,$$

y una comprobación rutinaria muestra que otra solución es $z(t) = te^{-at}$. A partir de ahí, el razonamiento precedente se adapta para justificar que la solución general es

$$z(t) = e^{-at}(A + Bt)$$

o, en términos de las condiciones iniciales:

$$z(t) = e^{-at}(z_0 + (az_0 + z'_0)t).$$

■

La ecuación de Euler-Cauchy Una *ecuación de Euler-Cauchy* es una ecuación diferencial de la forma

$$at^2 y''(t) + bty'(t) + cy(t) = 0,$$

donde a, b, c son constantes. Para resolverla basta hacer el cambio de variable $t = e^u$, de modo que

$$\frac{dy}{du} = \frac{dy}{dt}t, \quad \frac{d^2y}{du^2} = \frac{d^2y}{du^2}t^2 + \frac{dy}{dt}t = \frac{d^2y}{du^2}t^2 + \frac{dy}{du}.$$

Por lo tanto la ecuación se transforma en

$$a \left(\frac{d^2y}{du^2} - \frac{du}{du} \right) + b \frac{dy}{du} + cy = 0,$$

o también

$$ay''(u) + (b - a)y'(u) + cy(u) = 0,$$

que es del tipo resuelto en el ejemplo anterior.

Por ejemplo, las soluciones de la ecuación

$$t^2 y''(t) + 2ty'(t) = n(n+1)y(t)$$

se corresponden a través del cambio $t = e^u$ con las de la ecuación

$$y''(u) + y'(y) - n(n+1)y(u) = 0$$

que, con la notación del ejemplo anterior, corresponden al caso en que $a = 1/2$, $b = -n^2 - n$, $c = n + 1/2$. Por lo tanto, sus soluciones son las combinaciones lineales de $e^{-u/2}e^{\pm(n+1/2)u}$, es decir, de e^{nu} y $e^{-(n+1)u}$ y al deshacer el cambio de variable resulta que la solución general es

$$At^n + Bt^{-(n+1)}.$$

■

Campos de fuerzas El análogo de segundo orden a un campo de velocidades sería un campo de aceleraciones, pero en física resulta más natural hablar de *campos de fuerzas*. Es frecuente que la fuerza $F_t(x)$ que actúa sobre un cuerpo en un instante dado dependa únicamente de su posición en el espacio, con lo que tenemos una función $F : I \times D \rightarrow \mathbb{R}^n$. El campo de fuerzas será *estacionario* si no depende del tiempo. De acuerdo con la segunda ley de Newton, la trayectoria $x(t)$ de un cuerpo de masa m sometido a este campo vendrá determinada por la ecuación diferencial

$$F_t = m x''.$$

La solución depende de la posición inicial x_0 y la velocidad inicial v_0 . ■

7.3 Aplicaciones

7.3.1 Funciones armónicas con simetría esférica

Definición 7.11 El *laplaciano* de una función $f : A \rightarrow \mathbb{R}$ de clase C^2 en un abierto $A \subset \mathbb{R}^n$ es la función

$$\Delta f = \frac{\partial^2 f}{\partial x_1^2} + \cdots + \frac{\partial^2 f}{\partial x_n^2}.$$

La función f es *armónica* si cumple la ecuación en derivadas parciales $\Delta f = 0$.

Por ejemplo, es obvio que toda función afín $f(x) = a_0 + a_1 x_1 + \cdots + a_n x_n$ es armónica en \mathbb{R}^n . Aquí vamos a determinar todas las funciones armónicas con simetría esférica:

Teorema 7.12 *Las únicas funciones armónicas en \mathbb{R}^n de la forma $g(\|x\|)$ son las de la forma*

$$f(x) = \begin{cases} \frac{A}{\|x\|^{n-2}} + B & \text{si } n \neq 2, \\ A \log \|x\| + B & \text{si } n = 2. \end{cases}$$

DEMOSTRACIÓN: Sea f una función de la forma indicada. La función g es de clase C^2 en su dominio, pues f lo es y $g(r) = f(r, 0, \dots, 0)$. Por consiguiente

$$\frac{\partial f}{\partial x_i} = \frac{dg}{dr} \frac{x_i}{\|x\|}, \quad \frac{\partial^2 f}{\partial x_i^2} = \frac{d^2 g}{dr^2} \frac{x_i^2}{\|x\|^2} + \frac{dg}{dr} \left(\frac{1}{\|x\|} - \frac{x_i^2}{\|x\|^3} \right),$$

luego

$$\Delta f = \frac{d^2 g}{dr^2} + \frac{dg}{dr} \frac{n-1}{\|x\|} = 0.$$

Esta ecuación se cumple para todo $x \neq 0$ en el dominio de f , de donde se sigue claramente que

$$\frac{d^2 g}{dr^2} + \frac{dg}{dr} \frac{n-1}{r} = 0$$

para todo $r \neq 0$ en el dominio de g . En el ejemplo 2 de la sección anterior hemos visto que las únicas soluciones de esta ecuación son las de la forma

$$g(r) = \begin{cases} \frac{A}{r^{n-2}} + B & \text{si } n \neq 2 \\ A \log r + B & \text{si } n = 2 \end{cases}$$

de donde se sigue que f tiene la forma indicada. ■

7.3.2 La primera ley de Kepler

La *ley de la gravitación universal* de Newton afirma que la fuerza con que se atraen mutuamente dos cuerpos es directamente proporcional a sus masas e inversamente proporcional al cuadrado de la distancia que los separa. Consideremos una región del espacio donde haya un cuerpo S de masa M tan grande que la masa de cualquier otro cuerpo en las proximidades resulte despreciable. Es el caso del Sol, rodeado de planetas de masa insignificante a su lado, o de la Tierra y sus alrededores. En tal caso podemos suponer que la única fuerza que actúa sobre un cuerpo es la provocada por S . En efecto, cuando dejamos caer un objeto nuestro cuerpo lo atrae por gravedad, pero esta atracción es completamente inapreciable frente a la gravitación terrestre. Igualmente, Júpiter atrae gravitatoriamente a la Tierra, pero la fuerza con que lo hace es insignificante frente a la del Sol.

Si tomamos un sistema de referencia con origen en el punto donde se halla el objeto masivo S , la fuerza que experimenta otro cuerpo de masa m situado en un punto x viene dada por

$$F = -\frac{GMm}{\|x\|^3} x,$$

donde $G = 6.672 \cdot 10^{-11} \text{N} \cdot \text{m}^2/\text{Kg}^2$ es la *constante de gravitación universal*. El hecho de que su valor sea tan pequeño hace que la gravedad no se manifieste salvo en presencia de grandes masas, como las estrellas y los planetas.

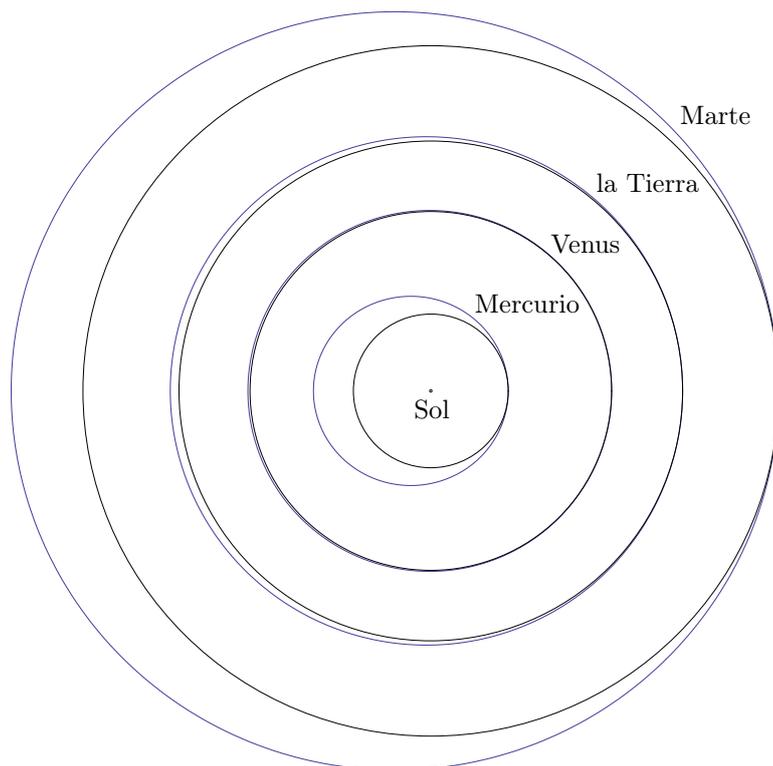
Tras estudiar minuciosamente una gran cantidad de observaciones astronómicas sobre Marte, en 1610 Kepler publicó su *astronomía nova*, donde concluía que los planetas se mueven siguiendo órbitas elípticas, de modo que el Sol ocupa uno de los focos. Ésta es la primera ley de Kepler.

La figura de la página siguiente muestra un mapa a escala del Sol y las órbitas de los primeros planetas del sistema solar: Mercurio, Venus, la Tierra y Marte. El diámetro del Sol está representando también a escala. Para cada planeta está representada su órbita y la órbita circular de radio igual al perihelio (el punto de la órbita más cercano al Sol). En la figura los perihelios están alineados, pero esto no sucede realmente.

Newton mostró que la primera ley de Kepler y muchos otros hechos sobre el movimiento de los astros pueden deducirse a partir de las leyes básicas de la dinámica y de su ley de gravitación. Vamos a comprobar que las trayectorias de

los objetos sometidos a un campo de fuerzas como el que estamos considerando son rectas o secciones cónicas.

En primer lugar, es claro que si un cuerpo se encuentra en un punto x_0 en las proximidades del Sol con una velocidad v_0 , su trayectoria no saldrá del plano determinado por los vectores x_0 y v_0 (o de la recta que determinan, si son linealmente dependientes). Por ello podemos tomar el sistema de referencia de modo que el eje Z sea perpendicular a este plano, con lo que la trayectoria cumplirá $z = 0$ y podemos trabajar únicamente con las coordenadas (x, y) .



Conviene introducir coordenadas polares, $r = (x, y) = (\rho \cos \theta, \rho \sin \theta)$. En general, cuando la trayectoria de un móvil viene dada en coordenadas polares por unas funciones $\rho(t)$, $\theta(t)$, se llama *velocidad angular* a la derivada $\omega = \theta'$. La segunda derivada $\alpha = \omega' = \theta''$ se conoce como *aceleración angular*.

La ecuación de Newton puede expresarse en términos de la *cantidad de movimiento*, definida como $p = mv$, donde $v = r'$ es la velocidad del móvil y m es su masa. Si ésta es constante, la segunda ley de Newton afirma que

$$F = \frac{dp}{dt},$$

donde F es la fuerza total que actúa sobre el cuerpo. En particular, la cantidad de movimiento de un cuerpo sobre el que no actúa ninguna fuerza permanece constante.

Existen magnitudes análogas a la cantidad de movimiento y la fuerza en coordenadas polares. Se llama *momento angular* de un móvil a la magnitud $L = r \times p = m r \times v$. Si la trayectoria está contenida en un plano, el vector L es perpendicular a él. Veamos su expresión en coordenadas polares. Para ello calculamos:

$$v = \rho'(\cos \theta, \sin \theta) + \rho\omega(-\sin \theta, \cos \theta), \quad (7.3)$$

de donde

$$L = m\rho(\cos \theta, \sin \theta, 0) \times \rho\omega(-\sin \theta, \cos \theta, 0) = (0, 0, m\rho^2\omega).$$

Para un movimiento plano podemos abreviar y escribir $L = m\rho^2\omega$. Se define el *momento* de una fuerza F que actúa sobre un móvil de posición r como $M = r \times F$. Es claro que si sobre un cuerpo actúan varias fuerzas el momento de la fuerza resultante es la suma de los momentos. La versión angular de la segunda ley de Newton es⁴

$$\frac{dL}{dt} = mv \times v + r \times ma = r \times F = M,$$

es decir, el momento total que actúa sobre un móvil es la derivada de su momento angular. En particular, si un cuerpo está libre de toda fuerza su momento angular permanece constante. Sin embargo, el momento angular se conserva incluso en presencia de fuerzas, con tal de que la fuerza resultante sea proporcional a la posición, como ocurre en el caso de la fuerza gravitatoria que el Sol ejerce sobre los planetas.

Esto ya nos da una información sobre el movimiento de los planetas: puesto que $m\rho^2\omega$ ha de ser constante, los planetas giran más rápidamente cuando están más cerca del Sol.

Veamos ahora la expresión de la segunda ley de Newton para la gravitación en coordenadas polares. Calculamos la aceleración

$$a = v' = (\rho'' - \rho\omega^2)(\cos \theta, \sin \theta) + (2\rho'\omega + \rho\alpha)(-\sin \theta, \cos \theta).$$

La fuerza gravitatoria es

$$F = -\frac{GMm}{\rho^2}(\cos \theta, \sin \theta).$$

Teniendo en cuenta que los vectores $(\cos \theta, \sin \theta)$ y $(-\sin \theta, \cos \theta)$ son ortogonales, la ecuación $F = ma$ equivale a las ecuaciones

$$\rho'' - \rho\omega^2 = -\frac{GM}{\rho^2}, \quad 2\rho'\omega + \rho\alpha = 0. \quad (7.4)$$

La solución es complicada, pero ahora estamos interesados únicamente en la forma de la trayectoria, aunque sea con otra parametrización. Por ello, en lugar

⁴Aquí usamos que la derivada del producto vectorial de dos vectores cumple una regla análoga a la derivada de un producto de funciones escalares, regla que no es difícil de probar en general.

de calcular las funciones $\rho(t)$ y $\theta(t)$ calcularemos la parametrización $\rho(\theta)$. Por supuesto hay un caso en el que esta parametrización es imposible: Si $\omega_0 = 0$ (lo cual equivale a que la velocidad inicial v_0 sea nula o proporcional a r_0) entonces es fácil ver que la solución de las ecuaciones es una recta de la forma $(\rho(t), \theta_0)$, donde ρ está determinada por la ecuación

$$\rho'' = -\frac{GM}{\rho^2}$$

con las condiciones iniciales ρ_0 y ρ'_0 , determinadas a su vez por r_0 y v_0 . En efecto, una trayectoria de este tipo cumple $\omega = \alpha = 0$ y satisface trivialmente las ecuaciones. En definitiva, el cuerpo se aleja del Sol en línea recta o bien cae sobre él. Hemos probado que si la trayectoria de un móvil cumple $\omega(t) = 0$ en un instante t entonces $\omega = 0$ en todo instante, luego, recíprocamente, si $\omega_0 \neq 0$ entonces ω no se anula nunca. Esto hace que la función $\theta(t)$ sea un cambio de parámetro, luego podemos considerar la reparametrización $\rho(\theta)$. Claramente

$$\rho' = \rho'_\theta \omega \quad \text{y} \quad \rho'' = \rho''_\theta \omega^2 + \rho'_\theta \alpha.$$

Sustituimos estas igualdades en (7.4) y eliminamos α en la primera usando la segunda. El resultado es la ecuación

$$\left(\rho'' - \frac{2\rho'^2}{\rho} - \rho \right) \omega^2 = -\frac{GM}{\rho^2},$$

donde ahora todas las derivadas son respecto de θ y no respecto de t . No obstante, la presencia de ω nos obliga a considerar ambos miembros como funciones de t (pues la expresión entre paréntesis es una función de θ compuesta con la función $\theta(t)$).

Sin embargo, al multiplicar ambos miembros por $m^2\rho^4$ queda

$$\frac{1}{\rho^2} \left(\rho'' - \frac{2\rho'^2}{\rho} - \rho \right) = -\frac{GMm^2}{L^2}, \quad (7.5)$$

donde $L = m\rho^2\omega$ es una constante, luego todo el segundo miembro es constante y la igualdad sigue siendo válida si consideramos el primer miembro como función de θ .

Si r es una recta y F un punto exterior, la cónica de directriz r y foco F está formada por los puntos tales que la razón entre las distancias a r y a F es constante (y recibe el nombre de *excentricidad* de la cónica). Toda cónica que no sea una circunferencia es de esta forma.⁵ Si suponemos que el foco es el origen y la directriz es la recta vertical $x = p > 0$, entonces la distancia de un punto de coordenadas polares (ρ, θ) a la directriz es $|p - \rho \cos \theta|$, luego la ecuación de la cónica de excentricidad ϵ es

$$\frac{\rho}{p - \rho \cos \theta} = \epsilon.$$

⁵Véase [G, sección 10.1].

En principio faltaría un valor absoluto. Si $\epsilon \leq 1$ tenemos una elipse o una parábola y podemos suprimirlo, pues la curva queda al mismo lado de la directriz que el foco. Si $\epsilon > 1$ tenemos una hipérbola, y al suprimir el valor absoluto estamos quedándonos con una de sus ramas (y hemos de restringir θ al intervalo en el que el denominador es positivo). Lo hacemos así porque los cuerpos que se mueven siguiendo trayectorias hiperbólicas tienen al Sol en el foco correspondiente a la rama que siguen o, dicho de otro modo, que la rama que eliminamos no va a ser solución de la ecuación diferencial. Despejando ρ y llamando $r = p\epsilon$ queda

$$\rho = \frac{r}{1 + \epsilon \cos(\theta + k)}.$$

La introducción de la constante k equivale a girar la cónica, de modo que ahora la directriz es arbitraria. Cualquier curva de esta forma es una cónica de excentricidad ϵ . Además esta expresión incorpora también a las circunferencias, para las que $\epsilon = 0$. Es fácil calcular

$$\rho' = \frac{\epsilon \rho^2}{r} \operatorname{sen}(\theta + k), \quad \rho'' = \frac{2\epsilon^2 \rho^3}{r^2} \operatorname{sen}^2(\theta + k) + \frac{\epsilon \rho^2}{r} \cos(\theta + k).$$

Al sustituir en el miembro izquierdo de (7.5) se obtiene sin dificultad el valor $-1/r$, luego concluimos que la cónica

$$\rho = \frac{L^2}{GMm^2} (1 + \epsilon \cos(\theta + k))^{-1}$$

satisface (7.5). Sólo queda probar que éstas son las únicas soluciones posibles o, lo que es equivalente, que hay una solución de esta forma cualesquiera que sean las condiciones iniciales $\rho_0, \theta_0, \rho'_0, \omega_0$. Basta ver que las ecuaciones

$$\rho_0 = \frac{L^2}{GMm^2} (1 + \epsilon \cos(\theta_0 + k))^{-1}, \quad \rho'_0 = \frac{GMm^2 \rho_0^2}{L^2} \epsilon \operatorname{sen}(\theta_0 + k)$$

tienen solución en $\epsilon > 0, k$ para todos los valores de $\rho_0, \theta_0, \rho'_0, \omega_0$, pero sustituyendo $L = m\rho_0^2\omega_0$ estas ecuaciones equivalen a

$$\epsilon(\cos(\theta_0 + k), \operatorname{sen}(\theta_0 + k)) = \left(\frac{\rho_0^3 \omega_0^2}{GM} - 1, \frac{\rho'_0 \rho_0^2 \omega_0^2}{GM} \right),$$

que obviamente tienen solución. Con esto hemos probado que las trayectorias de los objetos sometidos a la atracción de una masa fija puntual son rectas o secciones cónicas.

Si elegimos el sistema de referencia de modo que el perihelio corresponda a $\theta = 0$, entonces $k = 0$. Si v_0 es la velocidad del planeta el en perihelio y ρ_0 su distancia al Sol, entonces

$$\rho_0 = \frac{L^2}{GMm^2} (1 + \epsilon)^{-1} = \frac{m^2 \rho_0^2 v_0^2}{GMm^2} (1 + \epsilon)^{-1},$$

luego la excentricidad de la órbita en función de las condiciones iniciales en el perihelio es

$$\epsilon = \frac{\rho_0 v_0^2}{GM} - 1.$$

Para órbitas elípticas, el *afelio* o punto de mayor distancia al Sol se alcanza cuando $\theta = \pi$, luego

$$\rho_1 = \frac{L^2}{GMm^2} (1 - \epsilon)^{-1} = \rho_0 \frac{1 + \epsilon}{1 - \epsilon}.$$

El semieje mayor de la órbita es

$$a = \frac{\rho_0 + \rho_1}{2} = \frac{\rho_0}{1 - \epsilon}. \quad \blacksquare$$

7.3.3 El péndulo de Foucault

Supongamos que desplazamos un péndulo de su posición de equilibrio y lo dejamos caer sin ninguna velocidad inicial. Teóricamente oscilará siempre en el mismo plano, pero si tiene suficiente masa como para que se haga sensible la fuerza de Coriolis debida a la rotación de la Tierra, ya no será así, sino que el plano de la oscilación sufrirá un movimiento de precesión.

Aunque el hecho ya era conocido, la primera demostración pública de este fenómeno la realizó el físico y astrónomo francés Jean Bernard Léon Foucault, y por ello los péndulos sometidos a la atracción gravitatoria y a la aceleración de Coriolis se llaman *péndulos de Foucault*.

El primer péndulo de Foucault fue instalado en 1851 en la cúpula del *Panteón de París*, estaba sujeto de un cable de acero de 67 m de largo y 1.4 mm de diámetro y pesaba 28 kg. Osciló durante 6 horas con una amplitud máxima de 6 m y un periodo de 16.5 s. La velocidad de precesión fue de 11° por hora.

Vamos a analizar este resultado. Tomemos un sistema de referencia con origen en el punto donde se sujeta la cuerda del péndulo y supongamos que lo soltamos desde una posición $(0, x_0, z_0)$. Si de momento prescindimos de la fuerza de Coriolis, sobre el péndulo actúan dos fuerzas, su peso P , dirigido hacia abajo, y la tensión del cable T , en la dirección de éste. El peso de un objeto (es decir, la fuerza con la que es atraído por la Tierra) depende en principio de su distancia a la superficie terrestre, pero para variaciones de altura pequeñas con respecto al radio de la Tierra podemos suponerlo constante, de modo que

$$P = (0, 0, -mg),$$

donde la constante $g = G \frac{M_T}{R_T^2} = 9.8 \text{ N/kg}$ depende de la constante de gravitación universal G y la masa y el radio de la Tierra. Si $r(t) = (x(t), y(t), z(t))$ es la posición del péndulo en el instante t , la tensión de la cuerda será

$$\left(-\frac{T}{l}x, -\frac{T}{l}y, -\frac{T}{l}z \right),$$

donde $l = \sqrt{x^2 + y^2 + z^2}$ es la longitud de la cuerda, que introducimos en la expresión para que $T(t)$ sea precisamente el módulo de la tensión. Por lo tanto, la fuerza total que actúa sobre el péndulo (sin contar la fuerza de Coriolis) es

$$F = \left(-\frac{T}{l}x, -\frac{T}{l}y, -\frac{T}{l}z - mg \right).$$

La segunda ley de Newton proporciona el sistema de ecuaciones diferenciales que determinan el movimiento del péndulo:

$$mx''(t) = -\frac{T(t)}{l}x(t), \quad my''(t) = -\frac{T(t)}{l}y(t), \quad mz''(t) = -\frac{T(t)}{l}z(t) - mg.$$

No estamos en condiciones de resolver este sistema, pero pensemos en la función (desconocida) $T(t)$. Podemos calcular su valor en los instantes en los que el péndulo se encontraba a su altura máxima y mínima. En los instantes de altura mínima, es decir, cuando el cable se encontraba en posición vertical, la tensión contrarrestaba exactamente al peso, por lo que $T = 29 \cdot 9.8 = 274.4 \text{ N}$. En los instantes de altura máxima, en los que la velocidad del péndulo era 0, la tensión compensaba la componente normal del peso, con lo que resulta $T = 273.3 \text{ N}$. Éstos son los valores máximo y mínimo de $T(t)$ y, puesto que su diferencia es de 1.1 N , el error que cometemos si suponemos que toma siempre el valor $T = mg = 274.4 \text{ N}$ se mantiene siempre por debajo de un 0.4% .

En general, siempre que consideremos un péndulo con ángulo de oscilación pequeño, obtendremos una aproximación razonable si suponemos que $T \approx mg$, con lo que la tercera ecuación diferencial tiene solución trivial $z(t) = -l$ (lo que indica que estamos despreciando las variaciones en altura del péndulo) y las otras dos se convierten en

$$x''(t) = -\frac{g}{l}x(t), \quad y''(t) = -\frac{g}{l}y(t),$$

que son dos ecuaciones independientes, que podemos resolver por separado.

No perdemos generalidad si suponemos que el péndulo parte de la posición $(0, y_0)$ con velocidad inicial $x'_0 = y'_0 = 0$. (Para asegurarse de que la velocidad inicial era nula, Foucault sujetó el péndulo con un hilo y lo quemó una vez se hubo estabilizado.) Entonces la primera ecuación diferencial tiene solución trivial $x(t) = 0$, lo que significa que el péndulo oscilará en el plano YZ . La segunda ecuación es del tipo considerado en el ejemplo 1 de la página 287, y su solución es, pues, de la forma

$$y(t) = a \cos \left(\sqrt{\frac{g}{l}} t \right),$$

donde a es la amplitud máxima del péndulo. Es claro entonces que el tiempo que tarda el péndulo en volver a su posición inicial es⁶

$$T = 2\pi \sqrt{\frac{l}{g}}.$$

⁶Hasta ahora llamábamos T al módulo de la tensión, pero esta cantidad no va a aparecer de nuevo, así que a partir de ahora reservamos la T para referirnos al periodo de oscilación.

Es de destacar que el periodo de oscilación de un péndulo, para oscilaciones de amplitud pequeña, no depende de dicha amplitud, sino únicamente de la longitud de la cuerda que lo sostiene. Esta característica es el fundamento del funcionamiento de los relojes de péndulo.

En el caso del péndulo de París, el valor que obtenemos según nuestros cálculos es $T = 2\pi\sqrt{67/9.8} = 16.43$ s, que encaja muy bien con el valor experimental de 16.5 s.

Veamos ahora cómo incorporamos el efecto de la fuerza de Coriolis. En la página 233 obtuvimos que su valor es $\vec{f}_{\text{cor}} = -2m\vec{\omega} \times \vec{v}_r$, donde \vec{v}_r es la velocidad de un móvil respecto de la superficie terrestre, $\vec{\omega}$ el vector que apunta en la dirección del eje de rotación de la Tierra y de módulo ω igual a la velocidad angular de rotación de la Tierra. Concretamente,

$$\omega = \frac{2\pi}{24 \cdot 60 \cdot 60} = 7.2722 \cdot 10^{-5} \text{ rad/s.}$$

Podemos suponer que el sistema de coordenadas respecto al cual estamos describiendo el movimiento del péndulo tiene su eje X apuntando hacia el Norte, en cuyo caso el vector $\vec{\omega}$ tiene coordenadas $(\omega \cos \lambda, 0, \omega \sin \lambda)$, donde λ es la latitud (así, si estamos en el ecuador $k = (1, 0, 0)$, mientras que si estamos cerca del polo Norte $k \approx (0, 0, 1)$). La latitud de París es $\lambda = 48^\circ 50' 46.5''$. Respecto de nuestro sistema de referencia, la fuerza de Coriolis es

$$\begin{aligned} F_{\text{cor}} &= -2m\omega(\cos \lambda, 0, \sin \lambda) \times (x', y', z') = \\ &2m\omega(y' \sin \lambda, -x' \sin \lambda + z' \cos \lambda, -y' \cos \lambda). \end{aligned}$$

El mayor valor que toma el módulo de y' (en nuestro cálculo sin considerar la fuerza de Coriolis) es $a\sqrt{g/l}$, que en nuestro ejemplo es 2.3 m/s, luego el valor máximo de la componente vertical de la fuerza de Coriolis es

$$2m\omega y'_{\text{max}} \cos \lambda = 0.00615 \text{ N,}$$

que es totalmente despreciable en comparación con el peso de 274.4 N. Por otra parte, nuestro supuesto de que la tensión del cable es constante conlleva considerar z constante, con lo que, en coherencia con nuestros supuestos precedentes, debemos considerar $z' = 0$. En definitiva, una aproximación razonable para la fuerza de Coriolis es

$$\vec{f}_{\text{cor}} = 2m\omega(y' \sin \lambda, -x' \sin \lambda, 0),$$

lo que nos lleva a modificar como sigue las ecuaciones del movimiento que habíamos obtenido:

$$x''(t) = 2\omega y'(t) \sin \lambda - \frac{g}{l}x(t), \quad y''(t) = -2\omega x'(t) \sin \lambda - \frac{g}{l}y(t).$$

(Notemos que la m desaparece porque las ecuaciones resultan de dividir entre m las dadas por la segunda ley de Newton.) Ahora agrupamos las dos incógnitas

en una única función compleja $z(t) = x(t) + y(t)i$, con lo que el sistema de ecuaciones se combina en una única ecuación

$$z''(t) + 2i\omega \operatorname{sen} \lambda z'(t) + (g/l)z = 0,$$

que se corresponde con la que hemos estudiado en el ejemplo 3 de la página 288. Con la notación empleada allí, $q = \sqrt{g/l + \omega^2 \operatorname{sen}^2 \lambda}$, luego la solución (correspondiente a una velocidad inicial nula) es

$$z(t) = e^{-i\omega \operatorname{sen} \lambda t} z_0 \left(\cos qt + i \frac{\omega \operatorname{sen} \lambda}{q} \operatorname{sen} qt \right).$$

Para interpretar esta expresión supongamos por un momento que no estuviera la exponencial inicial. Es fácil ver que la expresión entre paréntesis es una elipse cuyos ejes son los ejes de coordenadas, el semieje mayor mide 1 y el menor $(\omega/q) \operatorname{sen} \lambda$. La multiplicación por z_0 gira la elipse un ángulo igual al argumento de z_0 , con lo que su semieje mayor pasa a estar en la dirección de z_0 , es decir, del desplazamiento inicial del péndulo, y cuyos semiejes están ahora multiplicados por la amplitud a de este desplazamiento inicial (6 m en el ejemplo que estamos considerando). En total, esta parte de la solución corresponde a una oscilación que ya no es plana, sino elíptica, con un diámetro de 12 metros y el otro de 1.7 mm. En la práctica esto no se distingue de una oscilación plana, con periodo $T = 2\pi/q \approx 2\pi\sqrt{l/g}$ (en nuestro ejemplo $T = 2\pi/q = 16.43$ no se distingue del periodo que habíamos calculado sin considerar la fuerza de Coriolis, ya que la corrección es insignificante).

Así pues, si no estuviera la primera exponencial, la solución que hemos obtenido sería prácticamente la misma que la que ya teníamos, pero dicho factor se interpreta como un giro de la elipse (es decir, del plano de oscilación) con velocidad angular $\omega_\lambda = \omega \operatorname{sen} \lambda$ (en sentido horario en el hemisferio Norte y antihorario en el hemisferio Sur). En nuestro ejemplo, la velocidad angular es

$$\omega_\lambda = 5.47558 \cdot 10^{-5} \text{ rad/s} = 11^\circ 17' 39.1''/\text{h},$$

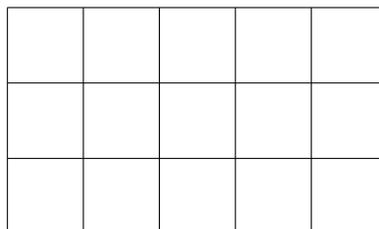
que cuadra con el resultado experimental. ■

Capítulo VIII

Teoría de la medida I

El problema de la medida consiste en definir un análogo en dimensiones superiores del concepto de longitud de un segmento. Cuando decimos que la longitud de un segmento es 5, esto significa que en él “cabem” cinco segmentos de longitud 1.

Similarmente, si el rectángulo de la figura tiene lados de longitudes 5 y 3, decimos que su área es 15, porque en él “cabem” 15 cuadrados de lado unitario. En otras palabras, tomamos como unidad de área la de un cuadrado de lado unitario, y entonces el área de una figura es el número de cuadrados unitarios que “cabem” en ella.



Ahora bien, esa idea de área que acabamos de describir no es lo suficientemente precisa como para que sea aceptable. Por una parte, cuando hablamos del número de cuadrados que “cabem” en una figura, se sobrentiende que éstos no deben solaparse, pero el concepto de “no solaparse” no coincide exactamente con el de “ser disjuntos”, pues los cuadrados de la figura no son disjuntos a menos que adoptemos unos convenios muy particulares sobre qué puntos de su frontera consideramos que pertenecen a cada cuadrado y cuáles no. Por otra parte, las figuras que pueden expresarse como uniones finitas de cuadrados que no se solapan son muy pocas, luego, con la definición de área que hemos dado, no tendría sentido plantearse, por ejemplo, cuál es el área de un triángulo. Sin embargo, el concepto de “área de un triángulo” tiene un sentido intuitivo muy claro. Por ejemplo, si trazamos una diagonal en el rectángulo de la figura, éste queda dividido en dos triángulos rectángulos congruentes, por lo que es razonable afirmar que cada uno de ellos tiene la mitad de área del rectángulo, es decir, que su área es de 7.5 unidades. Observemos que este argumento se basa en la misma idea original de “número de cuadrados que cabem”, pues, por ejemplo, si queremos fabricar una alfombra en forma de triángulo rectángulo de base 5 m y altura 3 m, podemos afirmar que necesitaremos 7.5 m^2 de tela, pues con el doble de esa cantidad podemos fabricar dos alfombras de esas medidas.

Si analizamos el argumento que nos ha llevado a concluir que el área del triángulo “debe ser” 7.5, veremos que utiliza implícitamente los postulados siguientes:

- a) El área de un cuadrado de lado unitario es 1.
- b) Dos figuras congruentes tienen la misma área.
- c) El área de la unión de dos figuras disjuntas es la suma de las áreas.

En realidad hemos usado algo más fuerte que c), pues los dos triángulos en que se descompone el rectángulo al trazar su diagonal no son disjuntos, sino que lo máximo que podemos decir es que “no se solapan” en el sentido de que su intersección es un segmento de recta (y lo mismo vale para el razonamiento que nos ha llevado a concluir que el área del rectángulo es 15, pues los cuadrados en que lo hemos descompuesto no son disjuntos).

Estas salvedades se traducen en una serie de problemas técnicos que tenemos que resolver si queremos llegar a una definición satisfactoria de área, pero antes de enfrentarnos al problema observemos que éste tiene su análogo en dimensiones superiores. Cuando consideramos figuras tridimensionales el concepto análogo al de área es el de “volumen”, de modo que la idea aproximada que pretende capturar este concepto es que el volumen de una figura es el número de cubos de arista unitaria que “cabén” en ella, con los matices correspondientes de que los cubos no deben solaparse, etc. El concepto análogo para figuras de dimensión arbitraria n recibe el nombre de “medida”, de manera que uno de los objetivos principales de este capítulo será el de dar una definición de medida de una figura n -dimensional de modo que se satisfagan las generalizaciones obvias de las propiedades a), b), c) anteriores.

8.1 La medida de Jordan

En esta sección daremos una primera solución, de carácter provisional, al problema de la medida, a la vez que despacharemos los principales problemas técnicos que plantea la definición, referentes al solapamiento de figuras. Para ello empezamos definiendo una familia de figuras especialmente simples, pero que bastan para dar una definición rudimentaria de “medida”:

Definición 8.1 Una *celda* en \mathbb{R}^n es un producto $C = \prod_{i=1}^n I_i$, donde cada I_i es un intervalo acotado en \mathbb{R} , es decir, un conjunto de la forma $[a_i, b_i]$, $]a_i, b_i[$, $[a_i, b_i[$ o $]a_i, b_i]$. Admitimos la posibilidad de que $a_i = b_i$, con lo que $I_i = \{a_i\}$ o bien $I_i = \emptyset$, y en este último caso $C = \emptyset$. Definimos el *contenido* de la celda C como $|C| = \prod_{i=1}^n (b_i - a_i)$.

Por ejemplo, las celdas en \mathbb{R}^2 son los rectángulos de lados paralelos a los ejes, los segmentos paralelos a los ejes, los puntos y el conjunto vacío.

Es inmediato que si $C = \prod_{i=1}^n I_i \subset C' = \prod_{i=1}^n I'_i$ son celdas no vacías, entonces $I_i \subset I'_i$ para todo i , luego $|C| \leq |C'|$.

En diversas ocasiones usaremos el siguiente argumento elemental de continuidad:

Teorema 8.2 *Si C es una celda en \mathbb{R}^n y $\epsilon > 0$, entonces existe una celda abierta C' tal que $C \subset C'$ y $|C'| - |C| < \epsilon$. Si además $|C| > 0$, existe una celda cerrada C'' tal que $C'' \subset C$ y $|C| - |C''| < \epsilon$.*

DEMOSTRACIÓN: Sea $C = \prod_{i=1}^n \|a_i, b_i\|$, donde $\|a_i, b_i\|$ representa a cualquiera de los intervalos $]a_i, b_i[$, $[a_i, b_i]$, $]a_i, b_i]$ o $[a_i, b_i[$. Sea $f : [0, +\infty[\rightarrow \mathbb{R}$ la función dada por

$$f(\delta) = \prod_{i=1}^n (b_i + \delta - (a_i - \delta)).$$

Claramente es continua y $f(0) = |C|$, luego existe un $\delta > 0$ tal que $f(\delta) < |C| + \epsilon$, y entonces la celda $C' = \prod_{i=1}^n]a_i - \delta, b_i + \delta[$ cumple lo requerido.

Si $|C| \neq 0$ entonces $a_i < b_i$, luego $m = \min_i (b_i - a_i) > 0$ y podemos considerar la función $g : [0, m/2] \rightarrow \mathbb{R}$ dada por

$$f(\delta) = \prod_{i=1}^n (b_i - \delta - (a_i + \delta)).$$

Tomando de nuevo un δ suficientemente pequeño, la celda $C'' = \prod_{i=1}^n [a_i + \delta, b_i - \delta]$ cumple lo requerido. ■

Definición 8.3 Una *figura elemental* en \mathbb{R}^n es una unión finita de celdas de \mathbb{R}^n . Llamaremos \mathcal{E}^n al conjunto de todas las figuras elementales en \mathbb{R}^n .

En general, es claro que toda figura elemental en \mathbb{R}^n es un subconjunto acotado de \mathbb{R}^n . Veamos algunas propiedades básicas:

Teorema 8.4 *Si A y B son figuras elementales, también lo son*

$$A \cup B, \quad A \cap B, \quad A \setminus B \quad \text{y} \quad a + A,$$

para todo $a \in \mathbb{R}^n$, donde $a + A = \{a + x \mid x \in A\}$ es la imagen de A por la traslación $x \mapsto a + x$.

DEMOSTRACIÓN: Sea $A = \bigcup_{k=1}^m C_k$, $B = \bigcup_{l=1}^r D_l$, donde los conjuntos C_k y D_l son celdas. Es trivial que $A \cup B$ es también una figura elemental. Por otra parte,

$$A \cap B = \bigcup_{k=1}^m \bigcup_{l=1}^r (C_k \cap D_l),$$

luego basta probar que la intersección de dos celdas C y D es una celda. Ahora bien, si

$$C = \prod_{i=1}^n I_i, \quad D = \prod_{i=1}^n J_i,$$

entonces $C \cap D = \prod_{i=1}^n (I_i \cap J_i)$, y es claro que una intersección de intervalos acotados es un intervalo acotado.

Como las figuras elementales son acotadas, podemos tomar un intervalo cerrado $[a, b]$ tal que $A \cup B \subset [a, b]^n$, y entonces

$$A \setminus B = A \cap ([a, b]^n \setminus B),$$

luego no perdemos generalidad si suponemos que $A = [a, b]^n$ y que $B \subset A$. Entonces, si llamamos $B_i = \{x \in A \mid x_i \notin B_i\}$, resulta que

$$A \setminus B = \bigcup_{i=1}^n B_i,$$

luego basta demostrar que cada B_i es una figura elemental, pero es claro que $[a, b] \setminus J_i = K_1 \cup K_2$, donde K_1 y K_2 son dos intervalos (que contendrán o no a sus extremos según si J_i los contiene o no), y entonces

$$B_i = \{x \in A \mid x_i \in K_1\} \cup \{x \in A \mid x_i \in K_2\},$$

y los dos conjuntos de la derecha son celdas.

Por último, $a + A = \bigcup_{k=1}^m (a + C_k)$, y es claro que el trasladado de una celda $C = \prod_{i=1}^n I_i$ es $a + C = \prod_{i=1}^n (a_i + I_i)$, donde a su vez $a_i + I_i$ es un intervalo, luego $a + C$ es una celda y $a + A$ es una figura elemental. ■

Conviene dar nombre a la situación dada por el teorema anterior:

Definición 8.5 Si X es un conjunto, un *anillo* de subconjuntos de X es una familia $\mathcal{A} \subset \mathcal{P}X$ que cumpla las propiedades siguientes:¹

- a) $\emptyset \in \mathcal{A}$.
- b) Si $A, B \in \mathcal{A}$, entonces $A \cup B, A \cap B, A \setminus B \in \mathcal{A}$.

Así, hemos probado que el conjunto de las figuras elementales es un anillo invariante por traslaciones, en el sentido de que cualquier trasladado de una figura elemental es una figura elemental.

Necesitamos un último resultado técnico:

¹El nombre se debe a que todo anillo de conjuntos es un anillo conmutativo en el sentido algebraico usual con las operaciones $A + B = (A \setminus B) \cup (B \setminus A)$ y $A \cdot B = A \cap B$, pero no vamos a necesitar este hecho.

Teorema 8.6 *Toda figura elemental puede expresarse de forma como unión de celdas disjuntas.*

DEMOSTRACIÓN: En primer lugar demostramos lo siguiente: Si I_1, \dots, I_k son intervalos acotados en \mathbb{R} , existen intervalos acotados J_1, \dots, J_l disjuntos dos a dos tales que $I_1 \cup \dots \cup I_k = J_1 \cup \dots \cup J_l$ y, para cada i, j , o bien $J_j \subset I_i$, o bien $J_j \cap I_i = \emptyset$.

En efecto, consideremos el conjunto de todos los extremos de los intervalos I_i y ordenemos sus elementos: $a_1 < a_2 < \dots < a_r$. Tomamos como intervalos J_j todos los de la forma $\{a_r\}$ o $]a_r, a_{r+1}[$ que estén contenidos en alguno de los I_i . Es claro entonces que $J_1 \cup \dots \cup J_l \subset I_1 \cup \dots \cup I_k$. Recíprocamente, dado $x \in I_1 \cup \dots \cup I_k$, existe un i tal que $x \in I_i$. Si $x = a_r$, para cierto r , entonces $\{x\} = \{a_r\} \subset I_i$, luego $\{x\}$ es un J_j y x está en la unión de todos ellos. La otra posibilidad es que $a_r < x < a_{r+1}$. En este caso los extremos de I_i deben cumplir $u \leq a_r < x < a_{r+1} \leq v$, por lo que $x \in]a_r, a_{r+1}[\subset I_i$, luego $]a_r, a_{r+1}[$ es un J_j y de nuevo x está en la unión.

Con esto tenemos probado que $J_1 \cup \dots \cup J_l = I_1 \cup \dots \cup I_k$, y es claro que los intervalos J_j son disjuntos dos a dos. Si $J_j \cap I_i \neq \emptyset$, tomamos $x \in J_j \cap I_i$ y distinguimos dos casos. O bien $J_j = \{a_r\}$, en cuyo caso $x = a_r$ y trivialmente $J_j \subset I_i$, o bien $J_j =]a_r, a_{r+1}[$, en cuyo caso I_i es un intervalo cuyos extremos cumplen $u \leq a_r < x < a_{r+1} \leq v$, por lo que $J_j \subset I_i$.

En segundo lugar demostramos que si C_1, \dots, C_k son celdas, existen otras celdas D_1, \dots, D_l disjuntas dos a dos de manera que $D_1 \cup \dots \cup D_l = C_1 \cup \dots \cup C_k$.

En efecto, si $C_v = \prod_{i=1}^n I_{iv}$, para cada $i = 1, \dots, n$ aplicamos el resultado anterior a la familia $\{I_{iv} \mid v = 1, \dots, k\}$, lo que nos da una familia de intervalos $\{J_{ju} \mid u = 1, \dots, l_i\}$ disjuntos dos a dos de modo que $\bigcup_{v=1}^k I_{iv} = \bigcup_{u=1}^{l_i} J_{iu}$ y cada I_{iv} está contenido en cada I_{iv} o bien es disjunto de él.

Llamamos D_r a las celdas de la forma $\prod_{i=1}^n J_{iu_i}$ contenidas en algún C_v . Claramente son disjuntas dos a dos, pues si dos de ellas $\prod_{i=1}^n J_{iu_i}$ y $\prod_{i=1}^n J_{iv_i}$ son distintas, entonces $J_{iu_i} \neq J_{iv_i}$ para algún i , luego $J_{iu_i} \cap J_{iv_i} = \emptyset$, luego también las celdas son disjuntas.

También es obvio que $D_1 \cup \dots \cup D_l \subset C_1 \cup \dots \cup C_k$ y si $x \in C_1 \cup \dots \cup C_k$, entonces $x \in C_v$, para cierto v , luego $x_i \in J_{iu_i} \subset I_{iv}$, para cierto u_i , luego $x \in \prod_{i=1}^n J_{iu_i} \subset C_v$, luego la celda es una de las D_u y $x \in D_1 \cup \dots \cup D_l$.

Lo que acabamos de probar equivale a que si $A = C_1 \cup \dots \cup C_k$ es una figura elemental, entonces A se expresa también como unión de celdas disjuntas dos a dos. ■

Pasamos ya al problema de asignar un “volumen” a cada figura elemental. La definición siguiente incorpora algunas de las condiciones que, según hemos discutido, debe cumplir tal asignación:

Definición 8.7 Una *medida finitamente aditiva* en un anillo \mathcal{A} en un conjunto X es una aplicación $\mu : \mathcal{A} \rightarrow [0, +\infty[$ que cumpla las propiedades siguientes:

- a) $\mu(\emptyset) = 0$.
- b) Si $A, B \in \mathcal{A}$ son disjuntos, entonces $\mu(A \cup B) = \mu(A) + \mu(B)$.

Toda medida finitamente aditiva cumple estas propiedades adicionales:

- a) Si $A \subset B$ son elementos de \mathcal{A} , entonces $\mu(A) \leq \mu(B)$.
En efecto, $\mu(B) = \mu(A \cup (B \setminus A)) = \mu(A) + \mu(B \setminus A) \geq \mu(A)$.
- b) Si $A, B \in \mathcal{A}$, entonces $\mu(A \cup B) \leq \mu(A) + \mu(B)$.
En efecto, $\mu(A \cup B) = \mu(A \cup (B \setminus A)) = \mu(A) + \mu(B \setminus A) \leq \mu(A) + \mu(B)$.
- c) Si $A, B \in \mathcal{A}$ tienen medida nula, entonces $\mu(A \cup B) = 0$.

Es un caso particular del apartado anterior.

Las propiedades de las medidas finitamente aditivas son parte de las que debe cumplir cualquier definición de “área” o “volumen” que sea acorde con el concepto intuitivo que queremos precisar. Una tercera es que sea invariante por isometrías, en particular por traslaciones. El teorema siguiente muestra que estas propiedades (junto con que tomamos la celda unitaria como unidad de medida) determinan completamente la medida de cualquier figura elemental:

Teorema 8.8 *Existe una única medida finitamente aditiva $m : \mathcal{E}^n \rightarrow [0, +\infty[$ que cumple las propiedades siguientes:*

- a) $m([0, 1]^n) = 1$,
- b) Si $a \in \mathbb{R}^n$ y $A \in \mathcal{E}^n$, entonces $m(a + A) = m(A)$.

DEMOSTRACIÓN: Veamos primero la unicidad. Para ello probaremos algunos hechos que debe cumplir una medida en las condiciones del enunciado:

- a) Si $C = \prod_{i=1}^n I_i$ cumple que algún $I_i = \{x\}$, entonces $m(C) = 0$.

En efecto, podemos tomar un intervalo $[a, b]$ tal que $C \subset [a, b]^n$ y $a < b$. Sea $e = (0, \dots, 1, \dots, 0) \in \mathbb{R}^n$, con el 1 en la posición i -ésima. Para cada $u \in [a, b]$ sea $C_u = (u - x)e + C$, es decir, la celda que coincide con C salvo que la componente i -ésima de sus puntos es u en lugar de x . Es claro que los conjuntos C_u son disjuntos dos a dos y, como son trasladados de C , todos tienen la misma medida. Supongamos que $m(C) = \epsilon > 0$. Entonces podemos tomar un $N > 0$ tal que $m([a, b]^n) < N\epsilon$, pero podemos tomar N puntos distintos u_1, \dots, u_N en $[a, b]$, con lo que

$$C^* = C_{u_1} \cup \dots \cup C_{u_N} \subset [a, b]^n$$

y $m(C^*) = m(C_{u_1}) + \dots + m(C_{u_N}) = N\epsilon > m([a, b]^n)$, contradicción.

b) $m(]0, 1[^n) = 1$.

En efecto, sea $C_{ij} = \{x \in [0, 1]^n \mid x_i = j\}$, que claramente es una celda de medida 0 (por el apartado anterior), y

$$[0, 1]^n =]0, 1[^n \cup \bigcup_{i=1}^n (C_{i0} \cup C_{i1}).$$

Como la unión de la derecha es nula y la primera unión es disjunta, vemos que $1 = m([0, 1]^n) = m(]0, 1[^n)$.

c) Si $C = \prod_{i=1}^n I_i$ es una celda y cada I_i es un intervalo de extremos 0 y 1, entonces $m(C) = 1$.

Basta tener en cuenta que $]0, 1[^n \subset C \subset [0, 1]^n$, y que ambas celdas tienen medida 1.

d) Si $C = \prod_{i=1}^n I_i$ es una celda y cada I_i es un intervalo de extremos a_i y $a_i + 1$, entonces $m(C) = 1$.

Porque C es un trasladado de una celda como la del apartado anterior.

e) Si $C = \prod_{i=1}^n I_i$ es una celda y cada I_i es un intervalo de extremos enteros, entonces $m(C) = |C|$.

Sea $k_i \in \mathbb{N}$ la diferencia entre los extremos de I_i . Es claro que cada intervalo I_i se descompone en unión disjunta de k_i intervalos disjuntos de longitud 1, luego C se descompone en $\prod_{i=1}^n k_i$ celdas disjuntas en las condiciones del apartado anterior, luego $m(C) = \prod_{i=1}^n k_i = |C|$.

f) Si $C = \prod_{i=1}^n I_i$ es una celda y cada I_i es un intervalo de extremos racionales, entonces $m(C) = |C|$.

En efecto, si los extremos de I_i son $a_i = a'_i/d$, $b_i = b'_i/d$, llamamos I'_i al intervalo $[a'_i, b'_i]$ o bien $]a'_i, b'_i]$, $[a'_i, b'_i[$, $]a'_i, b'_i[$ de modo que sea del mismo tipo que I_i . Entonces por el apartado anterior

$$m(C') = \prod_{i=1}^n (b'_i - a'_i) = d^n \prod_{i=1}^n (b_i - a_i) = d^n |C|.$$

Por otra parte, I'_i se descompone en unión de d intervalos disjuntos de longitud $b_i - a_i$, todos los cuales son trasladados de I_i , luego C' se descompone en d^n celdas disjuntas trasladadas de C , luego $m(C') = d^n m(C)$ y llegamos a la igualdad del enunciado.

g) Si $C = \prod_{i=1}^n I_i$ es una celda, entonces $m(C) = |C|$.

En efecto, si algún I_i es vacío o se reduce a un punto, ya sabemos que la medida es 0, luego se cumple la igualdad. En caso contrario, dado $\epsilon > 0$, usando que la función $\prod_{i=1}^n (y_i - x_i)$ es continua en \mathbb{R}^{2n} , es fácil construir celdas $C' \subset C \subset C''$ cuyos intervalos tengan extremos racionales y de modo que $|C| - |C'| < \epsilon$ y $|C''| - |C| < \epsilon$. Así, por el apartado anterior,

$$|C| - \epsilon < |C'| = m(C') \leq m(C) \leq m(C'') = |C''| < |C| + \epsilon,$$

luego $|m(C) - |C|| < \epsilon$ para todo $\epsilon > 0$, luego $m(C) = |C|$.

Ahora definimos $m : \mathcal{E}^n \rightarrow [0, +\infty[$ de modo que si A es una figura elemental que se expresa como unión de celdas disjuntas $A = \bigcup_{k=1}^m C_k$, entonces

$$m(A) = \sum_{k=1}^m |C_k|.$$

Tenemos que probar que $m(A)$ es independiente de la descomposición considerada. Admitiendo esto, es claro que m es una medida finitamente aditiva invariante por traslaciones, y por el apartado g) precedente m es la única medida que puede cumplir estas propiedades.

Empezamos observando que si I es un intervalo acotado de extremos $a \leq b$, entonces

$$b - a = \lim_m \frac{1}{m} |I \cap \frac{1}{m} \mathbb{Z}|.$$

En efecto, si m es suficientemente grande,

$$\begin{aligned} |I \cap \frac{1}{m} \mathbb{Z}| &= |\{z \in \mathbb{Z} \mid \frac{z}{m} \in I\}| \leq |\{z \in \mathbb{Z} \mid ma \leq z \leq mb\}| \\ &\leq |\{z \in \mathbb{Z} \mid E[ma] \leq z \leq E[mb] + 1\}| \leq mb + 1 - (ma - 1) + 1 \\ &= m(b - a) + 3. \end{aligned}$$

Igualmente

$$\begin{aligned} |I \cap \frac{1}{m} \mathbb{Z}| &= |\{z \in \mathbb{Z} \mid \frac{z}{m} \in I\}| \geq |\{z \in \mathbb{Z} \mid ma < z < mb\}| \\ &\geq |\{z \in \mathbb{Z} \mid E[ma] + 1 < z < E[mb]\}| \geq mb - 1 - (ma + 1) + 1 \\ &= m(b - a) - 1. \end{aligned}$$

Por lo tanto

$$b - a - \frac{1}{m} \leq \frac{1}{m} |I \cap \frac{1}{m} \mathbb{Z}| \leq b - a + \frac{3}{m},$$

y es claro que el límite es $b - a$.

En segundo lugar probamos que si $C = \prod_{i=1}^n I_i$ es una celda en \mathbb{R}^n , entonces

$$m(C) = \lim_m \frac{1}{m^n} |C \cap \frac{1}{m} \mathbb{Z}^n|.$$

En efecto,

$$|C \cap \frac{1}{m} \mathbb{Z}^n| = |\{x \in \mathbb{R}^n \mid x_i \in I_i \cap \frac{1}{m} \mathbb{Z}, i = 1, \dots, n\}| = \prod_{i=1}^n |I_i \cap \frac{1}{m} \mathbb{Z}|,$$

y por el apartado anterior el límite toma el valor indicado.

Como consecuencia inmediata, si $A \in \mathcal{E}^n$ se expresa como unión de celdas disjuntas $A = \bigcup_{k=1}^l C_k$, entonces

$$\sum_{k=1}^l m(C_k) = \lim_m \frac{1}{m^n} |A \cap \frac{1}{m} \mathbb{Z}^n|.$$

Como la expresión de la derecha no depende de la descomposición de A , concluimos que la definición de m es correcta. Equivalentemente, podríamos haber definido

$$m(A) = \lim_m \frac{1}{m^n} |A \cap \frac{1}{m} \mathbb{Z}^n|.$$

■

Así pues, la medida que hemos definido sobre las figuras elementales es la única definición de “medida” que cumple las propiedades que intuitivamente tiene que tener un área o un volumen. Sin embargo, las figuras elementales son una familia de figuras demasiado pobre. A continuación extendemos la medida a una familia mayor:

Definición 8.9 Sea $A \subset \mathbb{R}^n$ un conjunto acotado. Definimos la *medida exterior de Jordan* $J^*(A)$ y la *medida interior de Jordan* $J_*(A)$ como

$$J^*(A) = \inf\{m(B) \mid B \in \mathcal{E}^n, A \subset B\}, \quad J_*(A) = \sup\{m(B) \mid B \in \mathcal{E}^n, B \subset A\}.$$

Diremos que el conjunto A es *medible Jordan* si $J^*(A) = J_*(A)$, y en tal caso a este valor común lo llamaremos *medida de Jordan* de A , y lo representaremos por $m(A) = J^*(A) = J_*(A)$. Llamaremos \mathcal{J}^n al conjunto de todos los subconjuntos de \mathbb{R}^n medibles Jordan.

Notemos que todo conjunto acotado $A \subset \mathbb{R}^n$ cumple que $J_*(A) \leq J^*(A)$.

Más aún, si $B \subset A$ es una figura elemental, cualquier criterio que asigne a A una medida respetando los principios básicos que hemos convenido que debe cumplir la medida en el sentido geométrico intuitivo hará que $m(B) \leq m(A)$, y por definición de supremo será, de hecho, $J_*(A) \leq m(A)$. Igualmente, si $A \subset B$, tiene que cumplirse $m(A) \leq J^*(B)$, luego en general, las propiedades que intuitivamente debe cumplir la medida geométrica exigen que cualquier conjunto cumpla

$$J_*(A) \leq m(A) \leq J^*(A).$$

Por eso, si se da la igualdad $J_*(A) = J^*(A)$, podemos afirmar que dicho valor es la única asignación posible de una medida al conjunto A . Si se da la desigualdad $J_*(A) < J^*(A)$ nos quedamos con la duda de si podremos encontrar algún otro criterio más refinado que justifique una asignación de una medida a A que no pueda tacharse de arbitraria. Más adelante volveremos sobre ello.

El teorema siguiente nos proporciona una caracterización práctica de la medibilidad de Jordan:

Teorema 8.10 *Un subconjunto acotado $A \subset \mathbb{R}^n$ es medible Jordan si y sólo si para todo $\epsilon > 0$ existen figuras elementales A_0 y A_1 tales que $A_0 \subset A \subset A_1$ y $m(A_1 \setminus A_0) < \epsilon$.*

DEMOSTRACIÓN: Si A es medible Jordan, basta aplicar las definiciones de supremo e ínfimo para encontrar figuras elementales $A_0 \subset A \subset A_1$ tales que $m(A) - m(A_0) < \epsilon/2$, $m(A_1) - m(A) < \epsilon/2$, y entonces

$$\begin{aligned} m(A_1 \setminus A_0) &= m(A_1) - m(A_0) < m(A) + \epsilon/2 - m(A_0) \\ &< m(A_0) + \epsilon/2 + \epsilon/2 - m(A_0) = \epsilon. \end{aligned}$$

Recíprocamente, la hipótesis del teorema nos da que

$$0 \leq J^*(A) - J_*(A) \leq m(A_1) - m(A_0) = m(A_1 \setminus A_0) < \epsilon$$

para todo $\epsilon > 0$, luego $J^*(A) = J_*(A)$. ■

En ocasiones resulta útil la precisión siguiente:

Teorema 8.11 *Si $A \subset \mathbb{R}^n$ es medible Jordan y $\epsilon > 0$, existe una figura elemental compacta K y otra abierta U tales que $K \subset A \subset U$ tales que $K \subset A \subset U$ y $m(U \setminus K) < \epsilon$.*

DEMOSTRACIÓN: Tomamos $A_0 \subset A \subset A_1$ según el teorema anterior, de modo que $m(A_1 \setminus A_0) < \epsilon/3$. Pongamos que $A_0 = \bigcup_{k=1}^l C_k$, para ciertas celdas disjuntas C_k . Si eliminamos todas las de medida nula no alteramos la medida de A_0 , luego no perdemos generalidad si suponemos que todos los C_k tienen medida no nula. Por el teorema 8.2 podemos tomar celdas compactas $C'_k \subset C_k$ de modo que $|C_k| - |C'_k| < \epsilon/3l$.

Tomamos $K = \bigcup_{k=1}^l C'_k$, de modo que $K \subset A_0 \subset A$ y $m(A_0 \setminus K) < \epsilon/3$. Razonando igualmente con las celdas de A_1 , podemos formar una figura elemental abierta U tal que $A_1 \subset U$ y $m(U \setminus A_1) < \epsilon/3$. Así

$$\begin{aligned} m(U \setminus K) &= m(U) - m(K) = m(U) - m(A_1) + m(A_1) - m(A_0) + m(A_0) - m(K) \\ &= m(U \setminus A_1) + m(A_1 \setminus A_0) + m(A_0 \setminus K) < \epsilon. \end{aligned} \quad \blacksquare$$

Pasamos a probar los hechos relevantes sobre la medida de Jordan:

Teorema 8.12 \mathcal{J}^n es un anillo en \mathbb{R}^n que contiene a \mathcal{E}^n y la medida de Jordan $m : \mathcal{J}^n \rightarrow [0, +\infty[$ es la única medida finitamente aditiva que extiende a la medida de figuras elementales.

DEMOSTRACIÓN: Sean A y B dos conjuntos medibles Jordan. Dado $\epsilon > 0$, sean $A_0 \subset A \subset A_1$, $B_0 \subset B \subset B_1$ figuras elementales tales que $m(A_1 \setminus A_0) < \epsilon/2$, $m(B_1 \setminus B_0) < \epsilon/2$. Entonces $C_0 = A_0 \cup B_0 \subset A \cup B \subset A_1 \cup B_1 = C_1$ y

$$m(C_1 \setminus C_0) \leq m((A_1 \setminus A_0) \cup (B_1 \setminus B_0)) \leq m(A_1 \setminus A_0) + m(B_1 \setminus B_0) < \epsilon.$$

Por lo tanto, $A \cup B$ es medible Jordan.

Similarmente, $C'_0 = A_0 \cap B_0 \subset A \cap B \subset A_1 \cap B_1 = C'_1$ y

$$m(C'_1 \setminus C'_0) \leq m((A_1 \setminus A_0) \cup (B_1 \setminus B_0)) < \epsilon.$$

Por último tomamos $C''_0 = A_0 \setminus B_1 \subset A \setminus B \subset A_1 \setminus B_0 = C''_1$ y también

$$m(C''_1 \setminus C''_0) \leq m((A_1 \setminus A_0) \cup (B_1 \setminus B_0)) < \epsilon.$$

Así, $A \setminus B$ es medible Jordan y queda probado que \mathcal{J}^n es un anillo de conjuntos.

Es inmediato que si A es una figura elemental $J_*(A) = J^*(A) = m(A)$, luego $A \in \mathcal{J}^n$ y la medida de Jordan extiende a la que ya teníamos definida sobre las figuras elementales. En particular la medida de la celda unitaria es 1.

Si $\mu : \mathcal{J}^n \rightarrow [0, +\infty[$ es cualquier medida finitamente aditiva que cumpla estas propiedades, entonces a partir de las definiciones de medida interior y exterior es trivial que $m(A) = J_*(A) \leq \mu(A) \leq J^*(A) = m(A)$, luego $\mu = m$. ■

Teorema 8.13 La medida de Jordan es también la única medida finitamente aditiva en \mathcal{J}^n invariante por traslaciones que cumple $m([0, 1]^n) = 1$.

DEMOSTRACIÓN: Dado $a \in \mathbb{R}^n$ y un $A \in \mathcal{J}^n$, es inmediato comprobar que $J_*(a + A) = J_*(A) = J^*(A) = J^*(a + A)$, luego $a + A \in \mathcal{J}^n$ y $m(a + A) = m(A)$.

Si una medida $\mu : \mathcal{J}^n \rightarrow [0, +\infty[$ cumple esto mismo, entonces su restricción a \mathcal{E}^n es una medida finitamente aditiva invariante por traslaciones, luego tiene que ser la medida de las figuras elementales, y por el teorema anterior $\mu = m$. ■

La familia de los conjuntos medibles Jordan es mucho más amplia que la de las figuras elementales. Lo deduciremos del teorema siguiente:

Teorema 8.14 Un conjunto acotado $A \subset \mathbb{R}^n$ es medible Jordan si y sólo si lo es su frontera y $m(\partial A) = 0$.

DEMOSTRACIÓN: Supongamos que $m(\partial A) = 0$. Vamos a distinguir dos casos:

1) Si $J_*(A) = 0$, entonces $\overset{\circ}{A} = \emptyset$, pues si A contuviera un abierto, contendría una celda de medida no nula, y su medida interior no sería nula. Pero entonces $A \subset \partial A$, luego $J^*(A) \leq J^*(\partial A) = 0$, luego A es medible Jordan y $m(A) = 0$.

2) Si $J_*(A) > 0$ tomamos $\epsilon > 0$. Existe una figura elemental B tal que $\partial A \subset B$ y $m(B) < \epsilon$. Concretamente, $B = \bigcup_{k=1}^m C_k$, donde las celdas C_k son disjuntas dos a dos y $\sum_{k=1}^m m(C_k) < \epsilon$. Por el teorema 8.2, podemos agrandar ligeramente cada celda y suponer que los C_k son abiertos, aunque dejen de ser disjuntos, pero de modo que sigan cumpliendo que sus medidas sumen menos que ϵ .

Para cada $x \in \overset{\circ}{A}$ existe una celda abierta tal que $x \in C_x \subset \overset{\circ}{A}$. Por lo tanto,

$$\overline{A} \subset \bigcup_{k=1}^m C_k \cup \bigcup_{x \in \overset{\circ}{A}} C_x.$$

Como A es acotado, \overline{A} es compacto, luego existe un conjunto finito de puntos $x_1, \dots, x_l \in \overset{\circ}{A}$ de modo que

$$\overline{A} \subset \bigcup_{k=1}^m C_k \cup \bigcup_{j=1}^l C_{x_j}.$$

Sea $B = \bigcup_{j=1}^l C_{x_j} \subset A$. Entonces $J^*(A) \leq \sum_{k=1}^m m(C_k) + m(B) < \epsilon + J_*(A)$, para todo $\epsilon > 0$, luego $J^*(A) = J_*(A)$.

Ahora supongamos que A es medible Jordan pero que ∂A no tiene medida de Jordan nula. Eso es equivalente a que $J^*(A) > 0$, luego existe un $\epsilon > 0$ tal que toda figura elemental B que contenga a ∂A cumple $m(B) \geq \epsilon$. Como A es medible, tienen que existir figuras elementales $A_0 \subset A \subset A_1$ tales que $m(A_1 \setminus A_0) < \epsilon$. Expresemos $A_0 = \bigcup_{k=1}^m C_k$, $A_1 = \bigcup_{r=1}^l C'_r$ como uniones de celdas disjuntas, de modo que

$$\sum_{r=1}^l m(C'_r) < \sum_{k=1}^m m(C_k) + \epsilon.$$

Estas sumas no se alteran si cambiamos cada C'_r por su clausura y cada C_k por su interior (aunque los C'_r dejen de ser disjuntos). Así $A_0 \subset \overset{\circ}{A} \subset \overline{A} \subset A_1$, luego $\partial A = \overline{A} \setminus \overset{\circ}{A} \subset B = A_1 \setminus A_0$, con

$$m(B) = m(A_1) - m(A_0) \leq \sum_{r=1}^l m(C'_r) - \sum_{k=1}^m m(C_k) < \epsilon,$$

contradicción. ■

De aquí se sigue, en particular, que un conjunto A es medible Jordan si y sólo si lo es su clausura o su interior, y en tal caso los tres tienen la misma medida.

El teorema siguiente nos permitirá extraer muchas más consecuencias. Observemos que hasta ahora los únicos segmentos que sabemos que son medibles son los que son paralelos a algún eje de coordenadas:

Teorema 8.15 Si $n \geq 2$, todos los segmentos son medibles Jordan y tienen medida nula.

DEMOSTRACIÓN: El segmento de extremos a y b es el conjunto

$$S = \{(1 - \lambda)a + \lambda b \mid 0 \leq \lambda \leq 1\}.$$

Como la medida de Jordan es invariante por traslaciones, no perdemos generalidad si suponemos que $a = 0$. Sea

$$C_{m,j} = \prod_{i=1}^n \left[\frac{j-1}{m} b_i, \frac{j}{m} b_i \right], \quad j = 1, \dots, m,$$

donde hay que entender que si $b_i < 0$ entonces los extremos del intervalo van en el orden contrario. Es claro entonces que $m(C_{m,j}) = |b_1 \cdots b_n|/m^n$, luego la figura elemental

$$A_m = \bigcup_{j=1}^m C_{m,j}$$

cumple $m(A_m) \leq |b_1 \cdots b_n|/m^{n-1}$. Por otra parte es claro que $S \subset A$, pues dado $0 \leq \lambda \leq 1$, existe un j tal que $(j-1)/m \leq \lambda \leq j/m$, y entonces $\lambda b \in C_{m,j}$. Por consiguiente $J^*(S) \leq m(A_m) \leq |b_1 \cdots b_n|/m^{n-1}$ para todo natural m , luego $J^*(S) = 0$, y esto implica que S es medible Jordan y $m(S) = 0$. ■

Por consiguiente, en \mathbb{R}^2 , todas las figuras limitadas por segmentos son medibles Jordan. Otro resultado útil de medibilidad es el siguiente:

Teorema 8.16 Si $A \in \mathcal{J}^m$ y $B \in \mathcal{J}^n$, entonces $A \times B \in \mathcal{J}^{m+n}$ y se cumple que $m(A \times B) = m(A)m(B)$.

DEMOSTRACIÓN: El resultado es trivialmente cierto si A y B son celdas. Si $A = \bigcup_{k=1}^r C_k$, $B = \bigcup_{l=1}^s C'_l$ son uniones de celdas disjuntas, entonces

$$A \times B = \bigcup_{k=1}^r \bigcup_{l=1}^s (C_k \times C'_l)$$

es también una unión de celdas disjuntas, luego

$$m(A \times B) = \sum_{k=1}^r \sum_{l=1}^s m(C_k)m(C'_l) = m(A)m(B).$$

En el caso general tomamos figuras elementales $A_0 \subset A \subset A_1$, $B_0 \subset B \subset B_1$, de modo que $A_0 \times B_0 \subset A \times B \subset A_1 \times B_1$ y

$$\begin{aligned} m((A_1 \times B_1) \setminus (A_0 \times B_0)) &= m(A_1)m(B_1) - m(A_0)m(B_0) \\ &= m(A_1)m(B_1) - m(A_1)m(B_0) + m(A_1)m(B_0) - m(A_0)m(B_0) \\ &= m(A_1)(m(B_1) - m(B_0)) + (m(A_1) - m(A_0))m(B_0). \end{aligned}$$

Por lo tanto, dado $\epsilon > 0$, si exigimos que

$$m(A_1 \setminus A_0) < \frac{\epsilon}{2m(B)}, \quad m(A_1) < m(A) + 1, \quad m(B_1 \setminus B_0) < \frac{\epsilon}{2(m(A) + 1)},$$

concluimos que $m((A_1 \times B_1) \setminus (A_0 \times B_0)) < \epsilon$, luego $A \times B$ es medible Jordan y

$$m(A_0)m(B_0) \leq m(A \times B) \leq m(A_1)m(B_1),$$

y también $m(A_0)m(B_0) \leq m(A)m(B) \leq m(A_1)m(B_1)$, luego

$$|m(A \times B) - m(A)m(B)| \leq m(A_1)m(B_1) - m(A_0)m(B_0) < \epsilon,$$

luego $m(A \times B) = m(A)m(B)$. ■

En la introducción a este capítulo hemos observado que una definición de medida que se ajuste al concepto intuitivo de área o de volumen debe ser invariante por isometrías, pero hasta ahora sólo hemos considerado la invarianza por traslaciones. Teniendo en cuenta que la isometrías tienen determinante ± 1 , el teorema siguiente muestra que en realidad la medida de Jordan es invariante por isometrías:

Teorema 8.17 *Si $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ una aplicación lineal y $A \in \mathcal{J}^n$, entonces $f[A] \in \mathcal{J}^n$ y $m(f[A]) = |\det f| m(A)$.*

DEMOSTRACIÓN: Consideramos primero el caso en que f es biyectiva, con lo que $\det f \neq 0$. Aplicamos el teorema [Al 8.23], según el cual basta probar el teorema cuando f es una de las aplicaciones siguientes:

- a) La aplicación dada por $f(x_1, \dots, x_n) = f(x_1, \dots, x_j, \dots, x_i, \dots, x_n)$, donde x_j aparece en el lugar i -ésimo y viceversa.
- b) La aplicación $f(x_1, \dots, x_n) = (x_1, x_1 + x_2, x_3, \dots, x_n)$.
- c) La aplicación $f(x_1, \dots, x_n) = (ax_1, x_2, \dots, x_n)$, con $a \neq 0$.

En el caso a), el determinante es -1 . Es claro que f transforma celdas en celdas de la misma medida, luego figuras elementales en figuras elementales de la misma medida, de donde se sigue fácilmente que si $A \in \mathcal{J}^n$ entonces $f[A] \in \mathcal{J}^n$ y $m(f[A]) = m(A)$.

El caso c) es similar. Sigue siendo cierto que f transforma celdas en celdas, pero ahora $m(f[C]) = |a|m(C)$, de donde se sigue que lo mismo vale para figuras elementales y para conjuntos medibles Jordan cualesquiera, y eso es lo que había que probar, pues $\det f = a$.

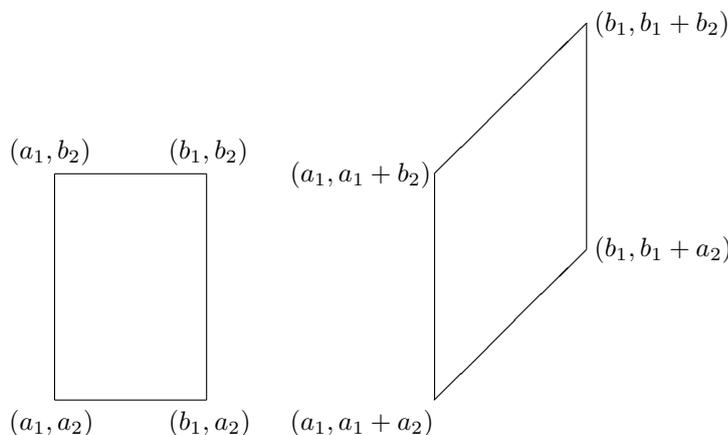
El caso b) es el más delicado, porque f ya no transforma celdas en celdas. En primer lugar veamos que podemos suponer $n = 2$. En efecto, sea $f' : \mathbb{R}^2 \rightarrow \mathbb{R}^2$ la aplicación dada por $f'(x_1, x_2) = (x_1, x_1 + x_2)$ y supongamos que el teorema es cierto para f' . Toda celda en \mathbb{R}^n (con $n > 2$) es de la forma $C_1 \times C_2$, donde

C_1 es una celda en \mathbb{R}^2 y C_2 es una celda en \mathbb{R}^{n-2} , y $f[C_1 \times C_2] = f'[C_1] \times C_2$, que es medible Jordan por el teorema anterior y

$$m(f[C_1 \times C_2]) = m(f'[C_1])m(C_2) = m(C_1)m(C_2) = m(C_1 \times C_2).$$

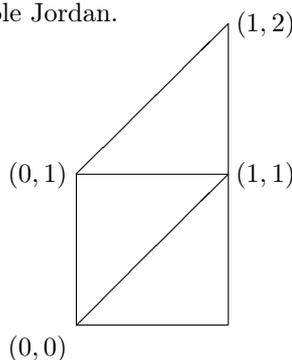
En definitiva, si C es cualquier celda, se cumple que $f[C] \in \mathcal{J}^n$ y además $m(f[C]) = m(C)$. De aquí se sigue inmediatamente que lo mismo vale para figuras elementales y a su vez para conjuntos medibles Jordan arbitrarios.

Así pues, suponemos que $n = 2$ y vamos a probar que las imágenes de las celdas son medibles Jordan con la misma medida. Partimos de una celda $C = [a_1, b_1] \times [a_2, b_2]$. Una comprobación rutinaria muestra que $f[C]$ es el paralelogramo de vértices $(a_1, a_1 + a_2), (b_1, b_1 + a_2), (b_1, b_1 + b_2), (a_1, a_1 + b_2)$.



Podemos afirmar que $f[C]$ es medible Jordan porque su frontera está formada por segmentos, luego es nula. De aquí se sigue, como en los casos anteriores, que la imagen de cualquier figura elemental es medible Jordan, y a su vez que la imagen de todo conjunto medible Jordan es medible Jordan.

Ahora consideramos concretamente el caso del cuadrado unitario $C = [0, 1]^2$, cuya imagen es el paralelogramo que muestra la figura. Vemos que podemos descomponerlo en unión de dos triángulos, y que el triángulo superior puede trasladarse para formar de nuevo el cuadrado unitario. Por consiguiente, $m(f[C]) = 1$.



Ahora, la función $\mu : \mathcal{J}^2 \rightarrow [0, +\infty[$ dada por $\mu(A) = m(f[A])$ es una medida finitamente aditiva, y es invariante por traslaciones, ya que $\mu(a + A) = m(f(a) + f[A]) = m(f[A]) = \mu(A)$.

Como además $\mu([0, 1]^2) = 1$, concluimos que $\mu = m$, luego $m(f[A]) = m(A)$ para todo conjunto medible Jordan.

Falta considerar el caso en que $\det f = 0$, y entonces hay que probar que todo conjunto medible Jordan cumple $m(f[A]) = 0$. Para ello nos basamos en que $f[\mathbb{R}^n] \subset H$, para cierto hiperplano H , y existe una aplicación lineal biyectiva $g: \mathbb{R}^n \rightarrow \mathbb{R}^n$ tal que $g[H]$ es el hiperplano de ecuación $x_1 = 0$. Así, $g[f[A]]$ es un subconjunto acotado de este hiperplano, luego está contenido en una celda de la forma $C = \{x_0\} \times C'$, luego $J^*(g[f[A]]) \leq J^*(C) = m(C) = 0$, luego $g[f[A]]$ es medible y tiene medida nula. Por la parte ya probada, aplicada a g^{-1} , concluimos que $f[A]$ es medible con medida nula. ■

Combinando el teorema anterior con la invarianza por traslaciones de la medida de Jordan, concluimos que ésta es invariante por isometrías afines (no necesariamente lineales) y, más en general, que el teorema anterior vale para biyecciones afines cualesquiera.

Otro caso interesante es el efecto de las homotecias. Como la homotecia de razón r tiene determinante r^n , se cumple la relación

$$m(rA) = |r|^n m(A).$$

Por el mismo argumento empleado en la parte final de la prueba del teorema anterior, toda figura acotada en \mathbb{R}^3 limitada por un número finito de caras planas es medible Jordan.

A partir de aquí es posible calcular el área de muchas figuras planas mediante argumentos sintéticos (véase [G 7.7]).

La familia de los conjuntos medibles Jordan es suficientemente amplia para trabajar en el contexto de la geometría sintética, pero no lo suficiente para trabajar con técnicas analíticas y topológicas. Por ejemplo, puede probarse que existen subconjuntos abiertos acotados en \mathbb{R}^n y, por consiguiente, también subconjuntos cerrados acotados, que no son medibles Jordan. Esto hace que al trabajar con funciones continuas arbitrarias sea fácil encontrarse con conjuntos cuya medibilidad no podamos asegurar.

Consideremos por ejemplo el conjunto $A = [0, 1]^2 \cap \mathbb{Q}^2$, es decir, el conjunto de puntos del cuadrado unitario con coordenadas racionales. Es fácil ver que $\overset{\circ}{A} = \emptyset$ y $\overline{A} = [0, 1]^2$, lo cual implica que A no es medible Jordan, ya que hemos visto que el interior y la clausura de un conjunto medible tienen la misma medida. No es difícil ver que $J_*(A) = 0$ y $J^*(A) = 1$. Así pues, A no tiene un área asignada en el sentido de la medida de Jordan, pero podemos preguntarnos si existe algún criterio adicional por el que asignarle un área. Los criterios que hemos empleado hasta ahora sólo exigen que si asignamos un área a A , ésta sea un número entre 0 y 1.

Conviene destacar que la pregunta ahora no es cuál es el área de A , sino más bien, qué área conviene que tenga. El área de un triángulo no es una cuestión de conveniencia: si le asignamos a un triángulo un área distinta de su base por su altura dividida entre 2, podemos afirmar que estamos considerando una noción de área distinta de la intuitiva, pues la tela necesaria para construir alfombras triangulares es la que es, y no la que nosotros queramos que sea. En cambio, no hay ninguna razón concluyente por la que pudiéramos decir que sería

“antiintuitivo” asignarle a A , por ejemplo, un área 0 (como a su interior) o un área 1 (como a su clausura), sino que el problema es qué criterios dan lugar a una asignación coherente de áreas y, si es que hay varios posibles, cuál es más conveniente en la práctica. Desde un punto de vista intuitivo, no hay nada que decir, pues no se pueden hacer alfombras de tela con forma de conjunto A . En este punto es relevante el teorema siguiente:

Teorema 8.18 *Sea $\{A_k\}_{k=0}^{\infty}$ una sucesión de conjuntos medibles Jordan en \mathbb{R}^n disjuntos dos a dos y supongamos que $A = \bigcup_{k=0}^{\infty} A_k$ es también medible Jordan. Entonces $m(A) = \sum_{k=0}^{\infty} m(A_k)$.*

DEMOSTRACIÓN: Por una parte,

$$\sum_{k=1}^m m(A_k) = m\left(\bigcup_{k=1}^l A_k\right) \leq m(A),$$

de modo que la serie es una sucesión monótona creciente acotada superiormente, luego converge y

$$\sum_{k=1}^{\infty} m(A_k) \leq m(A).$$

Sea $\epsilon > 0$. Por el teorema 8.11 podemos tomar una figura elemental $K \subset A$ compacta tal que $m(A \setminus K) < \epsilon/2$. Similarmente, para cada k existe una figura elemental abierta U_k tal que $A_k \subset U_k$ y $m(U_k \setminus A_k) < \epsilon/2^{k+2}$. Así $K \subset \bigcup_{k=0}^{\infty} U_k$,

luego por compacidad existe un l tal que $K \subset \bigcup_{k=0}^l U_k$. Por consiguiente

$$\begin{aligned} m(A) &= m(A \setminus K) + m(K) \leq \frac{\epsilon}{2} + \sum_{k=0}^l m(U_k) \leq \frac{\epsilon}{2} + \sum_{k=0}^l m(A_k) + \sum_{k=0}^l \frac{\epsilon}{2^{k+2}} \\ &\leq \frac{\epsilon}{2} + \sum_{k=0}^{\infty} m(A_k) + \sum_{k=0}^{\infty} \frac{\epsilon}{2^{k+2}} = \sum_{k=0}^{\infty} m(A_k) + \epsilon. \end{aligned}$$

Como esto vale para todo $\epsilon > 0$, se da la igualdad del enunciado. ■

Así pues, siempre que una unión disjunta de conjuntos medibles es medible, podemos afirmar que la medida es la suma infinita de las medidas de las partes, y esto nos plantea la posibilidad de tomar dicha suma de medidas como definición de medida cuando la unión no sea medible Jordan. Admitiendo que esto lleve a un concepto coherente de medida, vemos que, con este criterio, todo conjunto numerable tendrá medida 0, pues será unión numerable de conjuntos con un solo punto, que tienen medida 0. En particular esto resuelve que el conjunto $A = [0, 1]^2 \cap \mathbb{Q}^2$ tiene medida 0, pues es numerable.

En las secciones siguientes veremos que, en efecto, incorporando el criterio de que las uniones disjuntas numerables de conjuntos medibles son medibles y su medida es la suma infinita de las medidas, llegamos a una extensión de la medida de Jordan respecto a la cual todos los conjuntos abiertos y todos los cerrados son medibles.

8.2 Medidas

Empecemos dando nombre a la situación dada por el último teorema de la sección anterior:

Definición 8.19 Una *medida* en un anillo \mathcal{A} sobre un conjunto X es una aplicación $\mu : \mathcal{A} \rightarrow [0, +\infty]$ que cumple $\mu(\emptyset) = 0$ y, para toda familia $\{A_k\}_{k=0}^{\infty}$ de elementos de \mathcal{A} disjuntos dos a dos tal que $\bigcup_{k=0}^{\infty} A_k \in \mathcal{A}$, se cumple que

$$\mu\left(\bigcup_{k=0}^{\infty} A_k\right) = \sum_{k=0}^{\infty} \mu(A_k),$$

donde hay que entender que la suma toma el valor $+\infty$ si alguno de los términos $\mu(A_k)$ toma el valor $+\infty$ o, en caso contrario, si la sucesión de sumas parciales tiende a $+\infty$.

Observemos que toda medida es en particular una medida finitamente aditiva,² pues si tenemos dos conjuntos disjuntos A_0 y A_1 en \mathcal{A} , podemos aplicar la definición anterior a la sucesión que resulta de tomar $A_k = \emptyset$ para $k \geq 2$.

En estos términos tenemos que la medida de Jordan es una medida en \mathcal{J}^n , y también que la medida de figuras elementales es una medida en \mathcal{E}^n .

Observemos que hemos admitido la posibilidad de que la medida de un conjunto sea $+\infty$. La idea es que queremos extender la medida de Jordan de modo que todos los abiertos y todos los cerrados de \mathbb{R}^n sean medibles, lo cual exige en particular asignarle una medida, por ejemplo, a \mathbb{R}^n , y ésta no puede ser sino $m(\mathbb{R}^n) = +\infty$.

El teorema siguiente recoge las propiedades elementales de las medidas. Las demostraciones no ofrecen dificultad. Adoptamos el convenio de que si $a \in \mathbb{R}$ entonces $a + \infty = +\infty + \infty = +\infty$. Cuando hablamos de conjuntos *medibles* nos referimos a conjuntos del anillo en el que está definida una medida dada.

Teorema 8.20 μ una medida en un anillo \mathcal{A} sobre un conjunto X .

- Si $A \subset B$ son medibles entonces $\mu(A) \leq \mu(B)$.
- Si $A \subset B$ son medibles y $\mu(A) < +\infty$, entonces $\mu(B \setminus A) = \mu(B) - \mu(A)$.
- Si A y B son conjuntos medibles disjuntos $\mu(A \cup B) = \mu(A) + \mu(B)$.
- Si A y B son medibles entonces $\mu(A \cup B) \leq \mu(A) + \mu(B)$.
- Si $\{A_n\}_{n=0}^{\infty}$ son medibles y $\bigcup_{n=0}^{\infty} A_n$ también lo es, entonces

$$\mu\left(\bigcup_{n=0}^{\infty} A_n\right) \leq \sum_{n=0}^{\infty} \mu(A_n).$$

²Hubiera sido más natural llamar “medidas” a las medidas simplemente aditivas y “medidas numerablemente aditivas” a lo que hemos llamado “medidas”, pues así tanto unas como otras serían medidas. Sin embargo, como en lo sucesivo vamos a trabajar únicamente con medidas numerablemente aditivas, es preferible reservar para ellas el término “medida”.

f) Si $\{A_n\}_{n=0}^{\infty}$ y $\bigcup_{n=0}^{\infty} A_n$ son medibles y cada $A_n \subset A_{n+1}$, entonces

$$\mu\left(\bigcup_{n=0}^{\infty} A_n\right) = \sup_n \mu(A_n).$$

g) Si $\{A_n\}_{n=0}^{\infty}$ y $\bigcap_{n=0}^{\infty} A_n$ son medibles, cada $A_{n+1} \subset A_n$ y $\mu(A_0) < +\infty$, entonces

$$\mu\left(\bigcap_{n=0}^{\infty} A_n\right) = \inf_n \mu(A_n).$$

DEMOSTRACIÓN: Las primeras propiedades se demuestran igual que para medidas finitamente aditivas sobre anillos (el hecho de que las medidas puedan tomar el valor infinito no afecta a dichos argumentos).

e) Definimos $B_n = A_n \setminus \bigcup_{i=0}^{n-1} A_i$ (entendiendo que $B_0 = A_0$), de modo que $B_n \in \mathcal{A}$, $\bigcup_{n=0}^{\infty} B_n = \bigcup_{n=0}^{\infty} A_n \in \mathcal{A}$ y los conjuntos B_n son disjuntos dos a dos. Por consiguiente,

$$\mu\left(\bigcup_{n=0}^{\infty} A_n\right) = \mu\left(\bigcup_{n=0}^{\infty} B_n\right) = \sum_{n=0}^{\infty} \mu(B_n) \leq \sum_{n=0}^{\infty} \mu(A_n).$$

f) En este caso basta definir $B_n = A_n \setminus A_{n-1}$, con lo que

$$\begin{aligned} \mu\left(\bigcup_{n=0}^{\infty} A_n\right) &= \mu\left(\bigcup_{n=0}^{\infty} B_n\right) = \sum_{n=0}^{\infty} \mu(B_n) = \mu(A_0) + \lim_k \sum_{n=1}^k (\mu(A_n) - \mu(A_{n-1})) \\ &= \lim_k \mu(A_k) = \sup_k \mu(A_k), \end{aligned}$$

donde la última igualdad se debe a que la sucesión $\mu(A_k)$ es monótona creciente.

g) Aplicamos la propiedad anterior a los conjuntos $B_n = A_0 \setminus A_n$. Así $\bigcup_{n=0}^{\infty} B_n = A_0 \setminus \bigcap_{n=0}^{\infty} A_n \in \mathcal{A}$, luego

$$\mu(A_0) \setminus \mu\left(\bigcap_{n=0}^{\infty} A_n\right) = \mu\left(\bigcup_{n=0}^{\infty} B_n\right) = \sup_n \mu(B_n) = \mu(A_0) - \inf_n \mu(A_n),$$

de donde se sigue la conclusión. \blacksquare

Nuestro propósito es demostrar que toda medida en un anillo se puede extender a una medida en un conjunto mayor que vuelva redundantes las condiciones sobre la medibilidad de las uniones e intersecciones numerables. Esto nos lleva a introducir los conceptos siguientes:

Definición 8.21 Un σ -anillo sobre un conjunto X es un anillo \mathcal{A} sobre X con la propiedad adicional de que si $\{A_n\}_{n=0}^{\infty}$ es una familia de elementos de \mathcal{A} , entonces $\bigcup_{n=0}^{\infty} A_n \in \mathcal{A}$.

Observemos que esto implica que $\bigcap_{n=0}^{\infty} A_n \in \mathcal{A}$, pues

$$\bigcap_{n=0}^{\infty} A_n = A_0 \setminus \bigcup_{n=0}^{\infty} (A_0 \setminus A_n).$$

Un *álgebra* de subconjuntos de X es una familia \mathcal{A} de subconjuntos de X tal que:

- a) $\emptyset, X \in \mathcal{A}$.
- b) Si $A \in \mathcal{A}$, entonces $X \setminus A \in \mathcal{A}$.
- c) Si $A, B \in \mathcal{A}$, entonces $A \cup B, A \cap B \in \mathcal{A}$.

La propiedad b) hace que la propiedad c) para uniones implique la parte para intersecciones y viceversa. Es claro que toda álgebra de conjuntos es en particular un anillo, y un anillo \mathcal{A} de subconjuntos de X es un álgebra si y sólo si $X \in \mathcal{A}$.

Una σ -*álgebra* de subconjuntos de X es un álgebra \mathcal{A} en la que la condición c) se sustituye por la condición (más fuerte) de que la unión de toda familia numerable de elementos de \mathcal{A} está en \mathcal{A} (y, por consiguiente, lo mismo vale para las intersecciones numerables de elementos de \mathcal{A}). Equivalentemente, una σ -álgebra es un álgebra que además es un σ -anillo.

De este modo, el concepto de σ -álgebra “captura” dos de nuestras aspiraciones: extender una medida dada en un anillo de subconjuntos de un conjunto X a otra para la cual 1) las uniones e intersecciones numerables de conjuntos medibles sean medibles y 2) el propio conjunto X sea medible (aunque tal vez deba tener medida infinita).

Un *espacio medida* es una terna (X, \mathcal{A}, μ) , donde X es un conjunto, \mathcal{A} es una σ -álgebra de subconjuntos de X y μ es una medida en \mathcal{A} . Como es habitual en estos casos, en la práctica escribiremos X en lugar de (X, \mathcal{A}, μ) . A los elementos de \mathcal{A} los llamaremos *subconjuntos medibles* de X .

A medio plazo trabajaremos únicamente con espacios medida, pero de momento tenemos que trabajar con medidas sobre anillos arbitrarios para obtener resultados sobre existencia de medidas en determinadas σ -álgebras, pues las obtendremos extendiendo medidas definidas en anillos, como es el caso de la medida de Jordan.

Una medida que sólo tome valores finitos (y no sea idénticamente nula) es una *medida finita*. Si está definida sobre un álgebra de subconjuntos de X , esto equivale a que $0 < \mu(X) < +\infty$.

Diremos que una medida μ es *unitaria* si $\mu(X) = 1$. Diremos que μ es σ -*finita* si $\mu(X) > 0$ y existen conjuntos medibles $\{A_n\}_{n=0}^{\infty}$ de medida finita tales que

$$X = \bigcup_{n=0}^{\infty} A_n.$$

Éste es el caso del área en el plano o el volumen en el espacio. El área del plano es infinita, pero podemos descomponerlo en una unión numerable de bolas de área finita. También se habla de *espacios medida unitarios, finitos* o *σ -finitos* según sea la medida definida en ellos.

Los conjuntos de medida 0 se llaman también *conjuntos nulos*. Una medida es *completa* si todo subconjunto de un conjunto nulo es medible (y, por consiguiente, nulo). Es claro que la medida de Jordan es completa.

Antes de ocuparnos del problema de la extensión de una medida dada tenemos que estudiar la existencia de σ -álgebras de subconjuntos de un conjunto dado.

Trivialmente, el conjunto $\mathcal{P}X$ de todos los subconjuntos de X es una σ -álgebra de subconjuntos de X , pero en general es demasiado grande para que estemos en condiciones de definir medidas de interés en $\mathcal{P}X$.

Ahora bien, es inmediato comprobar que la intersección de cualquier familia de σ -anillos (o de σ -álgebras) de subconjuntos de un conjunto X cumple también la definición de σ -anillo (o de σ -álgebra). Por consiguiente, si $G \subset \mathcal{P}X$ es una familia de subconjuntos de X , podemos definir el σ -anillo (o la σ -álgebra) generada por G como la intersección de todos los σ -anillos (resp. σ -álgebras) de subconjuntos de X que contienen a G .

Observemos que si \mathcal{A} es un anillo sobre un conjunto X y se cumple que X es unión numerable de elementos de \mathcal{A} , entonces el σ -anillo generado por \mathcal{A} es de hecho una σ -álgebra, pues contiene a X , luego coincide con la σ -álgebra generada por \mathcal{A} .

El σ -anillo generado por un anillo tiene una descripción más simple que la que proporciona la definición y que es interesante conocer. Para ello tenemos que introducir una última estructura de familias de subconjuntos:

Una *clase monótona* \mathcal{M} en un conjunto X es una colección de subconjuntos de X tal que si $\{A_n\}_{n=0}^{\infty}$ es una familia creciente de conjuntos de \mathcal{M} (es decir, $A_n \subset A_{n+1}$) entonces $\bigcup_{n=0}^{\infty} A_n \in \mathcal{M}$ y si la familia es decreciente ($A_{n+1} \subset A_n$) entonces $\bigcap_{n=0}^{\infty} A_n \in \mathcal{M}$.

Es claro que la intersección de clases monótonas es de nuevo una clase monótona, por lo que podemos hablar también de la clase monótona generada por un conjunto, es decir, la menor clase monótona que lo contiene. El resultado fundamental sobre clases monótonas es el siguiente:

Teorema 8.22 *Si \mathcal{A} es un anillo de subconjuntos de un conjunto X , entonces la clase monótona generada por \mathcal{A} coincide con el σ -anillo generado por \mathcal{A} .*

DEMOSTRACIÓN: Sea \mathcal{M} la clase monótona generada por \mathcal{A} y sea \mathcal{S} el σ -anillo generado por \mathcal{A} . Como \mathcal{S} es trivialmente una clase monótona que contiene a \mathcal{A} , tiene que ser $\mathcal{M} \subset \mathcal{S}$, y si probamos que \mathcal{M} es un σ -anillo, entonces tendremos también que $\mathcal{S} \subset \mathcal{M}$.

Para cada $A \subset X$, definimos \mathcal{M}_A como el conjunto de todos los $B \subset X$ tales que $A \setminus B$, $B \setminus A$, $A \cup B \in \mathcal{M}$. Trivialmente

$$B \in \mathcal{M}_A \quad \text{si y sólo si} \quad A \in \mathcal{M}_B.$$

Veamos que \mathcal{M}_A es una clase monótona. Si $\{A_n\}_{n=0}^{\infty}$ es una sucesión creciente en \mathcal{M}_A , entonces

$$A \setminus \bigcup_{n=0}^{\infty} A_n = \bigcup_{n=0}^{\infty} (A \setminus A_n) \in \mathcal{M},$$

porque cada $A \setminus A_n \in \mathcal{M}$ y forma una sucesión decreciente. Igualmente

$$\bigcup_{n=0}^{\infty} A_n \setminus A = \bigcup_{n=0}^{\infty} (A_n \setminus A) \in \mathcal{M},$$

porque cada $A_n \setminus A \in \mathcal{M}$ y forman una sucesión creciente. Por último

$$A \cup \bigcup_{n=0}^{\infty} A_n = \bigcup_{n=0}^{\infty} (A \cup A_n) \in \mathcal{M},$$

porque $A \cup A_n \in \mathcal{M}$ y forman una sucesión creciente.

Esto prueba que $\bigcup_{n=0}^{\infty} A_n \in \mathcal{M}_A$. Igualmente se prueba que si la familia es decreciente su intersección está en \mathcal{M}_A .

Ahora, fijado $A \in \mathcal{A}$, observamos que, para todo $B \in \mathcal{A}$, se cumple que $B \in \mathcal{M}_A$, luego $\mathcal{A} \subset \mathcal{M}_A$, y como \mathcal{M}_A es una clase monótona que contiene a \mathcal{A} , tiene que ser $\mathcal{M} \subset \mathcal{M}_A$.

Por consiguiente, fijado $B \in \mathcal{M}$, para todo $A \in \mathcal{A}$ tenemos que $B \in \mathcal{M}_A$, luego $A \in \mathcal{M}_B$, luego $\mathcal{A} \subset \mathcal{M}_B$ y así \mathcal{M}_B es una clase monótona que contiene a \mathcal{A} , luego $\mathcal{M} \subset \mathcal{M}_B$.

Así pues, si $A, B \in \mathcal{M}$, como $\mathcal{M} \subset \mathcal{M}_B$, se cumple $A \setminus B$, $B \setminus A$, $A \cup B \in \mathcal{M}$, lo cual implica que \mathcal{M} es un anillo, y al ser un anillo y una clase monótona, es de hecho un σ -anillo, pues si $\{A_n\}_{n=0}^{\infty}$ es una familia de elementos de \mathcal{M} , llamamos $B_n = \bigcup_{i=0}^n A_i \in \mathcal{M}$ (porque \mathcal{M} es un anillo) y entonces

$$\bigcup_{n=0}^{\infty} A_n = \bigcup_{n=0}^{\infty} B_n \in \mathcal{M},$$

porque es una clase monótona. ■

El primer paso para extender una medida en un anillo al σ -anillo que genera es definir a partir de ella una medida exterior en el sentido siguiente:

Definición 8.23 Una aplicación $\mu^* : \mathcal{P}X \rightarrow [0, +\infty]$ es una *medida exterior* en X si cumple las propiedades siguientes:

- a) $\mu^*(\emptyset) = 0$.
- b) Si $A \subset B \subset X$, entonces $\mu^*(A) \leq \mu^*(B)$.

c) Si $\{A_k\}_{k=0}^{\infty}$ es una sucesión de subconjuntos de X , entonces

$$\mu^* \left(\bigcup_{k=0}^{\infty} A_k \right) \leq \sum_{k=0}^{\infty} \mu^*(A_k).$$

Teorema 8.24 Si $\mu : \mathcal{A} \rightarrow [0, +\infty[$ es una medida en un anillo \mathcal{A} de subconjuntos de X , entonces la aplicación $\mu^* : \mathcal{P}X \rightarrow [0, +\infty]$ dada por

$$\mu^*(A) = \inf \left\{ \sum_{k=0}^{\infty} \mu(A_k) \mid A_k \in \mathcal{A}, A \subset \bigcup_{k=0}^{\infty} A_k \right\}$$

(con el convenio de que $\inf \emptyset = +\infty$) es una medida exterior en X .

DEMOSTRACIÓN: Claramente $\mu^*(\emptyset) = 0$ y si $A \subset B \subset X$ se cumple que $\mu^*(A) \leq \mu^*(B)$.

Sea $\{A_k\}_{k=0}^{\infty}$ una sucesión de subconjuntos de X tal que $\mu^*(A_k) < +\infty$ para todo k . Dado $\epsilon > 0$ podemos tomar una sucesión $\{A_m^k\}_{m=0}^{\infty}$ en \mathcal{A} tal que $A_k \subset \bigcup_{m=0}^{\infty} A_m^k$ y

$$\sum_{m=0}^{\infty} \mu(A_m^k) < \mu^*(A_k) + \frac{\epsilon}{2^k}.$$

Al unir estas sucesiones para todo k obtenemos un cubrimiento de $\bigcup_{k=0}^{\infty} A_k$ del que deducimos que

$$\mu^* \left(\bigcup_{k=0}^{\infty} A_k \right) \leq \sum_{k=0}^{\infty} \mu^*(A_k) + \epsilon,$$

para todo $\epsilon > 0$, luego se cumple la desigualdad sin ϵ y μ^* es una medida exterior. (Para sucesiones con algún k tal que $\mu^*(A_k) = +\infty$ la desigualdad es trivial.)

Si $A \in \mathcal{A}$, como $A \subset A$, tenemos que $\mu^*(A) \leq \mu(A)$.

Si $A \subset \bigcup_{k=0}^{\infty} A_k$, con $A_k \in \mathcal{A}$, definimos

$$B_k = (A_k \cap A) \setminus \bigcup_{i < k} (A_i \cap A) \in \mathcal{A}.$$

Claramente, los B_k son disjuntos dos a dos y $A = \bigcup_{k=0}^{\infty} B_k$. Como μ es una medida,

$$\mu(A) = \sum_{k=0}^{\infty} \mu(B_k) \leq \sum_{k=0}^{\infty} \mu(A_k)$$

y, como esto vale para todo cubrimiento de A , tenemos que $\mu(A) \leq \mu^*(A)$, luego μ^* extiende a μ . ■

La medida exterior está definida en todos los subconjuntos de X , pero para garantizar que cumpla las propiedades de una medida tenemos que restringirla a una familia menor:

Definición 8.25 Si μ^* es una medida exterior en un conjunto X , llamaremos *conjuntos μ^* -medibles* a los conjuntos $A \subset X$ tales que, para todo $B \subset X$, se cumple que $\mu^*(B) = \mu^*(B \cap A) + \mu^*(B \setminus A)$.

Notemos que siempre se cumple que

$$\mu^*(B) = \mu^*((B \cap A) \cup (B \setminus A)) \leq \mu^*(B \cap A) + \mu^*(B \setminus A).$$

La medibilidad equivale a la desigualdad opuesta.

Teorema 8.26 Si μ^* es una medida exterior en un conjunto X , el conjunto \mathcal{M} de todos los subconjuntos de X que son μ^* -medibles es una σ -álgebra, y la restricción μ de μ^* a \mathcal{M} es una medida completa.

DEMOSTRACIÓN: Trivialmente $\emptyset, X \in \mathcal{M}$. También es claro que si $A \in \mathcal{M}$ entonces $X \setminus A \in \mathcal{M}$. Sean $A, B \in \mathcal{M}$ y $C \subset X$ arbitrario. Entonces

$$\mu^*(C) = \mu^*(C \cap A) + \mu^*(C \setminus A), \quad \mu^*(C \setminus A) = \mu^*((C \setminus A) \cap B) + \mu^*(C \setminus (A \cup B)),$$

luego

$$\begin{aligned} \mu^*(C) &= \mu^*(C \cap A) + \mu^*((C \setminus A) \cap B) + \mu^*(C \setminus (A \cup B)) \\ &\geq \mu^*(C \cap (A \cup B)) + \mu^*(C \setminus (A \cup B)), \end{aligned}$$

donde hemos usado que $C \cap (A \cup B) = (C \cap A) \cup ((C \setminus A) \cap B)$. Esto prueba que $A \cup B \in \mathcal{M}$. Hasta aquí tenemos probado que \mathcal{M} es un álgebra de conjuntos.

Supongamos ahora que $\{A_k\}_{k=0}^{\infty}$ es una familia de conjuntos μ^* -medibles disjuntos dos a dos y $S_m = \bigcup_{k=0}^m A_k$. Entonces, para todo $A \subset X$, se cumple que

$$\mu^*(A \cap S_m) = \sum_{k=0}^m \mu^*(A \cap A_k).$$

En efecto, lo probamos por inducción sobre m . Para $m = 0$ es trivial. Sabemos que $S_{m+1} \in \mathcal{M}$, luego

$$\begin{aligned} \mu^*(A \cap S_{m+1}) &= \mu^*(A \cap S_{m+1} \cap S_m) + \mu^*((A \cap S_{m+1}) \setminus S_m) \\ &= \mu^*(A \cap S_m) + \mu^*(A \cap A_{m+1}) = \sum_{k=0}^m \mu^*(A \cap A_k) + \mu^*(A \cap A_{m+1}) \\ &= \sum_{k=0}^{m+1} \mu^*(A \cap A_k). \end{aligned}$$

De aquí deducimos a su vez que

$$\mu^*(A \cap \bigcup_{k=0}^{\infty} A_k) = \sum_{k=0}^{\infty} \mu^*(A \cap A_k).$$

En efecto, llamando $S = \bigcup_{k=0}^{\infty} A_k$, como $S_m \subset S$,

$$\mu^*(A \cap S) \geq \mu^*(A \cap S_m) = \sum_{k=0}^m \mu^*(A \cap A_k),$$

luego la serie converge y $\mu^*(A \cap S) \geq \sum_{k=0}^{\infty} \mu^*(A \cap A_k)$. La desigualdad opuesta se da por definición de medida exterior.

A su vez de aquí obtenemos que $\bigcup_{k=0}^{\infty} A_k \in \mathcal{M}$. En efecto, si $A \subset X$ es arbitrario,

$$\mu^*(A) = \mu^*(A \cap S_m) + \mu^*(A \setminus S_m) \geq \sum_{k=0}^m \mu^*(A \cap A_k) + \mu^*(A \setminus S).$$

Como esto vale para todo m ,

$$\mu^*(A) \geq \sum_{k=0}^{\infty} \mu^*(A \cap A_k) + \mu^*(A \setminus S) = \mu^*(A \cap S) + \mu^*(A \setminus S).$$

No hemos probado todavía que \mathcal{M} es una σ -álgebra porque estamos suponiendo que los A_k son disjuntos dos a dos, pero, si no lo son, definimos $B_k = A_k \setminus \bigcup_{i=0}^{k-1} A_i \in \mathcal{M}$ y, por el caso ya probado para conjuntos disjuntos, tenemos que

$$\bigcup_{k=0}^{\infty} A_k = \bigcup_{k=0}^{\infty} B_k \in \mathcal{M}.$$

Haciendo $A = X$ en la igualdad

$$\mu^*(A \cap \bigcup_{k=0}^{\infty} A_k) = \sum_{k=0}^{\infty} \mu^*(A \cap A_k)$$

ya demostrada, concluimos que μ^* se restringe a una medida μ en \mathcal{M} .

Que μ sea completa significa que si $A \in \mathcal{M}$, $B \subset A$ y $\mu(A) = 0$, entonces $B \in \mathcal{M}$. En efecto, tenemos que $\mu^*(B) = 0$, y entonces es μ^* -medible, pues, para todo $C \subset X$, tenemos que $\mu^*(C \cap B) + \mu^*(C \setminus B) \leq \mu^*(B) + \mu^*(C) = \mu^*(C)$. ■

Combinando los dos últimos teoremas obtenemos:

Teorema 8.27 (Caratheodory) *Toda medida $\mu : \mathcal{A} \rightarrow [0, +\infty[$ en un anillo \mathcal{A} de subconjuntos de X se extiende a una medida completa $\bar{\mu} : \mathcal{M} \rightarrow [0, +\infty[$ definida sobre cierta σ -álgebra \mathcal{M} en X que contiene a \mathcal{A} .*

DEMOSTRACIÓN: Consideramos la medida exterior μ^* dada por 8.24, que extiende a μ y sea \mathcal{M} su σ -álgebra de conjuntos medibles. Sólo hemos de probar

que $\mathcal{A} \subset \mathcal{M}$. Para ello tomamos $A \in \mathcal{A}$ y $B \subset X$ arbitrario. Sea $\{A_k\}_{k=0}^{\infty}$ una familia de elementos de \mathcal{A} tales que $A \subset \bigcup_{k=0}^{\infty} A_k$.

Entonces $B \cap A \subset \bigcup_{k=0}^{\infty} (A_k \cap A)$ y $B \setminus A \subset \bigcup_{k=0}^{\infty} (A_k \setminus A)$, luego

$$\mu^*(B \cap A) + \mu^*(B \setminus A) \leq \sum_{k=0}^{\infty} \mu(A_k \cap A) + \sum_{k=0}^{\infty} \mu(A_k \setminus A) = \sum_{k=0}^{\infty} \mu(A_k),$$

porque μ es una medida. Como esto vale para todo cubrimiento de A , resulta que

$$\mu^*(B \cap A) + \mu^*(B \setminus A) \leq \mu^*(A),$$

luego $A \in \mathcal{M}$. ■

Observemos que la “cierta σ -álgebra” a la que alude el enunciado del teorema anterior es la σ -álgebra de los conjuntos μ^* -medibles, donde μ^* es la medida exterior definida por μ . Vamos a tratar de describir mejor esta σ -álgebra. En primer lugar demostramos un resultado de unicidad:

Teorema 8.28 *Toda medida $\mu : \mathcal{A} \rightarrow [0, +\infty[$ definida en un anillo \mathcal{A} de subconjuntos de X admite una única extensión al σ -anillo generado por \mathcal{A} .*

DEMOSTRACIÓN: Sea \mathcal{S} el σ -anillo generado por \mathcal{A} . La medida dada por el teorema anterior se restringe trivialmente a una medida $\mu_1 : \mathcal{S} \rightarrow [0, +\infty]$. Supongamos que existiera otra extensión $\mu_2 : \mathcal{S} \rightarrow [0, +\infty]$. Definimos

$$\mathcal{M} = \{A \in \mathcal{S} \mid \mu_1(A) = \mu_2(A)\}.$$

Claramente $\mathcal{A} \subset \mathcal{M}$. Basta probar que \mathcal{M} es una clase monótona, pues entonces $\mathcal{M} = \mathcal{S}$, luego $\mu_1 = \mu_2$. Ahora bien, si $\{A_n\}_{n=0}^{\infty}$ es una sucesión creciente en \mathcal{M} , entonces $\bigcup_{n=0}^{\infty} A_n \in \mathcal{S}$, porque \mathcal{S} es un σ -anillo, y

$$\mu_1 \left(\bigcup_{n=0}^{\infty} A_n \right) = \sup_n \mu_1(A_n) = \sup_n \mu_2(A_n) = \mu_2 \left(\bigcup_{n=0}^{\infty} A_n \right),$$

e igualmente se razona con sucesiones decrecientes. ■

En general, la σ -álgebra dada por el teorema 8.27 es mayor que el σ -anillo generado por el anillo de partida. El teorema siguiente da cuenta de la diferencia:

Teorema 8.29 *Sea X un conjunto $\mu : \mathcal{S} \rightarrow [0, +\infty]$ una medida definida sobre una σ -álgebra de subconjuntos de X . Sea*

$$\mathcal{M} = \{A \subset X \mid \text{existen } B, C \in \mathcal{S} \text{ tales que } B \subset A \subset C \text{ y } \mu(C \setminus B) = 0\}.$$

Entonces \mathcal{M} es una σ -álgebra de subconjuntos de X que contiene a \mathcal{S} y μ se extiende a una única medida en \mathcal{M} , que es completa.

DEMOSTRACIÓN: Si $A \in \mathcal{S}$ es claro que $A \in \mathcal{M}$. Basta tomar $B = C = A$. Así pues $\mathcal{S} \subset \mathcal{M}$, luego en particular $\emptyset, X \in \mathcal{M}$.

Si $A \in \mathcal{M}$, sean $B, C \in \mathcal{S}$ tales que $B \subset A \subset C$ y $\mu(C \setminus B) = 0$. Entonces $X \setminus C \subset X \setminus A \subset X \setminus B$, y claramente $X \setminus C, X \setminus A \in \mathcal{S}$ y

$$\mu((X \setminus B) \setminus (X \setminus C)) = \mu(C \setminus B) = 0.$$

Por lo tanto $X \setminus A \in \mathcal{M}$. Para probar que \mathcal{M} es una σ -álgebra basta probar que la unión numerable de elementos de \mathcal{M} está en \mathcal{M} . Sea, pues, $\{A_k\}_{k=0}^{\infty}$ una familia de elementos de \mathcal{M} . Sean $\{B_k\}_{k=0}^{\infty}$ y $\{C_k\}_{k=0}^{\infty}$ según la definición de \mathcal{M} . Entonces

$$\bigcup_{k=0}^{\infty} B_k \subset \bigcup_{k=0}^{\infty} A_k \subset \bigcup_{k=0}^{\infty} C_k,$$

los conjuntos de los extremos están en \mathcal{S} y

$$0 \leq \mu\left(\left(\bigcup_{k=0}^{\infty} C_k\right) \setminus \left(\bigcup_{k=0}^{\infty} B_k\right)\right) \leq \mu\left(\bigcup_{k=0}^{\infty} C_k \setminus B_k\right) = 0.$$

Por lo tanto $\bigcup_{k=0}^{\infty} A_k \in \mathcal{M}$. Esto prueba que \mathcal{M} es una σ -álgebra.

Si $A \in \mathcal{M}$ y B, C son los conjuntos dados por la definición, es claro que $\mu(B) = \mu(C)$. Veamos que podemos extender μ mediante $\mu(A) = \mu(B) = \mu(C)$. Obviamente es la única extensión posible. En efecto, si B' y C' también cumplen la definición, entonces $B \setminus B' \subset C' \setminus B'$, luego $\mu(B \setminus B') \leq \mu(C' \setminus B') = 0$, de donde $\mu(B \setminus B') = 0$ y $\mu(B) = \mu(B')$, luego B y B' dan lugar al mismo valor de $\mu(A)$.

Veamos que la extensión de μ que acabamos de definir es realmente una medida. Para ello tomamos una familia $\{A_k\}_{k=0}^{\infty}$ de elementos de \mathcal{M} disjuntos dos a dos. Sean $\{B_k\}$ y $\{C_k\}$ elementos de \mathcal{S} que satisfagan la definición de \mathcal{M} . Según hemos visto, los conjuntos $\bigcup_{k=0}^{\infty} B_k$ y $\bigcup_{k=0}^{\infty} C_k$ justifican que $\bigcup_{k=0}^{\infty} A_k \in \mathcal{M}$ y es claro que los B_k son disjuntos dos a dos, luego

$$\mu\left(\bigcup_{k=0}^{\infty} A_k\right) = \mu\left(\bigcup_{k=0}^{\infty} B_k\right) = \sum_{k=0}^{\infty} \mu(B_k) = \sum_{k=0}^{\infty} \mu(A_k).$$

La medida en \mathcal{M} es completa, pues si $A \in \mathcal{M}$ es nulo y $D \subset A$, tomamos $B, C \in \mathcal{S}$ tales que $B \subset A \subset C$ con $\mu(B) = \mu(C) = 0$. Claramente $\emptyset \subset D \subset C$ y $\mu(C \setminus \emptyset) = 0$, luego $D \in \mathcal{M}$. ■

La medida construida en el teorema anterior se conoce como *compleción* de μ .

Teorema 8.30 *Sea $\mu : \mathcal{A} \rightarrow [0, +\infty[$ una medida en un anillo \mathcal{A} de subconjuntos de X y supongamos que X se expresa como unión numerable de elementos de \mathcal{A} . Entonces la medida dada por la prueba del teorema 8.27 es la compleción de la única extensión de μ a la σ -álgebra generada por \mathcal{A} .*

DEMOSTRACIÓN: Sea \mathcal{S} el σ -anillo generado por \mathcal{A} y sea \mathcal{M} la σ -álgebra de los conjuntos μ^* -medibles, sobre la que está definida la medida $\bar{\mu}$ construida en el teorema 8.27. Observemos en primer lugar que la hipótesis sobre X implica que $X \in \mathcal{S}$, por lo que \mathcal{S} es, de hecho, una σ -álgebra. Ciertamente $\mathcal{S} \subset \mathcal{M}$ y $\bar{\mu}|_{\mathcal{S}}$ es la única extensión de μ a \mathcal{S} . Por otra parte, $\bar{\mu}$ es una medida completa, luego basta probar que \mathcal{M} es la σ -álgebra \mathcal{M}^* dada por el teorema anterior.

Si $A \in \mathcal{M}^*$, existen $B, C \in \mathcal{S}$ tales que $B \subset A \subset C$ y $\bar{\mu}(C \setminus B) = 0$, entonces $A = B \cup (A \setminus B)$, donde $A \setminus B \subset C \setminus B$ y, como $\bar{\mu}$ es completa, $A \setminus B \in \mathcal{M}$, luego $A \in \mathcal{M}$. Esto nos da una inclusión.

Tomemos ahora $A \in \mathcal{M}$. Fijemos $\{A_n\}_{n=0}^{\infty}$ en \mathcal{A} tales que $X = \bigcup_{n=0}^{\infty} A_n$. Entonces $A = \bigcup_{n=0}^{\infty} (A \cap A_n)$, luego si probamos que $A \cap A_n \in \mathcal{M}^*$, lo mismo valdrá para A . Equivalentemente, como $A \cap A_n \subset A_n$, podemos suponer que existe un $U \in \mathcal{A}$ tal que $A \subset U$.

Dado $m \in \mathbb{N}$ no nulo, por definición de la medida exterior μ^* , existen $\{A_{km}\}_{k=0}^{\infty}$ en \mathcal{A} tales que $A \subset \bigcup_{k=0}^{\infty} A_{km}$ y

$$\sum_{k=0}^{\infty} \mu(A_{km}) < \mu^*(A) + \frac{1}{m} = \bar{\mu}(A) + \frac{1}{m}.$$

Sea $C_m = \bigcup_{k=0}^{\infty} A_{km} \in \mathcal{S}$, de modo que $A \subset C_m$ y $\bar{\mu}(C_m) < \bar{\mu}(A) + 1/m$.

Cambiando C_m por $\bigcap_{i=0}^m C_i$, podemos suponer que $\{C_m\}_{m=1}^{\infty}$ es una sucesión decreciente, por lo que $C = \bigcap_{m=1}^{\infty} C_m \in \mathcal{S}$ cumple que $A \subset C$ y

$$\bar{\mu}(A) \leq \bar{\mu}(C) = \inf_m \bar{\mu}(C_m) \leq \bar{\mu}(A).$$

Por lo tanto, $\bar{\mu}(C) = \bar{\mu}(A)$. Razonamos igualmente con $U \setminus A$, con lo que obtenemos un $C' \in \mathcal{S}$ tal que $U \setminus A \subset C'$ y $\bar{\mu}(C') = \bar{\mu}(U \setminus A) = \bar{\mu}(U) - \bar{\mu}(A)$. Cambiando C' por $C' \cap U$, podemos suponer que $C' \subset U$. Sea $B = U \setminus C' \in \mathcal{S}$, con lo que $B \subset A$ y $\bar{\mu}(B) = \bar{\mu}(U) - \bar{\mu}(C') = \bar{\mu}(A)$. En definitiva, tenemos que $B \subset A \subset C$ y $\bar{\mu}(C \setminus B) = \bar{\mu}(C) - \bar{\mu}(B) = 0$, luego $A \in \mathcal{M}^*$. ■

8.3 La medida de Lebesgue

Vamos a aplicar a la medida de Jordan toda la teoría desarrollada en la sección anterior. En realidad, vamos a ver que la construcción de la medida de Lebesgue puede realizarse indistintamente a partir de la medida de Jordan o de la medida sobre las figuras elementales, por lo que, para llegar a la medida de Lebesgue, la construcción de la medida de Jordan es totalmente prescindible.

Definición 8.31 La *medida exterior de Lebesgue* en \mathbb{R}^n es la medida exterior m^* que se obtiene de aplicar el teorema 8.24 a la medida de figuras elementales $m : \mathcal{E}^n \rightarrow [0, +\infty[$.

No es difícil ver que si en 8.24 cambiamos las sumas infinitas por sumas finitas, lo que obtenemos a partir de la medida de figuras elementales es precisamente la medida exterior de Jordan, luego la diferencia esencial entre la medida de Jordan y la de Lebesgue es precisamente ésta: definir la medida exterior con sumas numerables en lugar de con sumas finitas.

La medida exterior de Lebesgue puede expresarse en términos ligeramente más simples que los que requiere la definición general:

Teorema 8.32 *La medida exterior de Lebesgue viene dada por*

$$m^*(A) = \inf \left\{ \sum_{k=0}^{\infty} m(C_k) \mid C_k \text{ es una celda abierta, } A \subset \bigcup_{k=0}^{\infty} C_k \right\}.$$

DEMOSTRACIÓN: Llamemos $m'(A)$ al ínfimo que aparece en el enunciado. La definición de medida exterior considera cubrimientos por figuras elementales en lugar de por celdas abiertas. Como toda celda abierta es una figura elemental, la desigualdad $m^*(A) \leq m'(A)$ es inmediata. Para probar la contraria tomamos $\epsilon > 0$ y consideramos un cubrimiento $A \subset \bigcup_{k=0}^{\infty} A_k$, donde cada A_k es una figura

elemental y $\sum_{k=0}^{\infty} m(A_k) < m^*(A) + \epsilon/2$. Como cada A_k es unión de un número finito de celdas disjuntas, de modo que $m(A_k)$ es la suma de las medidas de dichas celdas, podemos sustituir la sucesión $\{A_k\}_{k=0}^{\infty}$ por una enumeración de las celdas en que se descomponen los A_k y se sigue cumpliendo que la suma de las medidas es $< \epsilon/2$. En otras palabras, no perdemos generalidad si suponemos que cada A_k es una celda.

Ahora basta aplicar el teorema 8.2, que nos da celdas abiertas $A_k \subset C_k$ tales que $m(C_k) - m(A_k) < \epsilon/2^{k+2}$. De este modo, las celdas C_k cubren A y

$$m'(A) \leq \sum_{k=0}^{\infty} m(C_k) < \sum_{k=0}^{\infty} m(A_k) + \sum_{k=0}^{\infty} \frac{\epsilon}{2^{k+2}} < m^*(A) + \frac{\epsilon}{2} + \frac{\epsilon}{2} = m^*(A) + \epsilon.$$

Por lo tanto $m'(A) \leq m^*(A)$ y tenemos la igualdad. ■

Definición 8.33 La medida $m : \mathcal{M}^n \rightarrow [0, +\infty]$ que resulta de aplicar el teorema 8.27 a la medida de las figuras elementales de \mathbb{R}^n o, equivalentemente, de aplicar el teorema 8.26 a la medida exterior de Lebesgue, recibe el nombre de *medida de Lebesgue* en \mathbb{R}^n , y la σ -álgebra \mathcal{M}^n es la σ -álgebra de los conjuntos *medibles Lebesgue*.

Como $\mathbb{R}^n = \bigcup_{k=0}^{\infty} [-k, k]^n$ y las celdas tienen medida finita, vemos que la medida de Lebesgue es σ -finita. Más aún, esto muestra que es aplicable el

teorema 8.30, según el cual la medida de Lebesgue en \mathcal{M}^n es la completación de la única extensión de la medida de las figuras elementales a la σ -álgebra generada por \mathcal{E}^n . Vamos a describir esta σ -álgebra.

Si X es un espacio topológico la σ -álgebra generada por los conjuntos abiertos recibe el nombre de *σ -álgebra de Borel*. Una medida definida sobre la σ -álgebra de Borel de un espacio topológico X recibe el nombre de *medida de Borel* en X .

Llamaremos \mathcal{B}^n a la σ -álgebra de Borel en \mathbb{R}^n , y vamos a probar que es también la σ -álgebra \mathcal{S}^n generada por \mathcal{E}^n , con lo que en particular $\mathcal{B}^n \subset \mathcal{M}^n$.

Conviene observar en general que si $f : X \rightarrow Y$ es una aplicación continua entre espacios topológicos, entonces la antiimagen de un conjunto de Borel en Y es un conjunto de Borel en X . En efecto, basta considerar el conjunto

$$\mathcal{C} = \{A \subset Y \mid f^{-1}[A] \text{ es de Borel}\}.$$

Se comprueba trivialmente que es una σ -álgebra que contiene a los abiertos, por lo que contiene, de hecho, a todos los conjuntos de Borel de Y .

De aquí deducimos que toda celda en \mathbb{R}^n es un conjunto de Borel, pues si $C = \prod_{i=1}^n I_i$, podemos expresarlo también como $\bigcap_{i=1}^n p_i^{-1}[I_i]$, donde $p_i : \mathbb{R}^n \rightarrow \mathbb{R}$ es la proyección i -ésima, y todo intervalo I_i es un conjunto de Borel en \mathbb{R} (por ejemplo, $[a, b[= \{a\} \cup]a, b[$ es unión de un abierto y un cerrado, es decir, de un abierto y del complementario de un abierto, luego está en la σ -álgebra de Borel).

A su vez esto implica que $\mathcal{E}^n \subset \mathcal{B}^n$, luego $\mathcal{S}^n \subset \mathcal{B}^n$. Para probar la inclusión opuesta basta ver que \mathcal{S}^n contiene a todos los abiertos de \mathbb{R}^n . Para ello basta tener en cuenta que los productos de intervalos abiertos de extremos racionales forman una base numerable de \mathbb{R}^n (por el teorema 2.17), luego todo abierto en \mathbb{R}^n es unión numerable de celdas abiertas, y las celdas abiertas están en \mathcal{E}^n , luego en \mathcal{S}^n , y como ésta es una σ -álgebra, concluimos que todos los abiertos están en \mathcal{S}^n .

Por consiguiente, la restricción de la medida de Lebesgue a \mathcal{B}^n es una medida de Borel en \mathbb{R}^n , que recibe el nombre de *medida de Lebesgue incompleta*, o simplemente medida de Lebesgue cuando por el contexto está claro que se trata de dicha restricción.

Hay muchas formas distintas de construir la medida de Lebesgue, además de la que hemos empleado. El teorema siguiente implica que todas son equivalentes:

Teorema 8.34 *La medida de Lebesgue es la única medida de Borel en \mathbb{R}^n que sobre las celdas abiertas viene dada por*

$$m\left(\prod_{i=1}^n]a_i, b_i[\right) = \prod_{i=1}^n (b_i - a_i).$$

DEMOSTRACIÓN: Sea $\mu : \mathcal{B}^n \rightarrow [0, +\infty]$ una medida de Borel que a cada celda abierta le asigne su contenido. Ya hemos visto que \mathcal{B}^n es la σ -álgebra generada por \mathcal{E}^n . Teniendo en cuenta que

$$\prod_{i=1}^n [a_i, b_i] = \bigcap_{k=1}^{\infty} \prod_{i=1}^n]a_i - 1/k, b_i + 1/k[,$$

es claro que μ y m coinciden también sobre las celdas cerradas, y como toda celda contiene una celda abierta y está contenida en una celda cerrada de la misma medida, de hecho μ y m coinciden en todas las celdas. Esto implica trivialmente que coinciden sobre \mathcal{E}^n .

Sea $A \in \mathcal{B}^n$ y supongamos que $A \subset \bigcup_{k=0}^{\infty} A_k$, con $A_k \in \mathcal{E}^n$. Entonces la unión es un conjunto de Borel y

$$\mu(A) \leq \mu\left(\bigcup_{k=0}^{\infty} A_k\right) \leq \sum_{k=0}^{\infty} \mu(A_k) = \sum_{k=0}^{\infty} m(A_k).$$

Por definición de la medida exterior de Lebesgue $\mu(A) \leq m^*(A) = m(A)$. Sea $C_k = [-k, k]^n$ y sea $A_k = A \cap C_k$. Notemos que $\mu(A_k) \leq \mu(C_k) = (2k)^n < +\infty$. Por lo que acabamos de probar, $\mu(C_k \setminus A_k) \leq m(C_k \setminus A_k)$, luego

$$m(C_k) - \mu(A_k) \leq m(C_k) - m(A_k),$$

luego $m(A_k) \leq \mu(A_k)$ y la desigualdad contraria la tenemos probada en general, luego $\mu(A_k) = m(A_k)$. Por consiguiente,

$$\mu(A) = \lim_k \mu(A_k) = \lim_k m(A_k) = m(A). \quad \blacksquare$$

Tenemos así una caracterización sencilla de la medida de Lebesgue incompleta. A su vez, la medida de Lebesgue (completa) está determinada por ser la completación de ésta (por el teorema 8.30).

Veamos ahora que la medida de Lebesgue extiende, de hecho, a la medida de Jordan:

Teorema 8.35 *Todo conjunto medible Jordan es medible Lebesgue, y la medida de Lebesgue es la única extensión a \mathcal{M}^n de la medida de Jordan.*

DEMOSTRACIÓN: Si A es medible Jordan, el teorema 8.11 nos da figuras elementales $K_m \subset A \subset U_m$ de modo que K_m es compacta, U_m es abierta y $m(U_m \setminus K_m) < 1/m$. (Notemos que la medida de figuras elementales coincide con la medida de Lebesgue.) Sean $B = \bigcup_{m=1}^{\infty} K_m$, $C = \bigcap_{m=1}^{\infty} U_m$, que son conjuntos de Borel tales que $B \subset A \subset C$ y $m(C \setminus B) = 0$, luego A es medible Lebesgue, por la completitud de la medida de Lebesgue.

La restricción a \mathcal{J}^n de la medida de Lebesgue es una medida que extiende a la de las figuras elementales, luego por 8.12 coincide con la medida de Jordan.

Por último, si $\mu : \mathcal{M}^n \rightarrow [0, +\infty]$ es una medida que extiende a la de Jordan, entonces la restricción de μ a la σ -álgebra de Borel coincide con la medida de Lebesgue por el teorema anterior, pero la medida de Lebesgue es la única extensión a \mathcal{M}^n de la medida de Lebesgue incompleta, luego tiene que coincidir con μ . ■

De la unicidad extraemos esta consecuencia:

Teorema 8.36 *Sea $f : \mathbb{R}^n \rightarrow \mathbb{R}^n$ una biyección afín y $A \subset \mathbb{R}^n$ un conjunto de Borel (resp. medible Lebesgue). Entonces $f[A]$ es de Borel (resp. medible Lebesgue) y $m(f[A]) = |\det f| m(A)$. En particular, la medida de Lebesgue es invariante por isometrías.*

DEMOSTRACIÓN: Como f es un homeomorfismo, ya hemos observado que transforma conjuntos de Borel en conjunto de Borel. Sea $\mu : \mathcal{B}^n \rightarrow [0, +\infty]$ la aplicación dada por $\mu(A) = m(f[A])/|\det f|$. Es inmediato que se trata de una medida y por el teorema 8.17 su restricción a \mathcal{E}^n es m , luego el teorema anterior implica que $\mu = m$. Esto prueba el teorema para conjuntos de Borel.

Si A es medible Lebesgue, entonces $A = B \cup N$, con $N \subset C$, donde B y C son de Borel y C es nulo. Por lo tanto, $f[A] = f[B] \cup f[N]$, con $f[N] \subset f[C]$. Por la parte ya probada $f[B]$ y $f[C]$ son de Borel y $f[C]$ es nulo, luego $f[A]$ es medible Lebesgue y $m(f[A]) = m(f[B]) = |\det f| m(B) = |\det f| m(A)$. ■

Esto nos da otra caracterización de la medida de Lebesgue:

Teorema 8.37 *La única medida de Borel en \mathbb{R}^n invariante por traslaciones tal que $m([0, 1]^n) = 1$ es la medida de Lebesgue.*

DEMOSTRACIÓN: Por 8.8, la restricción de una tal medida a \mathcal{E}^n tiene que coincidir con la restricción de la medida de Lebesgue, luego el teorema 8.34 implica que es la medida de Lebesgue. ■

Terminamos la sección con algunas observaciones sobre los conjuntos nulos para la medida de Lebesgue. No es cierto en general que la imagen de un conjunto nulo por una aplicación continua sea un conjunto nulo, pero se cumple lo siguiente:

Teorema 8.38 *La imagen de un conjunto nulo por una función diferenciable entre dos abiertos de \mathbb{R}^n es un conjunto nulo.*

DEMOSTRACIÓN: Sea $g : U \rightarrow V$ diferenciable y sea $E \subset U$ un conjunto nulo. Si $x \in E$, la diferenciable de g en x significa que para $y \neq x$

$$g(y) - g(x) = \|y - x\| \left(dg(x) \left(\frac{y - x}{\|y - x\|} \right) + E(y - x) \right)$$

donde E es una cierta función continua en 0 con $E(0) = 0$. Como $dg(x)$ está acotada en la bola unidad, existen naturales k y p tales que

$$\text{si } y \in B_{1/p}(x) \quad \text{entonces} \quad \|g(y) - g(x)\| \leq k\|y - x\|.$$

Sea F_{kp} el conjunto de todos los $x \in E$ que cumplen esta relación. Hemos probado que E está contenido en la unión de estos conjuntos. Por consiguiente basta probar que $g[F_{kp}]$ es nulo.

Sea M el cociente entre la medida de una bola de radio 1 y la de un cubo $[a, b]^n$ de diámetro 1. Una homotecia de razón r los transforma en una bola de radio r y un cubo de diámetro r , cuyas medidas difieren de las anteriores en la constante r^n , luego la razón entre ambas sigue siendo M , es decir, que M es en realidad el cociente entre la medida de una bola de radio arbitrario r y la medida de un cubo de diámetro r .

Sea $\epsilon > 0$ y tomemos un abierto W tal que $F_{kp} \subset W$ y $m(W) < \epsilon/M$. Los cubos $[a, b]^n$ de extremos racionales y de diámetro menor que $1/p$ son una base numerable de \mathbb{R}^n , luego podemos expresar W como unión numerable de tales cubos. Desechamos los que no cortan a F_{kp} y así tenemos a éste cubierto por una familia numerable de cubos C_i cuyas medidas suman menos de ϵ/M . Cubrimos cada cubo por una bola $B_i = B_{r_i}(x_i)$ cuyo centro es un punto $x_i \in F_{kp} \cap C_i$ y su radio es el diámetro de C_i (menor que $1/p$). De este modo las bolas cubren a F_{kp} y la suma de sus medidas es menor que ϵ .

Si $x \in F_{kp} \cap B_i$, entonces $\|x - x_i\| < 1/p$ y $x_i \in F_{kp}$, luego

$$\|g(x) - g(x_i)\| \leq k\|x - x_i\| < nr_i,$$

luego $g[F_{kp} \cap B_i] \subset B_{kr_i}(g(x_i))$. Esto prueba que $g[F_{kp}]$ está cubierto por las bolas $B_{kr_i}(g(x_i))$, luego

$$m(g[F_{kp}]) \leq \sum_{i=1}^{\infty} m(B_{kr_i}(g(x_i))) = k^n \sum_{i=1}^{\infty} m(B_i) < k^n \epsilon.$$

Como ϵ es arbitrario, tenemos que $m(g[F_{kp}]) = 0$. ■

Como todo conjunto medible Lebesgue se puede expresar como unión de un conjunto de Borel con un conjunto nulo, el teorema anterior prueba que si $g : U \rightarrow V$ es un difeomorfismo entre dos abiertos de \mathbb{R}^n y $E \subset U$ es medible Lebesgue, entonces $g[E]$ es medible Lebesgue.

Los conjuntos nulos para la medida de Lebesgue y los conjuntos nulos para la medida de Jordan no son los mismos, pues un conjunto nulo para la medida de Jordan tiene clausura nula, mientras que $\mathbb{Q} \cap [0, 1]$ es nulo para la medida de Lebesgue y su clausura no es nula. Ahora, bien, es fácil ver que todo conjunto nulo para la medida de Jordan lo es para la de Lebesgue (si se puede cubrir con una cantidad finita de celdas de medida arbitrariamente pequeña, también se puede cubrir por una cantidad numerable de celdas de medida arbitrariamente pequeña), y que todo compacto nulo para la medida de Lebesgue lo es para la de Jordan (porque puede cubrirse por una cantidad numerable de celdas abiertas de medida arbitrariamente pequeña, y por compacidad podemos tomar un subcubrimiento finito). Como la frontera de un conjunto acotado es compacta, el teorema 8.14 puede reformularse así:

Teorema 8.39 *Un conjunto acotado $A \subset \mathbb{R}^n$ es medible Jordan si y sólo si ∂A es nula para la medida de Lebesgue.*

Esto nos permite identificar los conjuntos medibles Jordan en términos de la medida de Lebesgue, que es más práctica para hacer cálculos.

8.4 Funciones medibles

Si X es un espacio medido, a menudo tendremos que trabajar con subconjuntos de X definidos a partir de una aplicación $f : X \rightarrow Y$ y necesitaremos garantizar que dichos conjuntos son medibles. Esto lo lograremos mediante el concepto de aplicación medible.

Definición 8.40 Si X es un espacio medido e Y es un espacio topológico, una aplicación $f : X \rightarrow Y$ es *medible* si las antiimágenes por f de los abiertos de Y son conjuntos medibles.

Como las antiimágenes conservan las operaciones conjuntistas es muy fácil probar que los conjuntos de Y cuyas antiimágenes son medibles forman una σ -álgebra, que en el caso de una función medible contiene a los abiertos, luego contendrá a todos los conjuntos de Borel, es decir, una aplicación es medible si y sólo si las antiimágenes de los conjuntos de Borel son conjuntos medibles. El siguiente caso particular nos interesará especialmente:

Teorema 8.41 Una aplicación $f : X \rightarrow [-\infty, +\infty]$ es medible si y sólo si lo son todos los conjuntos $f^{-1}[x, +\infty]$, para todo $x \in \mathbb{R}$.

DEMOSTRACIÓN: Ya hemos comentado que los conjuntos con antiimagen medible forman una σ -álgebra \mathcal{A} . Hemos de ver que \mathcal{A} contiene a los abiertos de $[-\infty, +\infty]$. Por hipótesis contiene a los intervalos $]x, +\infty]$, luego también a sus complementarios $[-\infty, x]$. Todo intervalo $[-\infty, x[$ es intersección numerable de los intervalos $[-\infty, x + 1/n]$, luego también está en \mathcal{A} . De aquí se sigue que \mathcal{A} contiene también a los intervalos $]x, y[=]x, +\infty] \cap [-\infty, y[$. Finalmente, todo abierto de $[-\infty, +\infty]$ se expresa como unión numerable de intervalos abiertos, luego está en \mathcal{A} . ■

Es claro que la composición de una función medible con una función continua es una función medible. Esto nos da, por ejemplo, que si $f : X \rightarrow [-\infty, +\infty]$ es medible, también lo es $|f|$ y αf para todo número real α , así como $1/f$ si f no se anula.

Para probar resultados análogos cuando intervienen dos funciones (suma de funciones medibles, etc.) usaremos la observación siguiente:

Si los espacios topológicos Y, Z tienen bases numerables (como $[-\infty, +\infty]$ y sus subespacios) y $u : X \rightarrow Y, v : X \rightarrow Z$ son aplicaciones medibles, entonces la aplicación $u \times v : X \rightarrow Y \times Z$ dada por $(u \times v)(x) = (u(x), v(x))$ es medible.

Basta observar que los productos de abiertos básicos $A \times B$ forman una base numerable de $Y \times Z$, luego todo abierto de $Y \times Z$ es unión numerable de estos conjuntos, por lo que es suficiente que sus antiimágenes sean medibles, pero $(u \times v)^{-1}[A \times B] = u^{-1}[A] \cap v^{-1}[B]$.

Ahora, por ejemplo, si $u, v : X \rightarrow \mathbb{R}$ son aplicaciones medibles, también lo son $u + v$ y uv , pues son la composición de $u \times v$ con la suma y el producto, que son continuas.

Nos interesa extender este resultado a funciones $u, v : X \rightarrow [-\infty, +\infty]$, pero entonces tenemos el problema de que no es posible extender la suma y el producto de modo que sean continuas en los puntos $(+\infty, -\infty)$, $(-\infty, +\infty)$ en el caso de la suma y en los puntos $(\pm\infty, 0)$, $(0, \pm\infty)$ en el caso del producto.

Hacemos esto: definimos $u + v$ de modo que $+\infty - \infty = 0$ (por ejemplo) y ahora observamos lo siguiente:

Sea $u : X \rightarrow Y$ una función medible, A un subconjunto medible de X e $y \in Y$. Entonces la función $v : X \rightarrow Y$ que coincide con u fuera de A y toma el valor y en A es medible.

La razón es que

$$v^{-1}[B] = \begin{cases} u^{-1}[B] & \text{si } y \notin B \\ u^{-1}[B] \cup A & \text{si } y \in B \end{cases}$$

Así, dadas $u, v : X \rightarrow [-\infty, +\infty]$ medibles tales que donde una vale $+\infty$ la otra no vale $-\infty$, las modificamos para que valgan 0 donde toman valores infinitos, las sumamos y obtenemos una función medible, luego modificamos la suma para que tome el valor ∞ adecuado donde deba tomar dichos valores (claramente en un conjunto medible), con lo que obtenemos una función medible. Igualmente con el producto.

Otra consecuencia del teorema sobre el producto cartesiano de funciones medibles es que si tenemos dos funciones medibles $u, v : X \rightarrow [-\infty, +\infty]$, los conjuntos del estilo de $\{x \in X \mid u(x) < v(x)\}$ son medibles (por ejemplo en este caso se trata de la antiimagen por $u \times v$ del abierto $\{(x, y) \mid x < y\}$).

Dos operaciones definibles sobre todas las funciones $f, g : X \rightarrow [-\infty, +\infty]$ son las dadas por $(f \vee g)(x) = \max\{f(x), g(x)\}$ y $(f \wedge g)(x) = \min\{f(x), g(x)\}$.

Las funciones $\vee, \wedge : [-\infty, +\infty] \times [-\infty, +\infty] \rightarrow [-\infty, +\infty]$ son ambas continuas, pues lo son trivialmente cuando se restringen a los cerrados determinados por las condiciones $x \leq y$ e $y \leq x$, respectivamente. Esto implica que si $f, g : X \rightarrow [-\infty, +\infty]$ son medibles, también lo son $f \vee g$ y $f \wedge g$.

En particular si $f : X \rightarrow [-\infty, +\infty]$ es una función medible, definimos las funciones $f^+ = f \vee 0$ y $f^- = -(f \wedge 0)$, llamadas *parte positiva* y *parte negativa* de f , respectivamente.

Tenemos que si f es medible también lo son f^+ y f^- . El recíproco es cierto porque claramente $f = f^+ - f^-$. Además $|f| = f^+ + f^-$.

Seguidamente probaremos que la medibilidad se conserva al tomar límites. Si $\{a_n\}_{n=0}^{\infty}$ es una sucesión en $[-\infty, +\infty]$, definimos sus límites superior e inferior como

$$\overline{\lim}_n a_n = \inf_{k \geq 0} \sup_{n \geq k} a_n, \quad \underline{\lim}_n a_n = \sup_{k \geq 0} \inf_{n \geq k} a_n.$$

En el capítulo IV demostrábamos que el límite superior de una sucesión es el supremo de sus puntos adherentes. Igualmente se prueba que el límite inferior es el ínfimo de sus puntos adherentes.

Una sucesión converge si y sólo si tiene un único punto adherente (su límite), por lo que $\{a_n\}_{n=0}^\infty$ converge si y sólo si $\overline{\lim}_n a_n = \underline{\lim}_n a_n$ y entonces

$$\overline{\lim}_n a_n = \underline{\lim}_n a_n = \lim_n a_n.$$

Si $\{f_n\}_{n=0}^\infty$ es una sucesión de funciones $f_n : X \rightarrow [-\infty, +\infty]$, definimos puntualmente las funciones $\sup_n f_n$, $\inf_n f_n$, $\underline{\lim}_n f_n$ y $\overline{\lim}_n f_n$. Si la sucesión es puntualmente convergente las dos últimas funciones coinciden con la función límite puntual $\lim_n f_n$.

Teorema 8.42 *Si las funciones f_n son medibles, también lo son las funciones $\sup_n f_n$, $\inf_n f_n$, $\overline{\lim}_n f_n$ y $\underline{\lim}_n f_n$.*

DEMOSTRACIÓN: Claramente

$$\left(\sup_n f_n\right)^{-1} [x, +\infty] = \bigcup_{n=0}^{\infty} f_n^{-1} [x, +\infty],$$

es medible. Igualmente se prueba con ínfimos y de aquí se deducen los resultados sobre límites superiores e inferiores. En particular, el límite puntual de una sucesión de funciones medibles es una función medible. ■

Si X es un espacio medida y E es un subconjunto medible, entonces los subconjuntos medibles de E forman una σ -álgebra de subconjuntos de E y la medida de X restringida a esta σ -álgebra es una medida en E . En lo sucesivo consideraremos a todos los subconjuntos medibles de los espacios medida como espacios medida de esta manera.

Notar que si $X = \bigcup_{n=0}^{\infty} E_n$ es una descomposición de X en subconjuntos medibles (no necesariamente disjuntos), entonces $f : X \rightarrow Y$ es medible si y sólo si lo son todas las funciones $f|_{E_n}$, pues si f es medible y G es un abierto en Y , $(f|_{E_n})^{-1}[G] = f^{-1}[G] \cap E_n$, luego $(f|_{E_n})^{-1}[G]$ es medible, y si las $f|_{E_n}$ son medibles, entonces $f^{-1}[G] = (f|_{E_n})^{-1}[G]$, luego también es medible.

Si E es un subconjunto de X , su función característica χ_E es medible si y sólo si lo es E .

Si extendemos una función medible $f : E \rightarrow [-\infty, +\infty]$ asignándole el valor 0 fuera de E , obtenemos una función medible en X . Por ello identificaremos las funciones medibles $f : E \rightarrow [-\infty, +\infty]$ con las funciones medibles en X que se anulan fuera de E . En particular identificaremos la restricción a E de una función $f : X \rightarrow [-\infty, +\infty]$ con la función $f\chi_E$.

Los resultados que hemos dado son suficientes para garantizar que todas las funciones que manejaremos y los conjuntos definidos por ellas son medibles. No insistiremos en ello a menos que haya alguna dificultad inusual.

8.5 La integral de Lebesgue

Vamos a probar que en todo espacio medida X podemos definir una integral, de modo que todas las integrales que hemos calculado hasta ahora resultarán ser casos particulares de la integral asociada a la medida de Lebesgue en \mathbb{R} .

Definición 8.43 Una función *simple* en un espacio medida X es una función medible $s : X \rightarrow [0, +\infty[$ que sólo toma un número finito de valores $\alpha_1, \dots, \alpha_n$. Si llamamos $A_i = s^{-1}[\alpha_i]$, entonces los conjuntos A_i son medibles disjuntos y $s = \sum_{i=1}^n \alpha_i \chi_{A_i}$.

La base de nuestra construcción de la integral será el teorema siguiente:

Teorema 8.44 Si X es un espacio medida y $f : X \rightarrow [0, +\infty]$ es una función medible, entonces existe una sucesión $\{s_n\}_{n=1}^\infty$ de funciones simples en X tal que

$$0 \leq s_1 \leq s_2 \leq \dots \leq f \quad \text{y} \quad f = \lim_n s_n.$$

DEMOSTRACIÓN: Para cada número natural $n > 0$ y cada $t \in \mathbb{R}$ existe un único $k = k_n(t) \in \mathbb{N}$ tal que $k/2^n \leq t < (k+1)/2^n$. Sea $f_n : [0, +\infty] \rightarrow [0, +\infty]$ dada por

$$f_n(t) = \begin{cases} k_n(t)/2^n & \text{si } 0 \leq t < n \\ n & \text{si } n \leq t \leq +\infty \end{cases}$$

La figura muestra la función f_2 . Claramente f_n toma un número finito de valores y

$$f_1 \leq f_2 \leq f_3 \leq \dots \leq I,$$

donde I es la función identidad $I(t) = t$. Como $t - 1/2^n < f_n(t) \leq t$ para $0 \leq t \leq n$, es claro que $\{f_n\}_{n=1}^\infty$ converge puntualmente a I .

Sea $s_n = f \circ f_n$. Claramente s_n toma un número finito de valores (a lo sumo los que toma f_n) y las antiimágenes de estos valores son las antiimágenes por f de los intervalos donde los toma f_n , luego son conjuntos medibles. Además

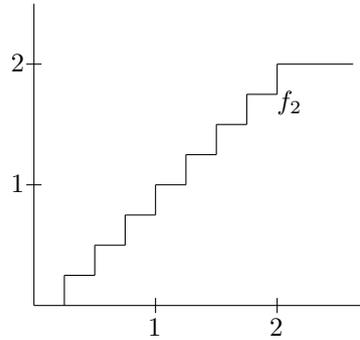
$$0 \leq s_1 \leq s_2 \leq \dots \leq f,$$

luego las funciones s_n son simples y, tomando límites, es obvio que la sucesión converge puntualmente a f . ■

Definición 8.45 Sea X un espacio medida y $s = \sum_{i=1}^n \alpha_i \chi_{A_i}$ una función simple en X . Definimos la *integral* de s en X como

$$\int_X s \, d\mu = \sum_{i=1}^n \alpha_i \mu(A_i) \in [0, +\infty],$$

con el convenio de que $+\infty \cdot 0 = 0$.



Si E es un subconjunto medible de X , entonces $s|_E = \sum_{i=1}^n \alpha_i \chi_{A_i \cap E}$, donde las funciones características se toman ahora sobre E . Por lo tanto

$$\int_E s|_E d\mu = \sum_{i=1}^n \alpha_i \mu(A_i \cap E).$$

Por otro lado $s\chi_E = \sum_{i=1}^n \alpha_i \chi_{A_i \cap E}$ (con las funciones características en X), luego concluimos que

$$\int_E s d\mu = \int_X s\chi_E d\mu = \sum_{i=1}^n \alpha_i \mu(A_i \cap E),$$

es decir, que a efectos de integración podemos adoptar consistentemente el convenio explicado antes por el que identificamos la función $s|_E$ con $s\chi_E$.

Ahora necesitamos el siguiente resultado técnico, que después generalizaremos notablemente.

Teorema 8.46 *Sea X un espacio medida.*

a) *Sea s una función simple en X . Para cada subconjunto medible E de X definimos $\nu(E) = \int_E s d\mu$. Entonces ν es una medida en X .*

b) *Si s y t son funciones simples en X se cumple*

$$\int_X (s+t) d\mu = \int_X s d\mu + \int_X t d\mu.$$

DEMOSTRACIÓN: a) Sea $s = \sum_{i=1}^n \alpha_i \chi_{A_i}$. Claramente $\nu(\emptyset) = 0$. Sea $E = \bigcup_{j=1}^{\infty} E_j$ una unión disjunta de conjuntos medibles. Entonces

$$\begin{aligned} \nu(E) &= \sum_{i=1}^n \alpha_i \mu(A_i \cap E) = \sum_{i=1}^n \alpha_i \sum_{j=1}^{\infty} \mu(A_i \cap E_j) \\ &= \sum_{j=1}^{\infty} \sum_{i=1}^n \alpha_i \mu(A_i \cap E_j) = \sum_{j=1}^{\infty} \nu(E_j). \end{aligned}$$

b) Sean $s = \sum_{i=1}^n \alpha_i \chi_{A_i}$ y $t = \sum_{j=1}^m \beta_j \chi_{B_j}$. Llamemos $E_{ij} = A_i \cap B_j$. Así, tanto s como t son constantes en los conjuntos E_{ij} (s toma el valor α_i y t el valor β_j). Por lo tanto

$$\int_{E_{ij}} (s+t) d\mu = (\alpha_i + \beta_j) \mu(E_{ij}) = \alpha_i \mu(E_{ij}) + \beta_j \mu(E_{ij}) = \int_{E_{ij}} s d\mu + \int_{E_{ij}} t d\mu.$$

Como los conjuntos E_{ij} son disjuntos dos a dos y su unión es X , la parte a) nos da que la igualdad se cumple para integrales en X . ■

En particular notamos que si $s \leq t$ son funciones simples en un espacio medido X , entonces $t - s$ también es una función simple y

$$\int_X s \, d\mu \leq \int_X s \, d\mu + \int_X (t - s) \, d\mu = \int_X t \, d\mu.$$

En particular se cumple que

$$\int_X t \, d\mu = \sup \left\{ \int_X s \, d\mu \mid s \text{ es una función simple, } s \leq t \right\}.$$

Esto hace consistente la siguiente definición:

Definición 8.47 Sea X un espacio medido y $f : X \rightarrow [0, +\infty]$ una función medible. Definimos la *integral* de f como

$$\int_X f \, d\mu = \sup \left\{ \int_X s \, d\mu \mid s \text{ es una función simple, } s \leq f \right\} \in [0, +\infty].$$

Observar que si E es un subconjunto medible de X , si s es una función simple en E por debajo de $f|_E$, su extensión a X (nula fuera de E) es una función simple bajo $f\chi_E$, y la restricción a E de una función simple en X bajo $f\chi_E$ es una función simple en E bajo $f|_E$. De aquí se sigue que

$$\int_E f \, d\mu = \int_X f\chi_E \, d\mu,$$

pues ambas integrales son el supremo del mismo conjunto de números reales.

Las propiedades siguientes son inmediatas a partir de la definición:

Teorema 8.48 Sea X un espacio medido y E un subconjunto medible de X .

- Si $0 \leq f \leq g$ son funciones medibles en X , entonces $\int_X f \, d\mu \leq \int_X g \, d\mu$.
- Si $f \geq 0$ es una función medible en X y $A \subset B$ son subconjuntos medibles de X , entonces $\int_A f \, d\mu \leq \int_B f \, d\mu$.
- Si $f \geq 0$ es una función medible en X y $f|_E = 0$, entonces $\int_E f \, d\mu = 0$ (aunque sea $\mu(E) = +\infty$).
- Si $f \geq 0$ es una función medible en X y $\mu(E) = 0$, entonces $\int_E f \, d\mu = 0$ (aunque sea $f|_E = +\infty$).

El resultado siguiente es uno de los más importantes del cálculo integral:

Teorema 8.49 (de la convergencia monótona de Lebesgue) Sea X un espacio medido y $\{f_n\}_{n=1}^\infty$ una sucesión de funciones medibles en X tal que

$$0 \leq f_1 \leq f_2 \leq \cdots \leq f \quad \text{y} \quad f = \lim_n f_n.$$

Entonces f es medible y

$$\int_X f \, d\mu = \lim_n \int_X f_n \, d\mu.$$

DEMOSTRACIÓN: Por el teorema anterior $\int_X f_n d\mu \leq \int_X f_{n+1} d\mu$. Toda sucesión monótona creciente en $[0, +\infty]$ converge a su supremo, luego existe $\alpha = \lim_n \int_X f_n d\mu \in [0, +\infty]$.

Sabemos que f es medible por ser límite puntual de funciones medibles. De nuevo por el teorema anterior $\int_X f_n d\mu \leq \int_X f d\mu$, luego $\alpha \leq \int_X f d\mu$.

Sea s una función simple $s \leq f$ y sea $0 < c < 1$. Definimos

$$E_n = \{x \in X \mid f_n(x) \geq cs(x)\}, \quad \text{para } n = 1, 2, 3, \dots$$

Claramente $E_1 \subset E_2 \subset E_3 \subset \dots$, son conjuntos medibles y, según veremos enseguida, $X = \bigcup_n E_n$.

En efecto, si $x \in X$ y $f(x) = 0$, entonces $x \in E_1$ y si, por el contrario, $f(x) > 0$ entonces $cs(x) < s(x) \leq f(x)$, luego $x \in E_n$ para algún n . Claramente

$$\int_X f_n d\mu \geq \int_{E_n} f_n d\mu \geq c \int_{E_n} s d\mu.$$

Ahora aplicamos el teorema 8.46 y el hecho de que la medida de la unión de una sucesión creciente de conjuntos es el supremo de las medidas, con lo que obtenemos

$$\alpha = \lim_n \int_X f_n d\mu \geq c \lim_n \int_{E_n} s d\mu = c \int_X s d\mu.$$

Como esto es cierto para todo $c < 1$ podemos concluir que $\alpha \geq \int_X s d\mu$ para toda función simple $s \leq f$, luego tomando el supremo de estas integrales resulta $\alpha \geq \int_X f d\mu$, con lo que tenemos la igualdad buscada. ■

El teorema de la convergencia monótona permite en particular reducir propiedades de la integral de funciones no negativas a propiedades de funciones simples, como ilustra el teorema siguiente, que muestra que la integral conserva las sumas, incluso las infinitas.

Teorema 8.50 *Sea X un espacio medida y sea $\{f_n\}_{n=1}^\infty$ una sucesión de funciones no negativas medibles en X . Entonces*

$$\int_X \sum_{n=1}^\infty f_n d\mu = \sum_{n=1}^\infty \int_X f_n d\mu.$$

DEMOSTRACIÓN: Probaremos en primer lugar que si f y g son medibles y no negativas, entonces

$$\int_X (f + g) d\mu = \int_X f d\mu + \int_X g d\mu.$$

Tomamos dos sucesiones monótonas $\{s_n\}_{n=1}^\infty$ y $\{t_n\}_{n=1}^\infty$ de funciones simples convergentes a f y g respectivamente (existen por el teorema 8.44).

Por el teorema 8.46 sabemos que $\int_X (s_n + t_n) d\mu = \int_X s_n d\mu + \int_X t_n d\mu$ y, tomando límites, el teorema de la convergencia monótona nos da la igualdad buscada. En el caso general sabemos, por lo que acabamos de probar, que

$$\int_X \sum_{n=1}^k f_n d\mu = \sum_{n=1}^k \int_X f_n d\mu \quad \text{para } k = 1, 2, 3, \dots$$

Las funciones $\sum_{n=1}^k f_n$ forman una sucesión monótona de funciones medibles, luego por el teorema de la convergencia monótona

$$\int_X \sum_{n=1}^{\infty} f_n d\mu = \sum_{n=1}^{\infty} \int_X f_n d\mu. \quad \blacksquare$$

Generalizamos ahora la primera parte del teorema 8.46.

Teorema 8.51 *Sea X un espacio medida y $f : X \rightarrow [0, +\infty]$ una función medible. Para cada subconjunto medible E de X definimos $\nu(E) = \int_E f d\mu$. Entonces ν es una medida en X .*

DEMOSTRACIÓN: Claramente $\nu(\emptyset) = 0$. Sea $E = \bigcup_{n=1}^{\infty} E_n$ una unión disjunta de conjuntos medibles. Es claro que $f\chi_E = \sum_{n=1}^{\infty} f\chi_{E_n}$. Aplicando el teorema anterior queda $\nu(E) = \sum_{n=1}^{\infty} \nu(E_n)$. \blacksquare

Después necesitaremos el hecho siguiente:

Teorema 8.52 (Lema de Fatou) *Sea X un espacio medida y sea $\{f_n\}_{n=1}^{\infty}$ una sucesión de funciones medibles no negativas en X . Entonces*

$$\int_X \liminf_n f_n d\mu \leq \liminf_n \int_X f_n d\mu.$$

DEMOSTRACIÓN: Sea $g_k = \inf_{n \geq k} f_n$. Entonces $g_k \leq f_n$ para $n \geq k$, luego $\int_X g_k d\mu \leq \inf_{n \geq k} \int_X f_n d\mu$. Además las funciones g_k forman una sucesión monótona creciente que converge a $\liminf_n f_n$, luego por el teorema de la convergencia monótona

$$\int_X \liminf_n f_n d\mu = \lim_k \int_X g_k d\mu = \sup_{k \geq 1} \int_X g_k d\mu \leq \sup_{k \geq 1} \inf_{n \geq k} \int_X f_n d\mu = \liminf_n \int_X f_n d\mu. \quad \blacksquare$$

Ahora extendemos la integral a funciones medibles no necesariamente mayores o iguales que 0.

Definición 8.53 Sea X un espacio medido y $f : X \rightarrow [-\infty, +\infty]$ una función medible. Entonces f^+ y f^- son funciones medibles no negativas y $f = f^+ - f^-$. Diremos que f es *integrable Lebesgue* en X si tanto $\int_X f^+ d\mu$ como $\int_X f^- d\mu$ son finitas. En tal caso definimos la *integral de Lebesgue* de f como

$$\int_X f d\mu = \int_X f^+ d\mu - \int_X f^- d\mu \in \mathbb{R}.$$

Llamaremos $L^1(\mu)$ al conjunto de las funciones integrables Lebesgue en X respecto a la medida μ .

Si una función f es no negativa, entonces $f^- = 0$ y su integral es la que ya teníamos definida. Las propiedades siguientes son todas inmediatas a partir de los resultados que ya hemos demostrado.

Teorema 8.54 Sea X un espacio medido y sean $f, g : X \rightarrow [-\infty, +\infty]$ funciones medibles.

a) f es integrable si y sólo si $\int_X |f| d\mu < +\infty$, y en tal caso

$$\left| \int_X f d\mu \right| \leq \int_X |f| d\mu.$$

b) Si $\alpha, \beta \in \mathbb{R}$, y f, g son integrables, entonces $\alpha f + \beta g$ es integrable y

$$\int_X (\alpha f + \beta g) d\mu = \alpha \int_X f d\mu + \beta \int_X g d\mu.$$

c) Si $f \leq g$ y ambas son integrables, entonces $\int_X f d\mu \leq \int_X g d\mu$.

d) Si E es un subconjunto medible de X y f es integrable en X , entonces f es integrable en E y $\int_E f d\mu = \int_X f \chi_E d\mu$. En particular, $\mu(E) = \int_E 1 d\mu$.

e) Si E y F son subconjuntos medibles disjuntos de X , entonces la función f es integrable en $E \cup F$ si y sólo si lo es en E y en F y, en tal caso,

$$\int_{E \cup F} f d\mu = \int_E f d\mu + \int_F f d\mu.$$

f) Si E es un subconjunto medible de X y $f|_E = 0$, entonces $\int_E f d\mu = 0$.

g) Si E es un subconjunto nulo de X , entonces $\int_E f d\mu = 0$.

h) Si f es integrable en X , entonces el conjunto de los puntos donde f toma los valores $\pm\infty$ es nulo.

i) Si $|f| \leq g$ y g es integrable entonces f también lo es.

j) Toda función medible y acotada sobre un conjunto de medida finita es integrable.

La propiedad e) sale de aplicar el teorema 8.51 a las partes positiva y negativa de f , la propiedad i) se deduce de a) y j) se deduce de i).

Otra observación de interés es la siguiente: si pasamos de una medida a su completión, es claro que las funciones simples para la primera lo son también para la segunda y las integrales coinciden. El teorema de la convergencia monótona implica entonces que toda función positiva integrable para una medida sigue siéndolo para su completión, y de aquí se sigue inmediatamente el resultado para funciones arbitrarias. De este modo la integral respecto a la completión extiende a la integral respecto a la medida de partida.

Veamos ahora un teorema de convergencia válido para funciones medibles arbitrarias.

Teorema 8.55 (de la convergencia dominada de Lebesgue) *Sea X un espacio medida y sean $\{f_n\}_{n=1}^{\infty}$ funciones medibles de X en $[-\infty, +\infty]$ que convergen puntualmente a una función f . Si existe una función integrable $g : X \rightarrow [-\infty, +\infty]$ tal que $|f_n| \leq g$ para todo n , entonces f es integrable y*

$$\int_X f d\mu = \lim_n \int_X f_n d\mu.$$

Se dice que las funciones f_n están dominadas por g .

DEMOSTRACIÓN: Claramente $|f| \leq g$, luego f es integrable. Puesto que $|f_n - f| \leq 2g$, podemos aplicar el lema de Fatou a las funciones no negativas $2g - |f_n - f|$, con lo que obtenemos que

$$\int_X 2g d\mu \leq \liminf_n \int_X (2g - |f_n - f|) d\mu = \int_X 2g d\mu + \liminf_n \int_X (-|f_n - f|) d\mu.$$

Es fácil ver que el signo negativo sale del límite, pero cambiando éste por un límite superior, así $-\liminf_n \int_X |f_n - f| d\mu \geq 0$, o sea, $\liminf_n \int_X |f_n - f| d\mu \leq 0$.

Pero es obvio que $0 \leq \liminf_n \int_X |f_n - f| d\mu \leq \overline{\lim}_n \int_X |f_n - f| d\mu = 0$, luego los límites superior e inferior coinciden, luego $\lim_n \int_X |f_n - f| d\mu = 0$. Ahora aplicamos que

$$\left| \int_X f_n d\mu - \int_X f d\mu \right| = \left| \int_X (f_n - f) d\mu \right| \leq \int_X |f_n - f| d\mu,$$

de donde se sigue el teorema. ■

Cuando una propiedad se verifica para todos los puntos de un espacio medida salvo los de un conjunto nulo se dice que la propiedad se verifica *para casi todo punto*, y lo abreviaremos p.c.t.p. Veamos un ejemplo:

Teorema 8.56 *Si X es un espacio medida y $f : X \rightarrow [0, +\infty]$ es una función medible tal que $\int_X f d\mu = 0$, entonces $f = 0$ p.c.t.p. de X .*

DEMOSTRACIÓN: Para cada natural $n > 0$ sea $E_n = \{x \in E \mid f(x) > 1/n\}$. Entonces

$$\frac{1}{n} \mu(E_n) \leq \int_{E_n} f d\mu \leq \int_X f d\mu = 0,$$

luego $\mu(E_n) = 0$. La unión de los E_n es el conjunto $E = \{x \in X \mid f(x) > 0\}$, luego $f = 0$ salvo en los puntos del conjunto nulo E . ■

Cuando digamos que una función $f : X \rightarrow [-\infty, +\infty]$ está definida p.c.t.p. esto significará que en realidad es $f : X \setminus E \rightarrow [-\infty, +\infty]$, donde E es un conjunto nulo. Diremos que f es medible si lo es al extenderla a X tomando el valor 0 en E . En tal caso podemos hablar de $\int_X f d\mu$ definida como la integral de dicha extensión.

Terminamos la sección con un importante teorema sobre integrales paramétricas:

Teorema 8.57 *Sea U abierto en \mathbb{R}^n , K un espacio métrico compacto, μ una medida de Borel finita en K , sean $f : U \times K \rightarrow \mathbb{R}$ una función continua y $g : K \rightarrow \mathbb{R}$ una función medible acotada. Definamos $F : U \rightarrow \mathbb{R}^n$ como la función dada por*

$$F(x) = \int_K f(x, y)g(y) d\mu(y),$$

donde $d\mu(y)$ indica que la integral se realiza respecto a la variable $y \in K$, considerando constante a x . Entonces F es continua en U y si existe

$$\frac{\partial f}{\partial x_i} : U \times K \rightarrow \mathbb{R}$$

y es continua en $U \times K$, entonces existe

$$\frac{\partial F}{\partial x_i} = \int_K \frac{\partial f}{\partial x_i}(x, y)g(y) d\mu(y)$$

y es continua en U .

DEMOSTRACIÓN: Tomemos $x_0 \in U$ y sea B una bola cerrada de centro x_0 contenida en U . Sea M una cota de g en K . Como f es uniformemente continua en $B \times K$, dado $\epsilon > 0$ existe un $\delta > 0$ tal que si $\|x - x_0\| < \delta$, entonces $|f(x, y) - f(x_0, y)| < \epsilon/M\mu(K)$, para todo $y \in K$. Por consiguiente, si $\|x - x_0\| < \delta$ se cumple

$$|F(x) - F(x_0)| \leq \int_K |f(x, y) - f(x_0, y)| |g(y)| d\mu(y) \leq \epsilon.$$

Esto prueba que F es continua en x_0 .

Supongamos ahora la hipótesis de derivabilidad respecto a x_i y sea e_i el i -ésimo vector de la base canónica de \mathbb{R}^n . Como $\partial f/\partial x_i$ es uniformemente continua en $B \times K$, existe un $\delta > 0$ tal que si $|h| < \delta$ entonces

$$\left| \frac{\partial f}{\partial x_i}(x_0 + he_i, y) - \frac{\partial f}{\partial x_i}(x_0, y) \right| < \frac{\epsilon}{M\mu(K)}, \quad \text{para todo } y \in K.$$

Si $|h| < \delta$ e $y \in K$, el teorema del valor medio nos da que existe un $r \in \mathbb{R}$ tal que $|r| < |h|$ y

$$f(x_0 + he_i, y) - f(x_0, y) = \frac{\partial f}{\partial x_i}(x_0 + re_i, y)h.$$

(Notemos que r depende de y .)

Por consiguiente,

$$\begin{aligned} & \left| \frac{f(x_0 + he_i, y)g(y) - f(x_0, y)g(y)}{h} - \frac{\partial f}{\partial x_i}(x_0, y)g(y) \right| \\ &= \left| \frac{\partial f}{\partial x_i}(x_0 + re_i, y) - \frac{\partial f}{\partial x_i}(x_0, y) \right| |g(y)| < \frac{\epsilon}{\mu(K)}. \end{aligned}$$

De aquí se sigue claramente que

$$\left| \frac{F(x_0 + he_i) - F(x_0)}{h} - \int_K \frac{\partial f}{\partial x_i}(x_0, y)g(y) d\mu(y) \right| < \epsilon.$$

siempre que $|h| < \delta$, luego existe

$$\frac{\partial F}{\partial x_i}(x_0) = \lim_{h \rightarrow 0} \frac{F(x_0 + he_i) - F(x_0)}{h} = \int_K \frac{\partial f}{\partial x_i}(x_0, y)g(y) d\mu(y).$$

Además la derivada es continua por la primera parte de este mismo teorema. \blacksquare

8.6 La integral de Lebesgue en \mathbb{R}

Todas las integrales que hemos calculado hasta ahora en la práctica han sido integrales de funciones continuas en un intervalo, y lo hemos hecho bajo el supuesto de que tenían primitiva, y tomando como definición de integral que

$$\int_a^b f(x) dx = F(b) - F(a),$$

donde $F : [a, b] \rightarrow \mathbb{R}$ es una función continua derivable en $]a, b[$ con derivada f . Vamos a probar que la integral de una función continua en este sentido no es más que la integral respecto de la medida de Lebesgue.

Observemos en primer lugar que toda función continua es medible Lebesgue (de hecho, una función continua en un espacio topológico X es medible para cualquier medida de Borel en X) y, como la medida de Lebesgue es finita sobre los intervalos cerrados, toda función continua f en un intervalo $[a, b]$ es integrable Lebesgue en $[a, b]$. En particular, si f es continua en un intervalo I y $a \in I$, podemos definir en I la función

$$F(x) = \int_a^x f(t) dt,$$

donde ahora el miembro derecho representa la integral de Lebesgue en \mathbb{R} .

Teorema 8.58 *Sea f una función continua en un intervalo I y sea $a \in I$. Entonces la función*

$$F(x) = \int_a^x f(t) dt$$

es derivable en el interior de I y $F' = f$.

DEMOSTRACIÓN: Sea x un punto interior de I y sea $J \subset I$ un intervalo cerrado y acotado que contenga a a y a x (a éste último en su interior). Por el teorema 3.34 la función f es uniformemente continua en J , luego para cada $\epsilon > 0$ existe un $\delta > 0$ tal que si $u, u' \in J$, $|u - u'| < \delta$ entonces $|f(u) - f(u')| < \epsilon$.

Sea $h \in \mathbb{R}$ tal que $|h| < \delta$ y $x + h \in J$. Sean m y M el mínimo y el máximo de f en el intervalo cerrado de extremos x y $x + h$. Si $h > 0$

$$mh = \int_x^{x+h} m \, dt \leq \int_x^{x+h} f(t) \, dt \leq \int_x^{x+h} M \, dt = Mh.$$

Si $h < 0$ se invierten las desigualdades, pero en ambos casos resulta

$$m \leq \frac{\int_x^{x+h} f(t) \, dt}{h} \leq M.$$

Por el teorema de los valores intermedios existe un α entre x y $x + h$ de modo que

$$f(\alpha) = \frac{\int_x^{x+h} f(t) \, dt}{h} = \frac{F(x+h) - F(x)}{h}.$$

Claramente $|\alpha - x| < |h| < \delta$, luego

$$\left| \frac{F(x+h) - F(x)}{h} - f(x) \right| = |f(\alpha) - f(x)| < \epsilon,$$

por lo que existe $F'(x) = f(x)$. ■

Como consecuencia obtenemos:

Teorema 8.59 (Regla de Barrow) *Si f es una función continua en un intervalo $[a, b]$ entonces F tiene una primitiva F en $[a, b]$ y*

$$\int_a^b f(x) \, dx = F(b) - F(a).$$

DEMOSTRACIÓN: Cuando decimos que F es una primitiva de f en $[a, b]$ queremos decir que F es continua en $[a, b]$, derivable en $]a, b[$ y $F' = f$ en $]a, b[$. Basta tomar como F la función

$$F(x) = \int_a^x f(t) \, dt.$$

Sólo falta probar que F es continua en $[a, b]$, lo cual es sencillo, pues si M es una cota de f en $[a, b]$, entonces

$$|F(y) - F(x)| = \left| \int_x^y f(t) \, dt \right| \leq \int_x^y |f(t)| \, dt \leq \int_x^y M \, dt = M(y - x),$$

luego basta tomar $|y - x| < \epsilon/M$ para garantizar $|F(y) - F(x)| < \epsilon$. ■

A partir de aquí es inmediato comprobar que todos los resultados de la sección 4.10 sobre cálculo de primitivas (como la fórmula de integración por partes o la fórmula de cambio de variable) son válidas para la integral de Lebesgue en un intervalo.

Ejemplo: El resto integral de Taylor En particular ahora está completa la prueba del teorema 4.47 sobre el resto integral del polinomio de Taylor, en el que usábamos que toda función continua tiene primitiva. Más aún, haciendo el cambio de variable $t = a + s(x - a)$, con lo que $x - t = (1 - s)(x - a)$, tenemos la expresión alternativa

$$R_n(f)(x) = \frac{(x - a)^{n+1}}{n!} \int_0^1 (1 - s)^n f^{(n+1)}(a + s(x - a)) ds,$$

de modo que

$$f(x) = \sum_{k=0}^n \frac{f^{(k)}(a)}{k!} (x - a)^k + r_n(f)(x)(x - a)^{n+1},$$

donde

$$r_n(f)(x) = \frac{1}{n!} \int_0^1 (1 - s)^n f^{(n+1)}(a + s(x - a)) ds$$

es una función continua en a y, más aún, por el teorema 8.57 es de clase C^k si f es de clase C^{n+1+k} . En particular es de clase C^∞ si f lo es. ■

Ejemplo *El producto de funciones integrables no es necesariamente integrable.*

Basta considerar $f(x) = 1/\sqrt{x}$ en $X =]0, 1]$. Se trata de una función integrable, pues aplicando el teorema de Barrow en el intervalo $[1/n, 1]$ vemos que

$$\int_{1/n}^1 \frac{1}{\sqrt{x}} dx = 2 - \frac{2}{\sqrt{n}},$$

luego por el teorema de la convergencia monótona

$$\int_0^1 \frac{1}{\sqrt{x}} dx = 2.$$

En cambio, f^2 no es integrable, ya que

$$\int_{1/n}^1 \frac{1}{x} dx = \log n, \quad \text{luego} \quad \int_0^1 \frac{1}{x} dx = +\infty. \quad \blacksquare$$

En el capítulo siguiente proporcionaremos técnicas potentes para calcular integrales de funciones de varias variables. Veamos ahora un ejemplo de integral de Lebesgue en un intervalo infinito:

La función factorial La *función factorial*³ es la función $\Pi :]-1, +\infty[\rightarrow \mathbb{R}$ dada por

$$\Pi(x) = \int_0^{+\infty} t^x e^{-t} dt.$$

³La función factorial fue descubierta y estudiada por Euler, aunque fue Gauss quien la expresó en la forma que aquí usamos como definición. Legendre introdujo el cambio de variable $\Gamma(x) = \Pi(x - 1)$, que inexplicablemente ha prevalecido sobre la notación de Gauss y actualmente la función es más conocida como “función Gamma”. Nosotros respetamos la notación de Gauss pues, como veremos, resulta mucho más natural.

Recordemos que $t^x = e^{x \log t}$. Vamos a probar que el integrando es realmente una función integrable en $]0, +\infty[$ para todo $x > -1$. Probamos por separado que es integrable en $]0, 1]$ y en $[1, +\infty[$.

Si $x \geq 0$ entonces $t^x e^{-t}$ es una función continua en $[0, 1]$, luego es integrable. Si $-1 < x < 0$ entonces (si $0 \leq t \leq 1$) se cumple $0 \leq t^x e^{-t} \leq t^x$, y la función t^x es integrable: su integral es

$$\int_0^1 t^x dt = \lim_n \int_{1/n}^1 t^x dt = \lim_n \left[\frac{t^{x+1}}{x+1} \right]_{1/n}^1 = \left[\frac{t^{x+1}}{x+1} \right]_0^1 = \frac{1}{x+1}.$$

La primera igualdad se sigue del teorema de la convergencia monótona. Siempre que tengamos una integral de una función no negativa definida sobre un intervalo y de modo que restringida a intervalos menores sea continua y acotada podemos aplicar esta técnica (aplicar la regla de Barrow en intervalos menores y tomar el límite). En situaciones similares pasaremos directamente del primer término al tercero sobrentendiendo el límite.

Consideremos ahora el intervalo $[1, +\infty[$. Observemos que

$$\lim_{t \rightarrow +\infty} \frac{t^x e^{-t}}{1/t^2} = \lim_{t \rightarrow +\infty} t^{x+2} e^{-t} = 0,$$

(acotamos $x+2$ por un número natural n y aplicamos n veces la regla de l'Hôpital.)

Esto implica que existe un $M > 0$ tal que si $t \geq M$ entonces $t^x e^{-t} \leq 1/t^2$. La función $t^x e^{-t}$ es continua en el intervalo $[1, M]$, luego es integrable, luego basta probar que también lo es en $[M, +\infty[$ y a su vez basta que lo sea $1/t^2$, pero

$$\int_M^{+\infty} \frac{1}{t^2} dt = \left[-\frac{1}{t} \right]_M^{+\infty} = \frac{1}{M}.$$

Con esto tenemos probada la existencia de la función Π . La figura muestra su gráfica. ■

El teorema siguiente recoge las propiedades más importantes de la función factorial, entre ellas la que le da nombre:

Teorema 8.60 *La función factorial*

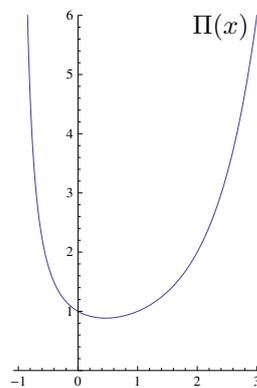
$$\Pi(x) = \int_0^{+\infty} t^x e^{-t} dt.$$

es continua en $] -1, +\infty[$ y cumple la ecuación funcional

$$\Pi(x+1) = (x+1)\Pi(x).$$

Además $\Pi(0) = 1$, de donde⁴ $\Pi(n) = n!$ para todo número natural n .

⁴Con la notación de Legendre las propiedades de la función factorial quedan distorsionadas. Por ejemplo, la función gamma cumple $\Gamma(x+1) = x\Gamma(x)$, $\Gamma(n) = (n-1)!$



DEMOSTRACIÓN: Para probar que Π es continua basta probar que lo es en un intervalo $I =]-1 + \epsilon, M[$. En este intervalo el integrando de Π está mayorado por la función $g(t) = t^{-1+\epsilon}e^{-t} + t^M e^{-t}$, que es integrable en $]0, +\infty[$, pues su integral es $\Pi(-1 + \epsilon) + \Pi(M)$. Si llamamos

$$\Pi_n(x) = \int_{-1+1/n}^n t^x e^{-x} dt,$$

el teorema 8.57 garantiza que Π_n es continua en I y basta probar que Π_n converge uniformemente a Π en I . Ahora bien, si $x \in I$ tenemos que

$$|\Pi(x) - \Pi_n(x)| \leq \int_0^{1/n} g(t) dt + \int_n^{+\infty} g(t) dt,$$

y es claro que el segundo miembro tiende a 0 con n .

La ecuación funcional se obtiene integrando por partes. En efecto:

$$\Pi(x+1) = \int_0^{+\infty} t^{x+1} e^{-t} dt = [-t^{x+1} e^{-t}]_0^{+\infty} + (x+1) \int_0^{+\infty} t^x e^{-t} dt = (x+1)\Pi(x).$$

Claramente $\Pi(0) = \int_0^{+\infty} e^{-t} dt = [-e^{-t}]_0^{+\infty} = 1 = 0!$, luego por inducción tenemos la igualdad $\Pi(n) = n!$ ■

Puede probarse que Π es de clase C^∞ en su dominio. La función factorial tiene numerosas aplicaciones en análisis real y complejo, desde la estadística hasta la teoría de números

Terminamos con un criterio de convergencia de series en términos de integrales:

Teorema 8.61 Sea $f : [1, +\infty[\rightarrow [0, +\infty[$ una función continua monótona decreciente y tal que $\lim_{x \rightarrow +\infty} f(x) = 0$. Sea $d_n = \sum_{k=1}^n f(k) - \int_1^n f(x) dx$.

a) $0 \leq f(n+1) \leq d_{n+1} \leq d_n \leq f(1)$.

b) Existe $\lim_n d_n$.

c) La serie $\sum_{k=1}^{\infty} f(k)$ converge si y sólo si $\int_1^{+\infty} f(x) dx < +\infty$.

DEMOSTRACIÓN: Es claro que a) implica b) y que b) implica c), luego basta probar a). Sean $s_n = \sum_{k=1}^n f(k)$ y $t_n = \int_1^n f(x) dx$. Se cumple

$$t_{n+1} = \int_1^{n+1} f(x) dx = \sum_{k=1}^n \int_k^{k+1} f(x) dx \leq \sum_{k=1}^n \int_k^{k+1} f(k) dx = \sum_{k=1}^n f(k) = s_n.$$

En consecuencia $0 \leq f(n+1) = s_{n+1} - s_n \leq s_{n+1} - t_{n+1} = d_{n+1}$, que es una parte de a). Además

$$\begin{aligned} d_n - d_{n+1} &= t_{n+1} - t_n - (s_{n+1} - s_n) = \int_n^{n+1} f(x) dx - f(n+1) \\ &\geq \int_n^{n+1} f(n+1) dx - f(n+1) = 0, \end{aligned}$$

luego $d_{n+1} \leq d_n \leq d_1 = f(1)$. ■

Como aplicación:

Teorema 8.62 Si $\alpha \in \mathbb{R}$, la serie $\sum_{n=1}^{\infty} \frac{1}{n^\alpha}$ converge si y sólo si $\alpha > 1$.

DEMOSTRACIÓN: Basta aplicar el teorema anterior con $f(x) = x^{-\alpha}$. Así,

$$\int_1^x f(x) dx = \begin{cases} \frac{x^{1-\alpha} - 1}{1-\alpha} & \text{si } \alpha \neq 1, \\ \log x & \text{si } \alpha = 1. \end{cases}$$

Por el teorema de la convergencia monótona tenemos que

$$\int_1^{+\infty} f(x) dx = \lim_n \int_1^n f(x) dx,$$

y es claro que este límite es finito si y sólo si $\alpha > 1$. ■

El teorema 8.61 nos indica cómo retocar la serie $\sum_n \frac{1}{n}$ para hacerla convergente:

Definición 8.63 Se llama *constante de Euler* al número real

$$\gamma = \lim_n \sum_{k=1}^n \frac{1}{k} - \log n,$$

que existe por el teorema 8.61.

Se desconoce si γ es un número racional o irracional. Sus cien primeras cifras decimales son

$$\begin{aligned} \gamma = & 0.5772156649015328606065120900824024310421593359399 \\ & 235988057672348848677267776646709369470632917467495 \dots \end{aligned}$$

Capítulo IX

Teoría de la medida II

En el capítulo anterior hemos presentado los conceptos y resultados más básicos de la teoría de la medida, y ahora vamos a exponer varios teoremas fundamentales a la hora de emplear dicha teoría. Como aplicación, en la última sección extenderemos la medida de Lebesgue a variedades diferenciables.

9.1 Producto de medidas

Vamos a ver que dos medidas en dos σ -álgebras en dos conjuntos X e Y permiten definir de manera natural una medida en una cierta σ -álgebra en $X \times Y$, a la que llamaremos medida producto de las dos medidas dadas, de modo que, por ejemplo, la medida de Lebesgue en \mathbb{R}^n es el producto de la medida de Lebesgue en \mathbb{R} consigo misma n -veces. Esto es relevante porque las integrales respecto de una medida producto pueden reducirse al cálculo de integrales respecto de los factores, de modo que así podremos reducir el cálculo de una integral de una función de n variables al cálculo de n integrales de funciones de una variable, lo que nos dará un método práctico de calcular integrales y medidas de conjuntos.

Definición 9.1 Sean X e Y dos conjuntos y \mathcal{A} , \mathcal{B} dos σ -álgebras de subconjuntos de X e Y respectivamente (a cuyos elementos llamaremos conjuntos medibles). Un *rectángulo medible* en $X \times Y$ es un conjunto de la forma $A \times B$, donde $A \in \mathcal{A}$ y $B \in \mathcal{B}$. Llamaremos *figuras elementales* a las uniones disjuntas de rectángulos medibles. Llamaremos $\mathcal{A} \times \mathcal{B}$ a la σ -álgebra generada por los rectángulos medibles. Cuando hablemos de conjuntos medibles en $X \times Y$ entenderemos que nos referimos a los de $\mathcal{A} \times \mathcal{B}$.

Observemos que las figuras elementales forman un anillo en $X \times Y$, pues las igualdades

$$\begin{aligned}(A_1 \times B_1) \cap (A_2 \times B_2) &= (A_1 \cap A_2) \times (B_1 \cap B_2) \\ (A_1 \times B_1) \setminus (A_2 \times B_2) &= ((A_1 \setminus A_2) \times B_1) \cup ((A_1 \cap A_2) \times (B_1 \setminus B_2))\end{aligned}$$

muestran que la intersección de dos rectángulos medibles es un rectángulo medible y que su diferencia es la unión de dos rectángulos medibles disjuntos, luego una figura elemental. De aquí se sigue claramente que la intersección y la diferencia de figuras elementales es una figura elemental. Lo mismo vale para la unión, pues $P \cup Q = (P \setminus Q) \cup Q$, y la unión es disjunta.

Si $E \subset X \times Y$, $x \in X$, $y \in Y$, definimos las *secciones* de E determinadas por x e y como

$$E_x = \{y \in Y \mid (x, y) \in E\}, \quad E^y = \{x \in X \mid (x, y) \in E\}.$$

Teorema 9.2 *En las condiciones anteriores, si E es medible en $X \times Y$, entonces E_x y E^y son medibles en Y y en X respectivamente.*

DEMOSTRACIÓN: Sea \mathcal{C} el conjunto de todos los $E \in \mathcal{A} \times \mathcal{B}$ tales que $E_x \in \mathcal{B}$ para todo $x \in X$. Si $E = A \times B$ entonces $E_x = B$ para todo $x \in A$, y $E_x = \emptyset$ si $x \notin A$ luego todos los rectángulos medibles están en \mathcal{C} . Ahora vemos que \mathcal{C} es una σ -álgebra, de donde se sigue que $\mathcal{A} \times \mathcal{B} \subset \mathcal{C}$, lo que prueba el teorema (para E_x , el caso E^y es análogo).

- a) Obviamente $X \times Y \in \mathcal{C}$.
- b) Si $E \in \mathcal{C}$, entonces $(X \times Y \setminus E)_x = Y \setminus E_x \in \mathcal{B}$, luego $X \times Y \setminus E \in \mathcal{C}$.
- c) Si $\{E_i\}_{i=1}^{\infty} \subset \mathcal{C}$ y $E = \bigcup_{i=1}^{\infty} E_i$, entonces $E_x = \bigcup_{i=1}^{\infty} E_{ix} \in \mathcal{B}$. Por lo tanto $E \in \mathcal{C}$.

■

Veamos ahora la relación entre la medibilidad de funciones en un producto y en los factores. Siempre en las mismas condiciones, si $f : X \times Y \rightarrow Z$, $x \in X$, $y \in Y$, definimos $f_x : Y \rightarrow Z$ y $f^y : X \rightarrow Z$ como las aplicaciones dadas por $f_x(y) = f(x, y)$, $f^y(x) = f(x, y)$.

Teorema 9.3 *Si $f : X \times Y \rightarrow Z$ es una función medible, entonces f_x es medible para todo $x \in X$ y f^y es medible para todo $y \in Y$.*

DEMOSTRACIÓN: Si V es un abierto en Z , claramente $f_x^{-1}[V] = f^{-1}[V]_x$, luego es medible. Por lo tanto f_x es medible. Igualmente se razona con f^y . ■

Con esto estamos casi a punto de definir el producto de medidas. La definición se apoyará en el teorema siguiente.

Teorema 9.4 *Sean X e Y espacios medida con medidas σ -finitas μ y ν . Sea E un subconjunto medible de $X \times Y$. Entonces las aplicaciones $\nu(E_x)$ y $\mu(E^y)$ son funciones medibles de x e y respectivamente. Además*

$$\int_X \nu(E_x) d\mu = \int_Y \mu(E^y) d\nu.$$

DEMOSTRACIÓN: Notar que por el teorema 9.2 los conjuntos E_x, E^y son medibles, luego tiene sentido considerar $\nu(E_x)$ y $\mu(E^y)$.

Llamemos \mathcal{C} a la familia de todos los subconjuntos medibles de $X \times Y$ para los que se cumple el teorema. Vamos a probar que \mathcal{C} tiene las propiedades siguientes:

- a) \mathcal{C} contiene a los rectángulos medibles.
- b) Si $\{Q_n\}_{n=1}^{\infty} \subset \mathcal{C}$ es creciente entonces $Q = \bigcup_{n=1}^{\infty} Q_n \in \mathcal{C}$.
- c) Si $\{Q_n\}_{n=1}^{\infty} \subset \mathcal{C}$ son disjuntos dos a dos entonces $Q = \bigcup_{n=1}^{\infty} Q_n \in \mathcal{C}$.
- d) Si $\{Q_n\}_{n=1}^{\infty} \subset \mathcal{C}$ es decreciente y $Q_1 \subset U \times V$, con $\mu(U), \nu(V) < +\infty$, entonces $Q = \bigcap_{n=1}^{\infty} Q_n \in \mathcal{C}$.

En efecto, si $U \times V$ es un rectángulo medible, entonces

$$(U \times V)_x = \begin{cases} V & \text{si } x \in U \\ \emptyset & \text{si } x \notin U. \end{cases}$$

Por lo tanto $\nu((U \times V)_x) = \nu(V)\chi_U$, que es una función medible. Igualmente $\mu((U \times V)_y) = \mu(U)\chi_V$. Las integrales valen ambas $\mu(U)\nu(V)$, luego $U \times V$ está en \mathcal{C} . Esto prueba a).

Para demostrar b) observamos que $Q_x = \bigcup_{n=1}^{\infty} (Q_n)_x$, y la sucesión es creciente, por lo que $\nu(Q_x) = \lim_n \nu((Q_n)_x)$. Como los conjuntos Q_n están en \mathcal{C} , las funciones $\nu((Q_n)_x)$ son medibles, luego su límite puntual $\nu(Q_x)$ también lo es. Igualmente ocurre con $\mu(Q_y)$. El teorema de la convergencia monótona da la igualdad de las integrales, luego $Q \in \mathcal{C}$.

La prueba de c) es similar, usando ahora que $\nu(Q_x)$ es la suma de las funciones $\nu((Q_n)_x)$ en lugar del límite.

En el caso d) tenemos también $\nu(Q_x) = \lim_n \nu((Q_n)_x)$, pero ahora la sucesión no es monótona creciente. La única diferencia es que en lugar del teorema de la convergencia monótona usamos el teorema de la convergencia dominada. La hipótesis $\nu(V) < +\infty$ garantiza que las funciones $\nu((Q_n)_x)$ están dominadas por la función integrable χ_V .

Estamos suponiendo que las medidas en X y en Y son σ -finitas, lo cual significa que podemos expresar $X = \bigcup_{n=1}^{\infty} X_n$ e $Y = \bigcup_{n=1}^{\infty} Y_n$ para ciertos conjuntos medibles de medida finita que además podemos suponer disjuntos dos a dos.

Sea E un conjunto medible en $X \times Y$. Definamos $E_{mn} = E \cap (X_m \times Y_n)$ y sea \mathcal{M} la familia de todos los conjuntos E tales que los E_{mn} así definidos están

en \mathcal{C} . las propiedades b) y d) muestran que \mathcal{M} es una clase monótona, mientras que a) y c) muestran que contiene a las figuras elementales. El teorema 8.22 implica ahora que \mathcal{M} contiene a todos los conjuntos medibles de $X \times Y$ (pues éstos forman el σ -anillo generado por las figuras elementales, que coincide con la clase monótona que generan).

Así pues, para todo conjunto medible E , los conjuntos E_{mn} están en \mathcal{C} , pero claramente E es unión disjunta de los E_{mn} , luego por c) tenemos $E \in \mathcal{C}$, es decir, todo conjunto medible E cumple el teorema. ■

Definición 9.5 Sean X e Y espacios con medidas σ -finitas μ y ν . Definimos la *medida producto* $\mu \times \nu$ como la dada por

$$(\mu \times \nu)(Q) = \int_X \nu(Q_x) d\mu(x) = \int_Y \mu(Q^y) d\nu(y).$$

Con el teorema 8.50 se prueba fácilmente que $\mu \times \nu$ es realmente una medida en la σ -álgebra producto. Además es claro que sobre los rectángulos medibles tenemos $(\mu \times \nu)(A \times B) = \mu(A)\nu(B)$ (con el convenio $0 \cdot \infty = 0$).

Conviene dar una caracterización de la medida producto que no dependa del teorema anterior:

Teorema 9.6 *Dados dos espacios X e Y con medidas σ -finitas μ y ν , la medida producto es la única que cumple que $(\mu \times \nu)(A \times B) = \mu(A)\nu(B)$ para todo rectángulo medible $A \times B$.*

DEMOSTRACIÓN: Supongamos que dos medidas λ_1 y λ_2 coinciden sobre los rectángulos medibles con la medida producto. Descompongamos $X = \bigcup_{n=1}^{\infty} X_n$ e $Y = \bigcup_{n=1}^{\infty} Y_n$, para ciertos conjuntos medibles de medida finita disjuntos dos a dos. Sea \mathcal{M} la familia de los conjuntos medibles E de $X \times Y$ tales que $\lambda_1(E \cap (X_m \times Y_n)) = \lambda_2(E \cap (X_m \times Y_n))$ para todo m, n . Es claro que \mathcal{M} es una clase monótona que contiene a las figuras elementales, luego por 8.22 tenemos que \mathcal{M} contiene a todos los conjuntos medibles. De aquí se sigue que las dos medidas coinciden sobre cualquier conjunto medible. ■

Veamos ahora que toda esta teoría es aplicable a la medida de Lebesgue en \mathbb{R}^n . Primero probemos un hecho general:

Teorema 9.7 *Si X e Y son dos espacios topológicos con bases numerables, entonces el producto de las σ -álgebras de Borel es la σ -álgebra de Borel del producto. En particular el producto de medidas de Borel es una medida de Borel.*

DEMOSTRACIÓN: Si U y V son conjuntos de Borel en X e Y respectivamente, entonces $U \times V$ es un conjunto de Borel en el producto, pues es la antiimagen de U por la proyección en X , que es continua, luego medible. Igualmente $X \times V$ es

un conjunto de Borel, y también lo es $U \times V$ por ser la intersección de ambos. De aquí se sigue que todas las figuras elementales son conjuntos de Borel, luego también lo son todos los conjuntos medibles en $X \times Y$. Recíprocamente, los productos de abiertos básicos $U \times V$ forman una base numerable de $X \times Y$, luego todo abierto de $X \times Y$ es unión numerable de estos abiertos básicos, luego todo abierto de $X \times Y$ es medible, luego todo conjunto de Borel es medible. ■

Teorema 9.8 *La medida de Lebesgue en \mathbb{R}^{m+n} (restringida a los conjuntos de Borel) es el producto de la medida de Lebesgue en \mathbb{R}^m por la medida de Lebesgue en \mathbb{R}^n (restringidas ambas a los conjuntos de Borel).*

En efecto, es claro que las celdas son rectángulos medibles, y la medida producto coincide sobre ellas con la medida de Lebesgue, luego es la medida de Lebesgue. ■

Ahora podemos dar una interpretación de la integral de Lebesgue:

Teorema 9.9 *Sea (X, \mathcal{A}, μ) un espacio medida, sea $f : X \rightarrow [0, +\infty[$ una función integrable y sea $A = \{(x, y) \in X \times \mathbb{R} \mid 0 \leq y \leq f(x)\}$. Entonces A es medible respecto del producto $\mu \times m$ de μ por la medida de Lebesgue en \mathbb{R} y*

$$\int_X f \, d\mu = (\mu \times m)(A).$$

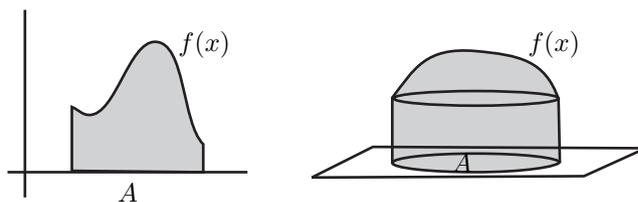
DEMOSTRACIÓN: Consideremos la función $g : X \times \mathbb{R} \rightarrow \mathbb{R}^2$ dada por $g(x, y) = (f(x), y)$ es medible, pues si U y V son dos abiertos en \mathbb{R} entonces $g^{-1}[U \times V] = f^{-1}[U] \times V$. Pero

$$A = g^{-1}[\{(u, v) \in \mathbb{R} \mid 0 \leq v \leq u\}]$$

y el conjunto de la derecha es cerrado, luego A es medible. Para cada $x \in X$, tenemos que $A_x = [0, f(x)]$, luego

$$(\mu \times m)(A) = \int_X m(A_x) \, d\mu = \int_X f(x) \, d\mu. \quad \blacksquare$$

Por ejemplo, la integral de una función $f : A \rightarrow [0, +\infty[$ respecto de la medida de Lebesgue es el área o el volumen sombreados en la figura, según si $A \subset \mathbb{R}$ o $A \subset \mathbb{R}^2$.



Si una función f toma valores positivos y negativos, la descomponemos como $f = f^+ - f^-$ y vemos que la integral de f es (en el caso de una variable) la diferencia entre el área situada bajo su gráfica por encima del eje X menos la situada sobre su gráfica por debajo del eje X , y análogamente en dimensiones superiores.

Ejemplo Vamos a calcular el área de la elipse E de semiejes a , y b , formada por los puntos que cumplen

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} \leq 1.$$

Despejamos $y^2 \leq \frac{b^2}{a^2}(a^2 - x^2)$, lo cual significa que

$$E_x = \left[-\frac{b}{a} \sqrt{a^2 - x^2}, \frac{b}{a} \sqrt{a^2 - x^2} \right]$$

siempre que $-a \leq x \leq a$ (y $E_x = \emptyset$ en caso contrario). Por consiguiente

$$m(E) = \int_{-a}^a m(E_x) dx = \frac{2b}{a} \int_{-a}^a \sqrt{a^2 - x^2} dx =$$

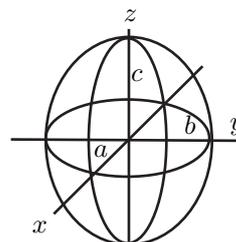
El cambio $x = a \sin t$ transforma la integral en

$$2ab \int_{-\pi/2}^{\pi/2} \cos^2 t dt = 2ab \int_{-\pi/2}^{\pi/2} \frac{1 + \cos 2t}{2} dt = 2ab \left[\frac{t}{2} + \frac{\sin 2t}{4} \right]_{-\pi/2}^{\pi/2} = \pi ab.$$

En particular, el área de un círculo de radio r es πr^2 . ■

Ejemplo Calculamos ahora el volumen del elipsoide E de semiejes a , b , c , que es la figura formada por los puntos que cumplen

$$\frac{x^2}{a^2} + \frac{y^2}{b^2} + \frac{z^2}{c^2} \leq 1.$$



Esto equivale a

$$\frac{x^2}{(a\sqrt{1 - z^2/c^2})^2} + \frac{y^2}{(b\sqrt{1 - z^2/c^2})^2} \leq 1,$$

lo cual significa que E_z es la elipse de semiejes $a\sqrt{1 - z^2/c^2}$ y $b\sqrt{1 - z^2/c^2}$ (siempre que $-c \leq z \leq c$). Por lo tanto

$$m(E) = \int_{-c}^c m(E_z) dz = \pi ab \int_{-c}^c (1 - z^2/c^2) dz = \pi ab \left[z - \frac{z^3}{3c^2} \right]_{-c}^c = \frac{4}{3} \pi abc.$$

En particular, el volumen de una esfera de radio r es $\frac{4}{3} \pi r^3$. ■

Ejercicio: Probar que el volumen de la figura de revolución que resulta de hacer girar sobre el eje X la función $f[a, b] \rightarrow [0, +\infty[$ es $V = \pi \int_a^b f^2(x) dx$.

Ejemplo Consideremos una bola cerrada n -dimensional de radio r :

$$B_r^n = \{x \in \mathbb{R}^n \mid x_1^2 + \cdots + x_n^2 \leq r^2\}.$$

Vamos a probar que $m(B_r^n) = v_n r^n$, para una cierta constante v_n que también calcularemos. Razonaremos por inducción sobre n . El ejemplo anterior nos da que el resultado es cierto para $n = 3$ con $v_3 = (4/3)\pi$. Claramente, también es cierto para $n = 2$ con $v_2 = \pi$. En general, la condición que define a B_r^n es

$$x_1^2 + \cdots + x_{n-1}^2 \leq \sqrt{r^2 - x_n^2},$$

lo cual significa que $(B_r^n)_{x_n} = B_{\sqrt{r^2 - x_n^2}}^{n-1}$. Por lo tanto, por hipótesis de inducción,

$$m(B_r^n) = \int_{-r}^r m v_{n-1} (r^2 - x_n^2)^{(n-1)/2} dx_n.$$

Hacemos el cambio de variable $x_n = r \sin \theta$, con lo que

$$m(B_r^n) = v_{n-1} r^n \int_{-\pi/2}^{\pi/2} \cos^n \theta d\theta.$$

Si llamamos $\kappa_n = \int_{-\pi/2}^{\pi/2} \cos^n \theta d\theta$, hemos probado que $m(B_r^n) = v_{n-1} \kappa_n r^n$, luego el resultado es cierto para n con $v_n = v_{n-1} \kappa_n$.

Con la notación del ejemplo de la página 190 tenemos que $\kappa_n = [I_{0,n}]_{-\pi/2}^{\pi/2}$. Si suponemos $n \geq 2$ los cálculos de dicho ejemplo nos dan que

$$\kappa_n = \frac{n-1}{n} \kappa_{n-2}.$$

Una simple inducción nos da ahora que

$$\kappa_n \kappa_{n-1} = \frac{2\pi}{n}.$$

En efecto, basta usar la relación anterior y tener en cuenta que

$$\kappa_0 = \int_{-\pi/2}^{\pi/2} d\theta = \pi, \quad \kappa_1 = \int_{-\pi/2}^{\pi/2} \cos \theta d\theta = 2.$$

Por consiguiente $v_n = v_{n-1} \kappa_n = v_{n-2} \kappa_n \kappa_{n-1}$. Así llegamos a las relaciones

$$v_1 = 2, \quad v_2 = \pi, \quad v_n = \frac{2\pi}{n} v_{n-2},$$

que nos permiten calcular fácilmente v_n para cualquier valor de n . ■

Ejemplo Dado un subconjunto A de un hiperplano $H \subset \mathbb{R}^n$ y un punto $V \in \mathbb{R}^n \setminus H$, definimos el *cono* de base A y vértice V como la unión de todos los segmentos que unen V con un punto de A . Podemos fijar un sistema de referencia en el que V sea el origen de coordenadas y el hiperplano H tenga ecuación $x_n = h$, donde h es la distancia de V a H .

Entonces, identificando a A con un subconjunto de \mathbb{R}^{n-1} ,

$$C = \{\lambda(a, h) \mid a \in A, 0 \leq \lambda \leq 1\}.$$

Si, por simplicidad, suponemos que A es cerrado, entonces C es también cerrado, luego es medible Lebesgue.

Para calcular su medida observamos que $x \in C_u$ si y sólo si $(x, u) \in C$, si y sólo si existe un λ en $[0, 1]$ tal que $(x, u) = \lambda(a, h)$, lo que obliga a que $\lambda = u/h$ (y esto a su vez a que $0 \leq u \leq h$ para que C_u no sea vacío) y entonces $x = (u/h)a$. Por lo tanto, $C_u = f_{u/h}[A]$, donde $f_{u/h}$ es la homotecia de centro V y razón u/h en \mathbb{R}^{n-1} . El teorema 8.36 nos da que $m(C_u) = (u/h)^{n-1}m(A)$. Por lo tanto

$$m(C) = \int_0^h m(C_u) du = \frac{m(A)}{h^n} \int_0^h u^{n-1} du = \frac{m(A)}{h^{n-1}} \left[\frac{u^n}{n} \right]_0^h = \frac{m(A)h}{n}.$$

En resumen, la medida de un cono es la medida de su base, por su altura, dividida entre la dimensión del espacio. Esto generaliza al hecho de que el área de un triángulo es su base por su altura dividida entre 2, y para un cono tridimensional tenemos que su volumen es un tercio de la superficie de su base por su altura. ■

El teorema siguiente permite calcular integrales arbitrarias en productos:

Teorema 9.10 (Teorema de Fubini) Sean X e Y dos espacios medida con medidas σ -finitas μ y ν y sea $f : X \times Y \rightarrow [-\infty, +\infty]$ una función medible.

a) Si $f \geq 0$, entonces las funciones

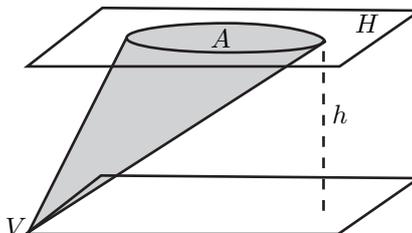
$$\int_Y f_x d\nu \quad (\text{como función de } x) \quad \text{e} \quad \int_X f_y d\mu \quad (\text{como función de } y)$$

son medibles, y se cumple

$$\int_{X \times Y} f d(\mu \times \nu) = \int_X \left(\int_Y f_x d\nu \right) d\mu = \int_Y \left(\int_X f_y d\mu \right) d\nu.$$

b) Si $\int_X \left(\int_Y |f_x| d\nu \right) d\mu < +\infty$, entonces $f \in L^1(\mu \times \nu)$.

c) Si $f \in L^1(\mu \times \nu)$, entonces $f_x \in L^1(\nu)$ p.c.t. x , $f_y \in L^1(\mu)$ p.c.t. y , y las funciones definidas en a) p.c.t.p. están en $L^1(\mu)$ y $L^1(\nu)$ respectivamente y sus integrales coinciden (según se afirma en a).



DEMOSTRACIÓN: a) Por el teorema 9.3, las funciones f_x y f_y son medibles, luego tienen sentido sus integrales. Si $f = \chi_Q$, para un cierto conjunto medible $Q \subset X \times Y$, entonces a) se reduce al teorema 9.4 y a la definición de la medida producto. De aquí se sigue que las igualdades de a) son válidas cuando f es una función simple. En general existe una sucesión creciente de funciones simples $\{s_n\}_{n=1}^{\infty}$ que converge puntualmente a f . Entonces es claro que $\{(s_n)_x\}_{n=1}^{\infty}$ converge puntualmente a f_x y, por el teorema de la convergencia monótona concluimos que

$$\lim_n \int_Y (s_n)_x d\nu = \int_Y f_x d\nu.$$

Como las funciones simples cumplen a) tenemos que las funciones $\int_Y (s_n)_x d\nu$ son medibles y

$$\int_{X \times Y} s_n d(\mu \times \nu) = \int_X \left(\int_Y (s_n)_x d\nu \right) d\mu.$$

Consecuentemente el límite $\int_Y f_x d\nu$ es medible y aplicando el teorema de la convergencia monótona a los dos miembros de la igualdad anterior queda

$$\int_{X \times Y} f d(\mu \times \nu) = \int_X \lim_n \left(\int_Y (s_n)_x d\nu \right) d\mu = \int_X \left(\int_Y f_x d\nu \right) d\mu.$$

La otra igualdad se prueba análogamente.

Las hipótesis de b) implican por a) que $|f|$ es integrable, luego f también lo es.

Para probar c) descompongamos $f = f^+ - f^-$. Tenemos que f^+ y f^- son integrables. Por a) la integrabilidad de f^+ significa que

$$\int_{X \times Y} f^+ d(\mu \times \nu) = \int_X \left(\int_Y f_x^+ d\nu \right) d\mu < +\infty,$$

luego $\int_Y f_x^+ d\nu$ ha de ser finita salvo a lo sumo en un conjunto nulo, y lo mismo el válido para $\int_Y f_x^- d\nu$. Salvo para los puntos x en la unión de los dos conjuntos nulos, tenemos que la integral

$$\int_Y f_x d\nu = \int_Y f_x^+ d\nu - \int_Y f_x^- d\nu$$

está definida y es finita, es decir, que la función $\int_Y f_x d\nu$ es integrable. Además

$$\begin{aligned} \int_{X \times Y} f d(\mu \times \nu) &= \int_{X \times Y} f^+ d(\mu \times \nu) - \int_{X \times Y} f^- d(\mu \times \nu) \\ &= \int_X \left(\int_Y f_x^+ d\nu \right) d\mu - \int_X \left(\int_Y f_x^- d\nu \right) d\mu = \int_X \left(\int_Y f_x d\nu \right) d\mu. \end{aligned}$$

La otra parte de c) es análoga. ■

En el caso de la medida de Lebesgue en \mathbb{R}^n sustituiremos dm por $dx_1 \cdots dx_n$. En estos términos el teorema de Fubini (por ejemplo para dos variables) se expresa como sigue:

$$\int_{\mathbb{R}^2} f(x, y) dx dy = \int_{\mathbb{R}} \left(\int_{\mathbb{R}} f(x, y) dx \right) dy = \int_{\mathbb{R}} \left(\int_{\mathbb{R}} f(x, y) dy \right) dx.$$

Si trabajamos con funciones continuas sobre un compacto no hemos de preocuparnos de la integrabilidad.

9.2 El teorema de Riesz

En el capítulo anterior hemos visto que toda medida tiene asociada una integral, y ahora vamos a ver que, en cierto sentido, toda integral tiene asociada una medida, es decir, que todo operador que cumpla ciertos requisitos necesarios para que pueda ser una integral es realmente la integral respecto de cierta medida. Necesitamos un concepto adicional sobre medidas en espacios topológicos:

Definición 9.11 Diremos que una medida μ en un espacio topológico X es *regular* si todos los abiertos son medibles, los subespacios compactos tienen medida finita y para todo conjunto medible E se cumple

$$\mu(E) = \inf\{\mu(V) \mid E \subset V, \quad V \text{ abierto}\}$$

y

$$\mu(E) = \sup\{\mu(K) \mid K \subset E, \quad K \text{ compacto}\}.$$

Diremos que la medida es *casi regular* si la segunda propiedad se cumple al menos cuando $\mu(E) < +\infty$ y cuando E es abierto.

En definitiva una medida es regular si la medida de todo conjunto medible puede aproximarse por la medida de un abierto mayor y de un compacto menor. El concepto de medida casi regular lo introducimos por cuestiones técnicas, pero seguidamente vamos a probar que en todos los espacios que nos van a interesar es equivalente a la regularidad.

Diremos que un espacio topológico es *σ -compacto* si es unión numerable de conjuntos compactos. Por ejemplo, todo abierto Ω en \mathbb{R}^n es σ -compacto, pues puede expresarse como unión de los compactos

$$\Omega_k = \{x \in \Omega \mid \|x\| \leq k, \quad d(x, \mathbb{R}^n \setminus \Omega) \geq 1/k\}, \quad k = 1, 2, 3, \dots$$

Teorema 9.12 *Toda medida casi regular en un espacio σ -compacto es regular.*

DEMOSTRACIÓN: Supongamos que X es la unión de los compactos $\{K_n\}_{n=1}^{\infty}$. Sustituyendo cada K_n por su unión con los precedentes podemos suponer que

la sucesión es creciente. Dado un conjunto de Borel B tal que $\mu(B) = +\infty$, tenemos que

$$B = \bigcup_{n=1}^{\infty} B \cap K_n,$$

y, como la unión es creciente, $\sup_n \mu(B \cap K_n) = \mu(B) = +\infty$. Dado $R > 0$ existe un n tal que $\mu(B \cap K_n) > R + 1$ y, como μ es casi regular, $B \cap K_n$ tiene medida finita y existe un compacto $K \subset B \cap K_n$ tal que $\mu(K) > R$, lo que prueba que μ es regular. ■

Ejercicio: Demostrar que la medida de Lebesgue es regular.

Teorema 9.13 (Teorema de representación de Riesz) *Sea X un espacio localmente compacto y sea $T : C_c(X) \rightarrow \mathbb{R}$ una aplicación lineal tal que si $f \geq 0$ entonces $T(f) \geq 0$. Entonces existe una única medida de Borel casi regular μ en X tal que para toda función $f \in C_c(X)$ se cumple*

$$T(f) = \int_X f d\mu.$$

DEMOSTRACIÓN: Veamos primero la unicidad. Es claro que una medida casi regular está completamente determinada por los valores que toma sobre los conjuntos compactos, luego basta probar que si μ_1 y μ_2 representan a T en el sentido del teorema entonces $\mu_1(K) = \mu_2(K)$ para todo compacto K .

Por la regularidad existe un abierto V tal que $K \subset V$ y $\mu_2(V) < \mu_2(K) + \epsilon$. Por el lema de Urysohn 3.16 existe una función $K \prec f \prec V$. Entonces

$$\begin{aligned} \mu_1(K) &= \int_X \chi_K d\mu_1 \leq \int_X f d\mu_1 = T(f) = \int_X f d\mu_2 \\ &\leq \int_X \chi_V d\mu_2 = \mu_2(V) < \mu_2(K) + \epsilon. \end{aligned}$$

Por consiguiente $\mu_1(K) \leq \mu_2(K)$ e igualmente se prueba la desigualdad contraria.

Para cada abierto V de X definimos $\mu(V) = \sup\{T(f) \mid f \prec V\}$. Es obvio que si $V_1 \subset V_2$ entonces $\mu(V_1) \leq \mu(V_2)$, luego si definimos

$$\mu(E) = \inf\{\mu(V) \mid E \subset V, V \text{ abierto}\}, \quad \text{para todo } E \subset X,$$

es claro que la medida de un abierto es la misma en los dos sentidos en que la tenemos definida.

Aunque hemos definido la medida de cualquier conjunto, ésta sólo cumplirá las propiedades de las medidas al restringirla a una cierta σ -álgebra que contiene a la σ -álgebra de Borel. Concretamente, definimos \mathcal{M}_F como la familia de subconjuntos E de X tales que $\mu(E) < +\infty$ y

$$\mu(E) = \sup\{\mu(K) \mid K \subset E, K \text{ compacto}\}.$$

Definimos \mathcal{M} como la familia de todos los $E \subset X$ tales que $E \cap K \in \mathcal{M}_F$ para todo compacto K . Probaremos que \mathcal{M} es una σ -álgebra que contiene a la σ -álgebra de Borel y que la restricción de μ a \mathcal{M} es una medida casi regular. Veremos también que \mathcal{M}_F está formada por los conjuntos de \mathcal{M} de medida finita. Dividimos la prueba en varios pasos.

1) Si $\{E_i\}_{i=1}^{\infty}$ son subconjuntos de X , entonces

$$\mu\left(\bigcup_{i=1}^{\infty} E_i\right) \leq \sum_{i=1}^{\infty} \mu(E_i).$$

Probamos primero que si V_1 y V_2 son abiertos $\mu(V_1 \cup V_2) \leq \mu(V_1) + \mu(V_2)$. Tomemos $g \prec V_1 \cup V_2$ arbitraria. Por el teorema 3.17 existen funciones h_1 y h_2 tales que $h_i \prec V_i$ y $h_1 + h_2$ vale 1 sobre los puntos del soporte de g . Por lo tanto $h_i g \prec V_i$, $g = h_1 g + h_2 g$.

$$T(g) = T(h_1 g) + T(h_2 g) \leq \mu(V_1) + \mu(V_2).$$

Como esto se cumple para toda $g \prec V_1 + V_2$, concluimos la desigualdad buscada.

Podemos suponer que $\mu(E_i) < +\infty$ para todo i , o la desigualdad que queremos probar se cumpliría trivialmente. Dado $\epsilon > 0$ la definición de μ implica que existen abiertos V_i que contienen a E_i de modo que $\mu(V_i) < \mu(E_i) + \epsilon/2^i$. Sea V la unión de todos los V_i y tomemos $f \prec V$. Como f tiene soporte compacto en realidad $f \prec V_1 \cup \dots \cup V_n$ para algún n , luego

$$T(f) \leq \mu(V_1 \cup \dots \cup V_n) \leq \mu(V_1) + \dots + \mu(V_n) \leq \sum_{i=1}^{\infty} \mu(E_i) + \epsilon.$$

Como esto vale para toda $f \prec V$, resulta que

$$\mu(V) \leq \sum_{i=1}^{\infty} \mu(E_i) + \epsilon.$$

Ahora bien, la unión de los E_i está contenida en V , y la función μ es claramente monótona por su definición, luego

$$\mu\left(\bigcup_{i=1}^{\infty} E_i\right) \leq \sum_{i=1}^{\infty} \mu(E_i) + \epsilon.$$

Como esto vale para todo ϵ tenemos la desigualdad buscada.

2) Si K es compacto, entonces $K \in \mathcal{M}_F$ y $\mu(K) = \inf\{T(f) \mid K \prec f\}$. En particular los compactos tienen medida finita.

Si $K \prec f$ y $0 < \alpha < 1$, sea $V_\alpha = \{x \in X \mid f(x) > \alpha\}$. Entonces $K \subset V_\alpha$ y si $g \prec V_\alpha$ se cumple $\alpha g \leq f$. Por lo tanto

$$\mu(K) \leq \mu(V_\alpha) = \sup\{T(g) \mid g \prec V_\alpha\} \leq \alpha^{-1} T(f).$$

Si hacemos que α tienda a 1 concluimos que $\mu(K) \leq T(f)$ y es obvio que K está en \mathcal{M}_F .

Dado $\epsilon > 0$ existe un abierto V tal que $K \subset V$ y $\mu(V) < \mu(K) + \epsilon$. Existe una función $K \prec f \prec V$, luego

$$\mu(K) \leq T(f) \leq \mu(V) < \mu(K) + \epsilon,$$

lo que prueba que $\mu(K) = \inf\{T(f) \mid K \prec f\}$.

3) \mathcal{M}_F contiene a todos los abiertos de medida finita.

Sea V un abierto de medida finita y α un número real tal que $\alpha < \mu(V)$. Existe $f \prec V$ tal que $\alpha < T(f)$. Entonces $T(f) \leq \mu(K)$, pues en caso contrario existiría un abierto W tal que $K \subset W$ y $\mu(K) \leq \mu(W) < T(f)$, pero $f \prec W$, luego $T(f) \leq \mu(W)$, contradicción. Así hemos encontrado un compacto $K \subset V$ tal que $\alpha < \mu(K)$, lo que prueba que $V \in \mathcal{M}_F$.

4) Si $\{E_i\}_{i=1}^{\infty}$ son elementos disjuntos de \mathcal{M}_F , y $E = \bigcup_{i=1}^{\infty} E_i$ entonces

$$\mu(E) = \sum_{i=1}^{\infty} \mu(E_i).$$

Si además $\mu(E) < +\infty$ entonces $E \in \mathcal{M}_F$.

Veamos primero que si K_1 y K_2 son compactos disjuntos $\mu(K_1 \cup K_2) = \mu(K_1) + \mu(K_2)$. Dado $\epsilon > 0$ existe $K_1 \prec f \prec X \setminus K_2$. Por el paso 2) existe g tal que $K_1 \cup K_2 \prec g$ y $T(g) < \mu(K_1 \cup K_2) + \epsilon$. Claramente $K_1 \prec fg$ y $K_2 \prec (1-f)g$, luego

$$\mu(K_1) + \mu(K_2) \leq T(fg) + T(g-fg) = T(g) < \mu(K_1 \cup K_2) + \epsilon,$$

luego tenemos $\mu(K_1) + \mu(K_2) \leq \mu(K_1 \cup K_2)$ y el paso 1) nos da la otra desigualdad.

Pasando al caso general, de nuevo por 1) basta probar una desigualdad. Ésta es trivial si $\mu(E) = +\infty$, luego podemos suponer que E tiene medida finita. Fijado $\epsilon > 0$, puesto que $E_i \in \mathcal{M}_F$ existen compactos $H_i \subset E_i$ tales que $\mu(H_i) > \mu(E_i) - \epsilon/2^i$. Sea $K_n = H_1 \cup \dots \cup H_n$. Entonces

$$\mu(E) \geq \mu(K_n) = \sum_{i=1}^n \mu(H_i) > \sum_{i=1}^n \mu(E_i) - \epsilon,$$

lo cual nos da claramente la desigualdad buscada. La desigualdad anterior muestra también que $\mu(K_n)$ tiende a $\mu(E)$ cuando n tiende a ∞ (una vez sabemos que la serie suma $\mu(E)$), lo que implica que $E \in \mathcal{M}_F$.

5) Si $E \in \mathcal{M}_F$ y $\epsilon > 0$, existen un compacto K y un abierto V tales que $K \subset E \subset V$ y $\mu(V \setminus K) < \epsilon$.

Por definición de \mathcal{M}_F y de μ existen K y V tales que

$$\mu(V) - \frac{\epsilon}{2} < \mu(E) < \mu(K) + \frac{\epsilon}{2}.$$

Puesto que $V \setminus K$ es abierto, por 3) tenemos que $V \setminus K \in \mathcal{M}_F$, luego 4) implica que

$$\mu(K) + \mu(V \setminus K) = \mu(V) < \mu(K) + \epsilon.$$

6) Si $A, B \in \mathcal{M}_F$ entonces $A \setminus B, A \cup B, A \cap B \in \mathcal{M}_F$.

Aplicamos el paso anterior a los conjuntos A y B , lo que nos da conjuntos K_i y V_i , para $i = 1, 2$, de modo que $K_1 \subset A \subset V_1$, $K_2 \subset B \subset V_2$ y $\mu(V_i \setminus K_i) < \epsilon$. Entonces

$$A \setminus B \subset V_1 \setminus K_2 \subset (V_1 \setminus K_1) \cup (K_1 \setminus V_2) \cup (V_2 \setminus K_2),$$

luego el paso 1) implica $\mu(A \setminus B) \leq \epsilon + \mu(K_1 \setminus V_2) + \epsilon$ y $K_1 \setminus V_2$ es un subconjunto compacto de $A \setminus B$, luego esto prueba que $A \setminus B \in \mathcal{M}_F$.

Ahora, $A \cup B = (A \setminus B) \cup B$, luego el paso 4) implica que $A \cup B \in \mathcal{M}_F$ y como $A \cap B = A \setminus (A \setminus B)$, también $A \cap B \in \mathcal{M}_F$.

7) \mathcal{M} es una σ -álgebra que contiene a la σ -álgebra de Borel.

Si $A \in \mathcal{M}$ y K es un compacto en X , entonces $(X \setminus A) \cap K = K \setminus (A \cap K)$, luego $(X \setminus A) \cap K$ es diferencia de dos elementos de \mathcal{M}_F , luego está en \mathcal{M}_F , luego $X \setminus A \in \mathcal{M}$.

Sea $A = \bigcup_{i=1}^{\infty} A_i$ una unión de elementos de \mathcal{M} . Si K es un compacto en X , tomamos $B_1 = A_1 \cap K$ y $B_n = (A_n \cap K) \setminus (B_1 \cup \dots \cup B_{n-1})$, con lo que cada $B_n \in \mathcal{M}_F$ y son disjuntos dos a dos. Por 4) tenemos que $A \cap K$ (la unión de los B_n) está en \mathcal{M}_F , luego $A \in \mathcal{M}$. Esto prueba que \mathcal{M} es una σ -álgebra.

Si C es un cerrado de X y K es compacto, entonces $C \cap K$ es compacto, luego está en \mathcal{M}_F , luego $C \in \mathcal{M}$. Por lo tanto \mathcal{M} contiene a todos los cerrados, luego a todos los abiertos, luego a todos los conjuntos de Borel.

8) \mathcal{M}_F está formado por los conjuntos de \mathcal{M} de medida finita.

Si $E \in \mathcal{M}_F$, los pasos 2) y 6) implican que $E \cap K \in \mathcal{M}_F$, luego $E \in \mathcal{M}$. Recíprocamente, si $E \in \mathcal{M}$ tiene medida finita, dado $\epsilon > 0$ existe un abierto V que contiene a E y tiene medida finita. Por 3) y 5) existe un compacto $K \subset V$ con $\mu(V \setminus K) < \epsilon$. Como $E \cap K \in \mathcal{M}_F$, existe un compacto $H \subset E \cap K$ con $\mu(E \cap K) < \mu(H) + \epsilon$.

Puesto que $E \subset (E \cap K) \cup (V \setminus K)$, resulta

$$\mu(E) \leq \mu(E \cap K) + \mu(V \setminus K) < \mu(H) + 2\epsilon,$$

luego $E \in \mathcal{M}_F$.

Tras estas comprobaciones ya estamos en condiciones de probar el teorema. Consideremos la restricción de μ a la σ -álgebra de Borel. Los pasos 4) y 8) justifican que esta restricción es una medida. Hemos probado que μ es finita sobre los compactos, por definición se aproxima por abiertos y por 8) se aproxima por

compactos en los conjuntos de medida finita. El argumento de 3) prueba de hecho que los abiertos de medida infinita contienen compactos de medida arbitrariamente grande, luego μ se aproxima por compactos en todos los abiertos. Así pues μ es casi regular.

Falta probar que μ representa a T . Dada $f \in C_c(X)$, basta probar que

$$T(f) \leq \int_X f d\mu,$$

pues aplicando esto mismo a $-f$ obtenemos la desigualdad opuesta. Sea K el soporte de f . Entonces $f[X] \subset f[K] \cup \{0\}$ es compacto, $f[X] \subset [a, b]$, para ciertos números reales a y b . Sea $\epsilon > 0$ y tomemos números

$$y_0 < a < y_1 < \cdots < y_n = b$$

tales que $y_{i+1} - y_i < \epsilon$.

Sea $E_i = \{x \in X \mid y_{i-1} < f(x) \leq y_i\} \cap K$. Como f es continua, f es medible respecto al álgebra de Borel, luego los conjuntos E_i son conjuntos de Borel disjuntos cuya unión es K . Existen abiertos V_i tales que $E_i \subset V_i$,

$$\mu(V_i) < \mu(E_i) + \frac{\epsilon}{n}$$

y $f(x) < y_i + \epsilon$ para todo $x \in V_i$. Por el teorema 3.17 existen funciones $h_i \prec V_i$ que suman 1 sobre K . Por lo tanto $f = h_1 f + \cdots + h_n f$ y de 2) se sigue que

$$\mu(K) \leq T(h_1 + \cdots + h_n) = T(h_1) + \cdots + T(h_n).$$

Como $h_i f \leq (y_i + \epsilon)h_i$ e $y_i - \epsilon < f(x)$ en E_i , tenemos que

$$\begin{aligned} T(f) &= \sum_{i=1}^n T(h_i f) \leq \sum_{i=1}^n (y_i + \epsilon)T(h_i) \\ &= \sum_{i=1}^n (|a| + y_i + \epsilon)T(h_i) - |a| \sum_{i=1}^n T(h_i) \\ &\leq \sum_{i=1}^n (|a| + y_i + \epsilon)(\mu(E_i) + \frac{\epsilon}{n}) - |a|\mu(K) \\ &= \sum_{i=1}^n (y_i - \epsilon)\mu(E_i) + 2\epsilon\mu(K) + \frac{\epsilon}{n} \sum_{i=1}^n (|a| + y_i + \epsilon) \\ &\leq \int_X f d\mu + \epsilon(2\mu(K) + |a| + b + \epsilon). \end{aligned}$$

Como ϵ es arbitrario, tenemos la desigualdad buscada. \blacksquare

Recordemos que el teorema 9.12 implica que si el espacio X es σ -compacto entonces las medidas que proporciona el teorema de Riesz son regulares. Una primera aplicación interesante del teorema de Riesz nos da que en espacios razonables (como \mathbb{R}^n) toda medida razonable es regular:

Teorema 9.14 *Sea X un espacio localmente compacto en el que todo abierto sea σ -compacto. Si μ es una medida de Borel en X tal que todo compacto tiene medida finita, entonces es regular.*

DEMOSTRACIÓN: Definamos el operador $T : C_c(X) \rightarrow \mathbb{R}$ definido por $T(f) = \int_X f d\mu$. Notemos que si f tiene soporte K y M es una cota de f en K entonces

$$|T(f)| \leq \int_X |f| d\mu \leq M\mu(K) < +\infty,$$

luego T está bien definido y claramente cumple las hipótesis del teorema de Riesz. Por lo tanto existe una medida de Borel regular ν tal que para toda función $f \in C_c(X)$ se cumple

$$\int_X f d\nu = \int_X f d\mu.$$

Basta probar que $\mu = \nu$. Si V es un abierto, por hipótesis lo podemos expresar como unión de compactos $\{K_n\}_{n=1}^\infty$. Tomemos funciones $K_n \prec f_n \prec V$ y sea $g_n = \max\{f_1, \dots, f_n\}$. Entonces $g_n \in C_c(X)$ y es claro que la sucesión $\{g_n\}_{n=1}^\infty$ es monótona creciente y converge puntualmente a χ_V . Por el teorema de la convergencia monótona resulta que

$$\nu(V) = \lim_n \int_X g_n d\nu = \lim_n \int_X g_n d\mu = \mu(V).$$

Así pues, μ y ν coinciden en los abiertos.

Dado un conjunto de Borel B , cortándolo con una familia creciente de compactos cuya unión sea X podemos expresarlo como unión creciente de conjuntos de Borel de medida finita para μ y ν , luego basta probar que ambas medidas coinciden sobre los conjuntos de Borel de medida finita. Por la regularidad de ν es fácil ver que existen un abierto V y un compacto K tales que $K \subset B \subset V$ y $\nu(V \setminus K) < \epsilon$, para un ϵ prefijado. Como $V \setminus K$ es abierto, esta última igualdad vale también para μ .

Por consiguiente:

$$\begin{aligned} \mu(B) &\leq \mu(V) = \nu(V) \leq \nu(B) + \epsilon, \\ \nu(B) &\leq \nu(V) = \mu(V) \leq \nu(B) + \epsilon, \end{aligned}$$

luego $|\mu(B) - \nu(B)| \leq \epsilon$ para todo $\epsilon > 0$. ■

En particular, esto implica que la medida de Lebesgue es regular, pues todo compacto en \mathbb{R}^n está contenido en una celda, luego tiene medida de Lebesgue finita.

Terminamos esta sección con un teorema importante sobre aproximación de funciones medibles por funciones continuas:

Teorema 9.15 (Teorema de Lusin) *Sea μ una medida de Borel regular en un espacio localmente compacto X y $f : X \rightarrow \mathbb{R}$ una función medible tal que $f = 0$ salvo en un conjunto de medida finita. Dado $\epsilon > 0$ existe una función $g \in C_c(X)$ tal que $f = g$ salvo en un conjunto de medida menor que ϵ . Además podemos exigir que*

$$\sup_{x \in X} |g(x)| \leq \sup_{x \in X} |f(x)|.$$

DEMOSTRACIÓN: Sea $A = \{x \in X \mid f(x) \neq 0\}$. Supongamos primero que A es compacto y que $0 \leq f \leq 1$. Sea $\{s_n\}_{n=1}^\infty$ la sucesión monótona creciente de funciones simples construida en el teorema 8.44. Definimos $t_1 = s_1$ y $t_n = s_n - s_{n-1}$ para $n > 1$. Entonces $s_n(x) = k_n(f(x))/2^n$ y es fácil ver que $k_n(f(x)) = 2k_{n-1}(f(x))$ o bien $k_n(f(x)) = 2k_{n-1}(f(x)) + 1$, con lo que $2^n t_n$ toma sólo los valores 0 y 1. En otros términos $s_n = \chi_{T_n}$, para ciertos conjuntos medibles $T_n \subset A$. Además

$$f(x) = \sum_{n=1}^{\infty} t_n(x), \quad \text{para todo } x \in X.$$

Sea V un abierto de clausura compacta que contenga a A . Por regularidad existen compactos K_n y abiertos V_n de manera que $K_n \subset T_n \subset V_n \subset V$ y $\mu(V_n \setminus K_n) < 2^{-n}\epsilon$. Sea $K_n \prec h_n \prec V_n$. Definimos

$$g(x) = \sum_{n=1}^{\infty} 2^{-n} h_n(x).$$

Por el criterio de Weierstrass tenemos que g es continua. Además su soporte está contenido en la clausura de V , luego es compacto. Como $2^{-n} h_n(x) = t_n(x)$ excepto en $V_n \setminus K_n$, tenemos que $f = g$ excepto en $\bigcup_n (V_n \setminus K_n)$, que es un conjunto de medida menor que ϵ .

Del caso anterior se deduce el teorema para el caso en que A es compacto y f está acotada. Para probar el caso general tomamos $B_n = \{x \in X \mid |f(x)| > n\}$. Estos conjuntos forman una familia decreciente con intersección vacía y todos tienen medida finita, luego $\mu(B_n)$ tiende a 0 con n . La función f coincide con la función acotada $h = (1 - \chi_{B_n})f$ salvo en B_n y, por otra parte, podemos tomar un compacto $K \subset A$ tal que $\mu(A \setminus K)$ sea arbitrariamente pequeño, luego basta aproximar $h\chi_K$ por una función continua, lo cual es posible por la parte ya probada.

Por último, sea K el supremo de $|f|$. Si K es finito consideramos la función $h: \mathbb{R} \rightarrow \mathbb{R}$ dada por

$$h(x) = \begin{cases} x & \text{si } |x| \leq K \\ \frac{Kx}{|x|} & \text{si } |x| > K \end{cases}$$

Claramente h es continua, $g_1 = g \circ h$ sigue cumpliendo el teorema y además su supremo no excede al de f . ■

9.3 Espacios L^p

En las secciones posteriores vamos a necesitar algunos resultados abstractos referentes a ciertos espacios de funciones integrables que estudiaremos aquí.

Diremos que dos números reales positivos p y q son *conjugados* si

$$\frac{1}{p} + \frac{1}{q} = 1.$$

Es obvio que cada $p > 1$ tiene un único conjugado $q > 1$. El caso $p = q = 2$ es especialmente importante. Como los pares de conjugados aparecen normalmente como exponentes, es frecuente llamarlos “exponentes conjugados”.

Teorema 9.16 (Desigualdad de Hölder) Sean p y q exponentes conjugados. Sea X un espacio medida y sean $f, g : X \rightarrow [0, +\infty]$ funciones medibles. Entonces

$$\int_X fg \, d\mu \leq \left(\int_X f^p \, d\mu \right)^{1/p} \left(\int_X g^q \, d\mu \right)^{1/q}.$$

DEMOSTRACIÓN: Llamemos A y B a los dos factores del segundo miembro. Si $A = 0$ entonces el teorema 8.56 implica que $f = 0$ p.c.t.p., luego $fg = 0$ p.c.t.p. y la desigualdad es clara. Si $A > 0$ y $B = +\infty$ de nuevo es obvio. Podemos suponer, pues, $0 < A < +\infty$ y $0 < B < +\infty$. Llamemos $F = f/A$, $G = g/B$. Entonces

$$\int_X F^p \, d\mu = \int_X G^q \, d\mu = 1. \quad (9.1)$$

Supongamos que $0 < F(x) < +\infty$, $0 < G(x) < +\infty$. Entonces existen números s y t tales que $F(x) = e^{s/p}$, $G(x) = e^{t/q}$. Como $1/p + 1/q = 1$ y la función exponencial es convexa, concluimos que

$$e^{s/p+t/q} \leq p^{-1}e^s + q^{-1}e^t.$$

Por consiguiente,

$$F(x)G(x) \leq p^{-1}F(x)^p + q^{-1}G(x)^q.$$

Esta desigualdad es trivialmente cierta si $G(x) = 0$ o $F(x) = 0$, luego vale para todo x . Integrando y usando (9.1) resulta

$$\int_X FG \, d\mu \leq p^{-1} + q^{-1} = 1.$$

De aquí se sigue inmediatamente la desigualdad de Hölder. ■

Teorema 9.17 (Desigualdad de Minkowski) Sea X un espacio medida y $f, g : X \rightarrow [0, +\infty]$ funciones medibles. Para todo $p \geq 1$ se cumple

$$\left(\int_X (f+g)^p \, d\mu \right)^{1/p} \leq \left(\int_X f^p \, d\mu \right)^{1/p} + \left(\int_X g^p \, d\mu \right)^{1/p}$$

DEMOSTRACIÓN: Podemos suponer que $p > 1$, que el primer miembro es mayor que 0 y que el segundo es menor que $+\infty$. Como la función x^p es convexa en $]0, +\infty[$ tenemos que

$$\left(\frac{f+g}{2} \right)^p \leq \frac{1}{2}(f^p + g^p),$$

con lo que el primer miembro también es finito, es decir, las tres integrales son finitas. Sea q el conjugado de p . Escribimos

$$(f+g)^p = f(f+g)^{p-1} + g(f+g)^{p-1}.$$

Aplicamos la desigualdad de Hölder junto con que $(p-1)q = p$ (porque p y q son conjugados). El resultado es

$$\int_X f(f+g)^{p-1} d\mu \leq \left(\int_X f^p d\mu \right)^{1/p} \left(\int_X (f+g)^p d\mu \right)^{1/q}.$$

Intercambiando los papeles de f y g y sumando las desigualdades resulta

$$\int_X (f+g)^p d\mu \leq \left(\int_X (f+g)^p d\mu \right)^{1/q} \left(\left(\int_X f^p d\mu \right)^{1/p} + \left(\int_X g^p d\mu \right)^{1/p} \right).$$

Dividiendo entre el primer factor del segundo miembro y teniendo en cuenta que $1 - 1/q = 1/p$ tenemos la desigualdad buscada. ■

Definición 9.18 Sea X un espacio medida, $1 \leq p < +\infty$ y $f : X \rightarrow \mathbb{R}$ una función medible. Sea

$$\|f\|_p = \left(\int_X |f|^p d\mu \right)^{1/p}.$$

Llamaremos $L^p(\mu)$ al conjunto de todas las funciones medibles f tales que $\|f\|_p < +\infty$. Notemos que $L^1(\mu)$ es el conjunto de todas las funciones integrables, tal y como ya lo teníamos definido.

Por la desigualdad de Minkowski tenemos que si $f, g \in L^p(\mu)$ entonces $\|f+g\|_p \leq \|f\|_p + \|g\|_p$. En particular $f+g \in L^p(\mu)$. Por otra parte es claro que $\|\alpha f\|_p = |\alpha| \|f\|_p$, con lo que $\alpha f \in L^p(\mu)$. En particular vemos que $L^p(\mu)$ es un espacio vectorial sobre \mathbb{R} .

No es cierto que $\|\cdot\|_p$ sea una norma en $L^p(\mu)$, porque existen funciones no nulas f tales que $\|f\|_p = 0$ (las que son nulas p.c.t.p.). Ahora bien, es claro que las funciones de "norma" nula forman un subespacio vectorial de $L^p(\mu)$. Usaremos también la notación $L^p(\mu)$ para referirnos al espacio vectorial cociente. Si dos funciones f y g están en la misma clase entonces $f = g + h$, donde $\|h\|_p = 0$, luego $\|f\|_p \leq \|g\|_p + 0$ e igualmente tenemos la desigualdad contraria, luego $\|f\|_p = \|g\|_p$.

Podemos definir la norma de una clase de funciones como la norma de cualquiera de sus miembros. Al considerar clases de equivalencia si tenemos un espacio normado, pues las funciones de norma 0 forman una única clase.

Teorema 9.19 Sea X un espacio medida y $1 \leq p < +\infty$. Entonces $L^p(\mu)$ es un espacio de Banach.

DEMOSTRACIÓN: Sea $\{f_n\}_{n=1}^\infty$ una sucesión de Cauchy en $L^p(\mu)$. Basta probar que tiene una subsucesión convergente. Extrayendo una subsucesión podemos suponer que $\|f_{n+1} - f_n\| < 2^{-n}$. Sea

$$g_k = \sum_{n=1}^k |f_{n+1} - f_n|, \quad g = \sum_{n=1}^{\infty} |f_{n+1} - f_n|.$$

Claramente $\|g_k\|_p < 1$ y aplicando el lema de Fatou a $\{g_k^p\}_{k=1}^\infty$ concluimos que $\|g\|_p \leq 1$. En particular $g(x) < +\infty$ p.c.t.x. Así pues, la serie

$$f(x) = f_1(x) + \sum_{n=1}^{\infty} (f_{n+1}(x) - f_n(x))$$

converge absolutamente p.c.t.x. Definamos $f(x) = 0$ en los puntos donde no converja. Teniendo en cuenta quiénes son las sumas parciales de la serie, es claro que

$$f(x) = \lim_n f_n(x) \quad \text{p.c.t.x.}$$

Veamos que $f \in L^p(\mu)$ y que es el límite para la norma de la sucesión dada. Dado $\epsilon > 0$ existe un k tal que si $m, n > k$ entonces $\|f_n - f_m\|_p < \epsilon$. Por el lema de Fatou tenemos

$$\int_X |f - f_m|^p d\mu \leq \liminf_n \int_X |f_n - f_m|^p d\mu \leq \epsilon^p.$$

Esto significa que $\|f - f_m\|_p \leq \epsilon$, de donde $\|f\|_p \leq \|f_k\|_p + \epsilon$ y por lo tanto $f \in L^p(\mu)$. También es claro ahora que f es el límite en $L^p(\mu)$ de la sucesión dada. ■

En la prueba del teorema anterior hemos visto lo siguiente:

Teorema 9.20 *Toda sucesión que converge en un espacio $L^p(\mu)$ a una función f , tiene una subsucesión que converge puntualmente a f salvo en un conjunto nulo.*

En el caso de los espacios $L^2(\mu)$ todavía podemos decir más:

Teorema 9.21 *Sea X un espacio medida. Entonces $L^2(\mu)$ es un espacio de Hilbert con el producto escalar dado por*

$$\langle f, g \rangle = \int_X fg d\mu.$$

DEMOSTRACIÓN: La integral que define el producto escalar es finita, pues por la desigualdad de Hölder cumple en realidad que $|fg| \leq \|f\|_2 \|g\|_2$. Claramente es bilineal y la norma que induce es precisamente la de $L^2(\mu)$. ■

Ejercicio: Sea μ la medida en $\{1, \dots, n\}$ en la que cada punto tiene medida 1. Probar que $L^p(\mu) = \mathbb{R}^n$ y que las normas $\|\cdot\|_1$ y $\|\cdot\|_2$ son las definidas en el capítulo II.

Veamos ahora una generalización de la desigualdad de Hölder:

Teorema 9.22 (Desigualdad de Hölder generalizada) *Sea X un espacio medida, consideremos números positivos tales que $\sum_{i=1}^n \frac{1}{p_i} = 1$ así como funciones*

$f_i \in L^{p_i}(\mu)$. Entonces $\prod_{i=1}^n f_i \in L^1(\mu)$ y

$$\left\| \prod_{i=1}^n f_i \right\|_1 \leq \prod_{i=1}^n \|f_i\|_{p_i}.$$

DEMOSTRACIÓN: El caso $n = 2$ es consecuencia inmediata de la desigualdad de Hölder que ya hemos probado, pues si $f_1, f_2 \in L^2(\mu)$, podemos aplicarla a $|f_1|$ y $|f_2|$. Ahora razonamos por inducción sobre n . Sean

$$q_n = \frac{p_n}{p_n - 1}, \quad r_i = p_i \left(1 - \frac{1}{p_n}\right),$$

de modo que

$$\frac{1}{p_n} + \frac{1}{q_n} = 1, \quad \sum_{i=1}^{n-1} \frac{1}{r_i} = 1, \quad q_n r_i = p_i.$$

Como $f_i \in L^{p_i}(\mu)$, es decir, que $f_i^{p_i} = (f_i^{q_n})^{r_i}$ es integrable, tenemos que $f_i^{q_n} \in L^{r_i}(\mu)$, y además

$$\|f_i^{q_n}\|_{r_i}^{1/q_n} = \left(\int_X |f_i|^{q_n r_i}\right)^{1/q_n r_i} = \left(\int_X |f_i|^{p_i}\right)^{1/p_i} = \|f_i\|_{p_i}.$$

Por hipótesis de inducción tenemos que $\prod_{i=1}^{n-1} f_i^{q_n} \in L^1(\mu)$ o, equivalentemente,

que $\prod_{i=1}^{n-1} f_i \in L^{q_n}(\mu)$, y además

$$\left\| \prod_{i=1}^{n-1} f_i \right\|_{q_n} = \left(\int_X \prod_{i=1}^{n-1} |f_i|^{q_n} \right)^{1/q_n} = \left\| \prod_{i=1}^{n-1} f_i^{q_n} \right\|_1^{1/q_n} \leq \prod_{i=1}^{n-1} \|f_i^{q_n}\|_{r_i}^{1/q_n} = \prod_{i=1}^{n-1} \|f_i\|_{p_i}.$$

Ahora aplicamos el caso $n = 2$ a este producto y a f_n , con lo que

$$\left\| \prod_{i=1}^n f_i \right\|_1 \leq \left\| \prod_{i=1}^{n-1} f_i \right\|_{q_n} \|f_n\|_{p_n} \leq \prod_{i=1}^n \|f_i\|_{p_i}. \quad \blacksquare$$

Terminamos el estudio de los espacios $L^p(\mu)$ con dos teoremas de densidad:

Teorema 9.23 *Sea X un espacio medida. Sea S la clase de las funciones simples que son nulas salvo en un conjunto de medida finita. Entonces S es un subconjunto denso de $L^p(\mu)$ para $1 \leq p < +\infty$.*

DEMOSTRACIÓN: Es claro que $S \subset L^p(\mu)$. Tomemos primero una función $f \geq 0$ en $L^p(\mu)$. Sea $\{s_n\}_{n=1}^\infty$ una sucesión monótona creciente de funciones simples que converja a f . Como $0 \leq s_n \leq f$ es claro que $s_n \in L^p(\mu)$, luego $s_n \in S$. Como $|f - s_n|^p \leq f^p$, el teorema de la convergencia dominada implica que $\|f - s_n\|_p$ converge a 0, luego f está en la clausura de S . Para el caso general aplicamos la parte ya probada a f^+ y f^- . \blacksquare

Teorema 9.24 *Sea μ una medida de Borel regular en un espacio localmente compacto X . Entonces $C_c(X)$ es denso en $L^p(\mu)$ para $1 \leq p < +\infty$.*

DEMOSTRACIÓN: Consideremos la clase S del teorema anterior. Basta ver que toda función $s \in S$ puede aproximarse por una función de $C_c(X)$. Sea $\epsilon > 0$ y K el supremo de s , que claramente es finito. Por el teorema de Lusin existe una función $g \in C_c(X)$ tal que $g = s$ salvo en un conjunto de medida menor que ϵ . Por consiguiente $\|g - s\|_p \leq 2K\epsilon^{1/p}$. \blacksquare

9.4 Medidas signadas

Para enunciar más adecuadamente los próximos resultados conviene que modifiquemos nuestra definición de medida o, con más exactitud, que introduzcamos otro tipo de medidas distintas de las medidas positivas. Aunque pronto veremos la utilidad del nuevo concepto desde un punto de vista puramente matemático, quizá ahora sea más conveniente motivarlo mediante un ejemplo físico: la función que a cada región del espacio le asigna la cantidad de materia que contiene es un ejemplo de medida positiva (que podemos suponer finita), sin embargo, la aplicación que a cada región del espacio le asigna la carga eléctrica que contiene ya no se ajusta a nuestra definición de medida, porque puede tomar valores negativos, y pese a ello puede tratarse de forma muy similar.

Definición 9.25 Sea \mathcal{A} una σ -álgebra en un conjunto X . Una *medida signada (finita)* en \mathcal{A} es una aplicación $\mu : \mathcal{A} \rightarrow \mathbb{R}$ tal que $\mu(\emptyset) = 0$ y si $\{E_n\}_{n=1}^{\infty}$ es una familia de conjuntos de \mathcal{A} disjuntos dos a dos entonces

$$\mu\left(\bigcup_{n=1}^{\infty} E_n\right) = \sum_{n=1}^{\infty} \mu(E_n).$$

Observemos que en la definición está implícita la hipótesis de que las series de medidas de conjuntos disjuntos son convergentes (en el caso de las medidas positivas donde admitíamos el valor $+\infty$ esto era evidente). Más aún, la serie ha de converger absolutamente. En efecto, la serie (finita o infinita) formada por los términos correspondientes a los conjuntos E_n con medida negativa ha de converger a la medida de su unión, y obviamente la serie de los valores absolutos converge al valor absoluto de la suma, es decir, los términos negativos convergen absolutamente. Lo mismo vale para los términos positivos, luego la serie completa también converge absolutamente.

Conviene saber que toda la teoría que vamos a exponer sobre medidas signadas se generaliza con cambios mínimos a medidas con valores complejos, pero no vamos a necesitar nada al respecto. Con esta definición, las medidas signadas sobre una σ -álgebra fija en un conjunto X forman un espacio vectorial real con las operaciones dadas por

$$(\mu + \nu)(E) = \mu(E) + \nu(E), \quad (\alpha\mu)(E) = \alpha\mu(E).$$

En particular, las medidas signadas de Borel en un espacio topológico X forman un espacio vectorial real.

Ejemplo Sea X un espacio topológico y $x \in X$. Definimos la *delta de Dirac* de soporte x como la medida de Borel dada por

$$\delta_x(E) = \begin{cases} 1 & \text{si } x \in E \\ 0 & \text{si } x \notin E \end{cases}.$$

Claramente se trata de una medida signada positiva. Si una región del espacio está ocupada por partículas puntuales en las posiciones x_1, \dots, x_n con

cargas eléctricas q_1, \dots, q_n , entonces la distribución de carga viene dada por la medida signada

$$\mu = \sum_{i=1}^n q_i \delta_{x_i}.$$

También podemos considerar la medida positiva

$$|\mu| = \sum_{i=1}^n |q_i| \delta_{x_i},$$

que a cada región del espacio le asigna la cantidad total de carga que contiene, haciendo abstracción de su signo. Vamos a probar que a toda medida signada μ le podemos asignar una medida positiva $|\mu|$ con una interpretación análoga a la de este ejemplo. ■

Definición 9.26 Sea μ una medida signada en un conjunto X . Llamaremos *variación total* de μ a la aplicación definida sobre la misma σ -álgebra que μ dada por

$$|\mu|(E) = \sup \sum_{n=1}^{\infty} |\mu(E_n)|,$$

donde el supremo se toma sobre todas las particiones $\{E_n\}_{n=1}^{\infty}$ de E en conjuntos medibles disjuntos dos a dos.

Tomando la partición formada únicamente por E obtenemos la relación $|\mu(E)| \leq |\mu|(E)$. Es inmediato comprobar que la medida $|\mu|$ construida en el ejemplo anterior es la variación total de μ en el sentido de la definición anterior.

Teorema 9.27 *La variación total de una medida compleja es una medida positiva finita.*

DEMOSTRACIÓN: Obviamente $|\mu|(\emptyset) = 0$. Sea $\{E_n\}_{n=1}^{\infty}$ una partición de un conjunto medible E en conjuntos medibles disjuntos dos a dos. Sea $r_n < |\mu|(E_n)$. Entonces cada E_n tiene una partición $\{E_{nm}\}_{m=1}^{\infty}$ de modo que

$$r_n < \sum_{m=1}^{\infty} |\mu(E_{nm})|.$$

La unión de todas las particiones forma una partición de E , con lo que

$$\sum_{n=1}^{\infty} r_n \leq \sum_{m,n=1}^{\infty} |\mu(E_{nm})| \leq |\mu|(E).$$

Tomando el supremo en todas las posibles elecciones de $\{r_n\}_{n=1}^{\infty}$ resulta que

$$\sum_{n=1}^{\infty} |\mu|(E_n) \leq |\mu|(E).$$

Sea ahora $\{A_m\}_{m=1}^{\infty}$ una partición de E en conjuntos medibles disjuntos dos a dos. Entonces $\{E_n \cap A_m\}_{n=1}^{\infty}$ es una partición de A_m y $\{E_n \cap A_m\}_{m=1}^{\infty}$ es una partición de E_n , luego

$$\begin{aligned} \sum_{m=1}^{\infty} |\mu(A_m)| &= \sum_{m=1}^{\infty} \left| \sum_{n=1}^{\infty} \mu(A_m \cap E_n) \right| \leq \sum_{m=1}^{\infty} \sum_{n=1}^{\infty} |\mu(A_m \cap E_n)| \\ &= \sum_{n=1}^{\infty} \sum_{m=1}^{\infty} |\mu(A_m \cap E_n)| \leq \sum_{n=1}^{\infty} |\mu(E_n)|. \end{aligned}$$

Como esto vale para toda partición de E , tenemos

$$|\mu|(E) \leq \sum_{n=1}^{\infty} |\mu|(E_n).$$

Falta probar que $|\mu|$ es finita. Supongamos que existe un conjunto medible E tal que $|\mu|(E) = +\infty$. Sea $t = 2(1 + |\mu|(E))$. Puesto que $|\mu|(E) > t$, por definición de variación total existen conjuntos medibles E_n contenidos en E y disjuntos dos a dos tales que

$$t < \sum_{n=1}^k |\mu(E_n)|.$$

Sea P la suma de los términos $|\mu(E_n)|$ tales que $\mu(E_n) \geq 0$ y sea N la suma de los $|\mu(E_n)|$ tales que $\mu(E_n) < 0$. Por la desigualdad anterior tenemos $t < P + N$, luego $t < 2P$ o bien $t < 2N$, según si $N \leq P$ o $P \leq N$. Sea A la unión de los E_n correspondientes a P o N según el caso, de modo que $A \subset E$ y $t < 2|\mu(A)|$, luego $|\mu(A)| > t/2 > 1$.

Sea ahora $B = E \setminus A$. Entonces

$$|\mu(B)| = |\mu(B) - \mu(A)| \geq |\mu(A)| - |\mu(E)| > \frac{t}{2} - |\mu(E)| = 1.$$

Así pues, hemos partido E en dos conjuntos disjuntos A y B tales que $|\mu(A)| > 1$ y $|\mu(B)| > 1$. Obviamente $|\mu|(A) = +\infty$ o bien $|\mu|(B) = +\infty$.

Supongamos ahora que el espacio total X tiene variación total infinita. Por el argumento anterior podemos partirlo en dos conjuntos medibles disjuntos $X = A_1 \cup B_1$ tales que $|\mu|(B_1) = +\infty$ y $|\mu(A_1)| > 1$. Aplicando el mismo razonamiento a B_1 obtenemos $B_1 = A_2 \cup B_2$ con $|\mu|(B_2) = +\infty$ y $|\mu(A_2)| > 1$. Procediendo de este modo construimos una familia numerable de conjuntos medibles disjuntos $\{A_n\}_{n=1}^{\infty}$ tales que $|\mu(A_n)| > 1$ para todo n . Debería cumplirse

$$\mu \left(\bigcup_{n=1}^{\infty} A_n \right) = \sum_{n=1}^{\infty} \mu(A_n),$$

pero la serie no converge, porque su término general no tiende a 0. Esta contradicción prueba que $|\mu|(X) < +\infty$ y por lo tanto $|\mu|$ es una medida finita. ■

Ejercicio: Probar que el conjunto de todas las medidas signadas sobre una σ -álgebra en un conjunto X es un espacio normado con la norma dada por $\|\mu\| = |\mu|(X)$.

Definición 9.28 Si μ es una medida signada en un conjunto X , llamaremos *variación positiva* y *variación negativa* de μ a las medidas (definidas sobre la misma σ -álgebra) dadas por

$$\mu^+ = \frac{|\mu| + \mu}{2}, \quad \mu^- = \frac{|\mu| - \mu}{2}.$$

Claramente son dos medidas positivas finitas y cumplen las relaciones

$$\mu = \mu^+ - \mu^-, \quad |\mu| = \mu^+ + \mu^-.$$

Por ejemplo, si μ representa la carga eléctrica contenida en una región del espacio, μ^+ y μ^- representan, respectivamente, la carga positiva y la carga negativa que contiene dicha región.

Diremos que una función $f : X \rightarrow \mathbb{R}$ es *integrable* respecto a una medida signada μ si lo es respecto a μ^+ y μ^- , y definimos su integral como

$$\int_X f d\mu = \int_X f d\mu^+ - \int_X f d\mu^-.$$

Es fácil ver que el conjunto $L^1(\mu)$ de las funciones integrables es un espacio vectorial y la integral determina sobre él una aplicación lineal. Además se cumple la desigualdad

$$\left| \int_X f d\mu \right| \leq \int_X |f| d|\mu|.$$

También es claro que la aplicación dada por $\nu(E) = \int_E f d\mu$ es una medida signada en X .

Nos encaminamos a probar ahora uno de los teoremas más importantes de la teoría de la medida. Para ello necesitamos algunos conceptos y resultados previos.

Definición 9.29 sea μ una medida positiva en un conjunto X y λ una medida arbitraria (positiva o signada) en la misma σ -álgebra. Diremos que λ es *absolutamente continua* respecto a μ , y lo representaremos por $\lambda \ll \mu$, si todos los conjuntos nulos para μ son nulos para λ .

Si existe un conjunto medible A tal que para todo conjunto medible E se cumple $\lambda(E) = \lambda(A \cap E)$ se dice que λ *está concentrada* en A . Esto equivale a que $\lambda(E) = 0$ siempre que $E \cap A = \emptyset$.

Diremos que dos medidas arbitrarias (sobre una misma σ -álgebra) son *mutuamente singulares*, y lo representaremos por $\lambda_1 \perp \lambda_2$, si existen conjuntos medibles disjuntos A y B tales que λ_1 está concentrada en A y λ_2 está concentrada en B .

Ejemplo Si admitimos como principio que toda masa ocupa un volumen, entonces la medida μ que a cada región del espacio le asigna la masa que contiene es absolutamente continua respecto a la medida de Lebesgue m . Por el contrario, a veces es más conveniente trabajar con masas puntuales, es decir, suponiéndolas localizadas en puntos del espacio sin volumen. Éste sería el caso de una distribución de masas de la forma $\mu = \sum_{n=1}^k m_n \delta_{x_n}$. Es fácil ver que entonces $\mu \perp m$. ■

He aquí algunas propiedades elementales:

Teorema 9.30 Sean λ , λ_1 y λ_2 medidas arbitrarias en un conjunto X y μ una medida positiva, todas ellas con los mismos conjuntos medibles. Entonces

- a) Si λ está concentrada en un conjunto medible A , también lo está $|\lambda|$.
- b) Si $\lambda_1 \perp \lambda_2$ entonces $|\lambda_1| \perp |\lambda_2|$.
- c) Si $\lambda_1 \perp \mu$ y $\lambda_2 \perp \mu$ entonces $\lambda_1 + \lambda_2 \perp \mu$.
- d) Si $\lambda_1 \ll \mu$ y $\lambda_2 \ll \mu$ entonces $\lambda_1 + \lambda_2 \ll \mu$.
- e) Si $\lambda \ll \mu$, entonces $|\lambda| \ll \mu$.
- f) Si $\lambda_1 \ll \mu$ y $\lambda_2 \perp \mu$ entonces $\lambda_1 \perp \lambda_2$.
- g) Si $\lambda \ll \mu$ y $\lambda \perp \mu$ entonces $\lambda = 0$.

DEMOSTRACIÓN: a) Si $E \cap A = \emptyset$ y $\{E_n\}_{n=1}^{\infty}$ es cualquier partición de E , entonces $\lambda(E_n) = 0$ para todo n , luego $|\lambda|(E) = 0$.

b) Es consecuencia inmediata de a).

c) Existen conjuntos disjuntos A_1 y B_1 tales que λ_1 está concentrada en A_1 y μ está concentrada en B_1 e igualmente existen conjuntos disjuntos A_2 y B_2 tales que λ_2 está concentrada en A_2 y μ está concentrada en B_2 . Entonces $\lambda_1 + \lambda_2$ está concentrada en $A = A_1 \cup A_2$ y μ está concentrada en $B = B_1 \cup B_2$.

d) Obvio.

e) Si $\mu(E) = 0$ y $\{E_n\}_{n=1}^{\infty}$ es una partición de E en conjuntos medibles disjuntos, entonces $\mu(E_n) = 0$ para todo n , luego $\lambda(E_n) = 0$ y $|\lambda|(E) = 0$.

f) Tenemos que λ_2 está concentrada en un conjunto A tal que $\mu(A) = 0$. Como $\lambda_1 \ll \mu$ ha de ser $\lambda_1(E) = 0$ para todo conjunto medible $E \subset A$. Por consiguiente λ_1 está concentrada en $X \setminus A$.

g) Por f) tenemos $\lambda \perp \mu$, pero esto implica que $\lambda = 0$. ■

Ahora probamos dos hechos elementales sobre medidas positivas que nos harán falta a continuación.

Teorema 9.31 Si μ es una medida positiva σ -finita en un conjunto X , entonces existe una función $w \in L^1(\mu)$ tal que $0 < w(x) < 1$ para todo $x \in X$.

DEMOSTRACIÓN: Sea $\{E_n\}_{n=1}^\infty$ una partición de X en conjuntos disjuntos de medida finita. Definamos

$$w_n(x) = \begin{cases} \frac{1}{2^n(1+\mu(E_n))} & \text{si } x \in E_n \\ 0 & \text{si } x \in X \setminus E_n \end{cases}$$

La función $w = \sum_{n=1}^\infty w_n$ cumple lo pedido. ■

Teorema 9.32 *Sea μ una medida positiva finita, $f \in L^1(\mu)$ y $C \subset \mathbb{R}$ un conjunto cerrado. Si*

$$P_E(f) = \frac{1}{\mu(E)} \int_E f d\mu \in C$$

para todo conjunto medible E no nulo, entonces $f(x) \in C$ p.c.t. $x \in X$.

DEMOSTRACIÓN: Sea $I = [x - r, x + r]$ un intervalo cerrado disjunto de C . Puesto que $\mathbb{R} \setminus C$ es unión de una familia numerable de tales intervalos, basta probar que $E = f^{-1}[I]$ es nulo. En caso contrario

$$|P_E(f) - x| = \frac{1}{\mu(E)} \left| \int_E (f - x) d\mu \right| \leq \frac{1}{\mu(E)} \int_E |f - x| d\mu \leq r,$$

lo cual es imposible, pues $P_E(f) \in C$. ■

Finalmente podemos probar:

Teorema 9.33 (de Lebesgue-Radon-Nikodým) *Sea μ una medida positiva σ -finita en un conjunto X y sea λ una medida signada sobre la misma σ -álgebra. Entonces*

a) *Existe un único par de medidas signadas λ_a y λ_s tales que*

$$\lambda = \lambda_a + \lambda_s, \quad \lambda_a \ll \mu, \quad \lambda_s \perp \mu.$$

Si λ es positiva (y finita) también lo son λ_a y λ_s .

b) *Existe una única $h \in L^1(\mu)$ tal que para todo conjunto medible E*

$$\lambda_a(E) = \int_E h d\mu.$$

La parte a) se conoce como Teorema de Lebesgue. La parte b) es el Teorema de Radon-Nikodým.

DEMOSTRACIÓN: La unicidad de a) es clara a partir de 9.30, pues si λ'_a, λ'_s es otro par que cumpla lo mismo, entonces $\lambda'_a - \lambda_a = \lambda_s - \lambda'_s$, $\lambda'_a - \lambda_a \ll \mu$ y $\lambda_s - \lambda'_s \perp \mu$, luego $\lambda'_a - \lambda_a = \lambda_s - \lambda'_s = 0$.

La unicidad de h en b) (como función de $L^1(\mu)$, es decir, p.c.t.p.) es fácil de probar: si existe otra h' en las mismas condiciones entonces $f = h - h'$ tiene integral nula sobre todo conjunto. Tomamos $E = \{x \in X \mid f(x) > 0\}$ y, como f tiene integral nula sobre E , por el teorema 8.56 concluimos que $f = 0$ p.c.t.p.

Supongamos primero que λ es positiva (y finita). Sea w según el teorema 9.31. Sea ν la medida positiva finita dada por

$$\nu(E) = \lambda(E) + \int_E w d\mu.$$

En otros términos, si $f = \chi_E$ se cumple

$$\int_X f d\nu = \int_X f d\lambda + \int_X fw d\mu.$$

Claramente, la misma relación vale cuando f es una función simple y, en consecuencia, para funciones medibles no negativas. Si $f \in L^2(\nu)$ la desigualdad de Hölder implica

$$\left| \int_X f d\lambda \right| \leq \int_X |f| d\lambda \leq \int_X |f| d\nu \leq \left(\int_X |f|^2 d\nu \right)^{1/2} \nu(X)^{1/2} = \nu(X)^{1/2} \|f\|_2.$$

Según el teorema 3.35, esto significa que

$$f \mapsto \int_X f d\lambda$$

es una aplicación lineal continua de $L^2(\nu)$ en \mathbb{R} . Como $L^2(\nu)$ es un espacio de Hilbert, el teorema 3.42 implica que existe $g \in L^2(\nu)$ tal que para toda $f \in L^2(\nu)$ se cumple

$$\int_X f d\lambda = \int_X fg d\nu. \quad (9.2)$$

Si aplicamos esto a una función característica $f = \chi_E$, donde $\nu(E) > 0$ el miembro izquierdo es $\lambda(E)$ y así

$$0 \leq \frac{1}{\nu(E)} \int_E g d\nu = \frac{\lambda(E)}{\nu(E)} \leq 1.$$

Según el teorema 9.32, resulta que $g(x) \in [0, 1]$ p.c.t.x (respecto a ν). Puesto que g sólo está determinada como elemento de $L^2(\nu)$, podemos modificarla en un conjunto nulo y suponer que $0 \leq g(x) \leq 1$ para todo $x \in X$. Entonces (9.2) puede reescribirse como

$$\int_X (1-g)f d\lambda = \int_X fgw d\mu. \quad (9.3)$$

Sean

$$A = \{x \in X \mid 0 \leq g(x) < 1\}, \quad B = \{x \in X \mid g(x) = 1\}.$$

Definimos las medidas λ_a y λ_s mediante

$$\lambda_a(E) = \lambda(A \cap E), \quad \lambda_s(E) = \lambda(B \cap E).$$

Haciendo $f = \chi_B$ en (9.3) el miembro izquierdo es 0 y el derecho es $\int_B w d\mu$, y como $w > 0$ concluimos que $\mu(B) = 0$, luego $\lambda_s \perp \mu$.

Ahora aplicamos (9.3) a $(1 + g + \dots + g^n)\chi_E$, con lo que tenemos

$$\int_E (1 - g^{n+1}) d\lambda = \int_E g(1 + g + \dots + g^n)w d\mu.$$

El integrando de la izquierda es nulo en B y converge a 1 en A , luego la integral de la izquierda converge a $\lambda(A \cap E) = \lambda_a(E)$. Por otra parte, el integrando de la derecha converge a una función medible no negativa h (quizá con valores infinitos), luego tomando límites en n resulta que

$$\lambda_a(E) = \int_E h d\mu.$$

En particular esto vale para $E = X$, lo que prueba que h toma valores finitos p.c.t.p., luego modificándola si es necesario en un conjunto nulo podemos suponer que $h \in L^1(\mu)$.

Esto prueba el teorema cuando λ es positiva. Si es una medida signada arbitraria basta aplicar la parte ya probada a λ^+ y λ^- . ■

Definición 9.34 Si μ una medida positiva σ -finita en un conjunto X y λ es una medida signada sobre la misma σ -álgebra tal que $\lambda \ll \mu$, entonces la función h cuya existencia afirma el teorema de Radon-Nikodým se llama la *derivada de Radon-Nikodým* de λ respecto a μ . La relación que expresa el teorema se representa también por $d\lambda = h d\mu$.

Veamos una aplicación del teorema de Radon-Nikodým. Una interpretación del teorema siguiente es que si μ representa la distribución de carga eléctrica en el espacio, entonces el espacio puede dividirse en dos regiones, una íntegramente ocupada por cargas positivas y otra por cargas negativas.

Teorema 9.35 (Teorema de descomposición de Hahn) *Sea μ una medida signada en un conjunto X . Entonces X se descompone en unión de dos conjuntos medibles disjuntos A y B tales que para todo conjunto medible E se cumple*

$$\mu^+(E) = \mu(A \cap E), \quad \mu^-(E) = -\mu(B \cap E).$$

DEMOSTRACIÓN: Obviamente $\mu \ll |\mu|$, luego existe una función $h \in L^1(|\mu|)$ tal que $d\mu = h d|\mu|$. Veamos que h toma los valores ± 1 p.c.t.p. Dado un número real r , sea $A_r = \{x \in X \mid |h(x)| < r\}$. Para toda partición $\{E_n\}_{n=1}^\infty$ de A_r se cumple

$$\sum_{n=1}^\infty |\mu(E_n)| = \sum_{n=1}^\infty \left| \int_{E_n} h d|\mu| \right| \leq \sum_{n=1}^\infty r|\mu|(E_n) = r|\mu|(A_r).$$

Para $r < 1$ esto implica que $|\mu|(A_r) = 0$, luego $|h| \geq 1$ p.c.t.p. Por otra parte, si $|\mu|(E) > 0$ tenemos que

$$\left| \frac{1}{|\mu|(E)} \int_E h d|\mu| \right| = \frac{|\mu(E)|}{|\mu|(E)} \leq 1,$$

luego el teorema 9.32 implica que $|h| \leq 1$ p.c.t.p. Modificando h en un conjunto nulo podemos suponer que $h = \pm 1$. Ahora basta definir

$$A = \{x \in X \mid h(x) = 1\}, \quad B = \{x \in X \mid h(x) = -1\}.$$

En efecto, por definición de μ^+ tenemos que para todo conjunto medible E se cumple

$$\mu^+(E) = \frac{1}{2} \int_E (1+h) d|\mu| = \int_{E \cap A} h d|\mu| = \mu(E \cap A).$$

Igualmente se razona con la variación negativa. ■

Se dice que el par (A, B) es una *partición de Hahn* de μ . Como aplicación obtenemos un par de hechos de interés sobre la derivada de Radon-Nikodým de una medida signada:

Teorema 9.36 *Sea μ una medida positiva σ -finita en un conjunto X y sea $f \in L^1(\mu)$. Sea λ la medida signada determinada por $d\lambda = f d\mu$. Entonces $|\lambda| = |f| d\mu$ y si $g \in L^1(\lambda)$ entonces $gf \in L^1(\mu)$ y*

$$\int_X g d\lambda = \int_X gf d\mu.$$

DEMOSTRACIÓN: Claramente $\lambda \ll \mu$, luego $|\lambda| \ll \mu$, luego $\lambda^+ \ll \mu$ y $\lambda^- \ll \mu$. Por el teorema de Radon-Nikodým existen funciones $f_+, f_- \in L^1(\mu)$ tales que $d\lambda^+ = f_+ d\mu$, $d\lambda^- = f_- d\mu$. Si (A, B) es una partición de Hahn para λ , podemos exigir que f_+ se anule en B y f_- se anule en A . Es claro que $f = f_+ - f_-$ p.c.t.p., luego $|f| = f_+ + f_-$ p.c.t.p. Por consiguiente $d|\lambda| = (f_+ + f_-) d\mu = |f| d\mu$.

La segunda parte del teorema es obvia si g es una función simple. Si g es positiva tomamos una sucesión creciente $\{s_n\}$ de funciones simples $0 \leq s_n \leq g$ que converja puntualmente a g . Entonces

$$\int_X s_n |f| d\mu = \int_X s_n d|\lambda|.$$

Por el teorema de la convergencia monótona concluimos que

$$\int_X g |f| d\mu = \int_X g d|\lambda| < +\infty,$$

luego $gf \in L^1(\mu)$. También tenemos

$$\int_X s_n d\lambda^+ - \int_X s_n d\lambda^- = \int_X s_n f d\mu$$

Aplicando el teorema de la convergencia monótona a la izquierda y el de la convergencia dominada a la derecha llegamos a la igualdad del enunciado. Si g no es positiva aplicamos la parte ya probada a g^+ y g^- . ■

Para terminar probaremos una versión del teorema de Riesz para medidas signadas. Sea K un espacio topológico compacto y $C(K)$ el espacio de todas las funciones reales continuas sobre K . Sabemos que $C(K)$ es un espacio de Banach con la norma supremo. Sea $C(K)'$ el espacio de las aplicaciones lineales continuas de $C(K)$ en \mathbb{R} , que también es un espacio de Banach con la norma

$$\|T\| = \sup\{|T(f)| \mid \|f\|_\infty \leq 1\}.$$

Además, para toda $f \in C(K)$ se cumple $|T(f)| \leq \|T\| \|f\|_\infty$. Llamemos $M(K)$ al conjunto de todas las medidas signadas de Borel en K , que claramente es un espacio normado con la norma $\|\mu\| = |\mu|(K)$.

Teorema 9.37 (Teorema de representación de Riesz) *Si K es un espacio compacto, a cada funcional lineal continuo $T \in C(K)'$ le corresponde una única medida signada $\mu \in M(K)$ tal que para toda función $f \in C(K)$ se cumple*

$$T(f) = \int_K f d\mu.$$

Además esta correspondencia es una isometría $C(K)' \rightarrow M(K)$.

DEMOSTRACIÓN: Si T fuera positivo, es decir, si $T(f) \geq 0$ cuando $f \geq 0$, la versión del teorema de Riesz que probamos en el capítulo anterior nos daría la medida que buscamos. En el caso general vamos a descomponer T en diferencia de dos funcionales positivos. Sea $C^+(K) = \{f \in C(K) \mid f \geq 0\}$ y definamos

$$T^+(f) = \sup\{T(u) \mid u \in C^+(K), u \leq f\}, \quad \text{para } f \in C^+(K).$$

Notemos que si $0 \leq u \leq f$ se cumple $|T(u)| \leq \|T\| \|u\|_\infty \leq \|T\| \|f\|_\infty$, luego T^+ es finito y $|T^+(f)| \leq \|T\| \|f\|_\infty$.

Es claro que si $\alpha \geq 0$ se cumple $T^+(\alpha f) = \alpha T^+(f)$. Además $T^+(f+g) = T^+(f) + T^+(g)$. En efecto, si $0 \leq u \leq f$ y $0 \leq v \leq g$ entonces $0 \leq u+v \leq f+g$, luego $T(u)+T(v) \leq T^+(f+g)$. Tomando supremos $T^+(f)+T^+(g) \leq T^+(f+g)$. Recíprocamente, si $w \leq f+g$ es claro que $u = f \wedge w$ y $v = w - u$ son funciones continuas y $0 \leq u \leq f$, $0 \leq v \leq g$, $w = u+v$, luego $T(w) \leq T^+(f) + T^+(g)$ y, tomando supremos, $T^+(f+g) \leq T^+(f) + T^+(g)$.

Dada $f \in C(K)$ definimos $T^+(f) = T^+(f^+) - T^+(f^-)$. Es fácil probar que $T^+ : C(K) \rightarrow \mathbb{R}$ es un funcional lineal continuo positivo. Lo mismo vale para $T^- = T^+ - T$. Por el teorema de Riesz para funcionales positivos existen medidas positivas μ_+ y μ_- tales que

$$T^+(f) = \int_K f d\mu_+, \quad T^-(f) = \int_K f d\mu_-.$$

Ambas medidas son finitas (basta aplicar las fórmulas a $f = 1$). Por lo tanto podemos definir $\mu = \mu_+ - \mu_-$, que es una medida signada en K y claramente representa a T .

Veamos que $\|\mu\| = \|T\|$. Si $\|f\|_\infty \leq 1$ tenemos

$$|T(f)| = \left| \int_K f d\mu \right| \leq \int_K |f| d|\mu| \leq |\mu|(K) = \|\mu\|,$$

luego $\|T\| \leq \|\mu\|$. Dado $\epsilon > 0$ consideramos una partición de Hahn (A, B) para μ . Sean K_1 y K_2 conjuntos compactos tales que $K_1 \subset A$, $K_2 \subset B$, $\mu^+(A) - \mu^+(K_1) + \mu^-(B) - \mu^-(K_2) < \epsilon$. Es fácil construir una función continua $f : X \rightarrow [-1, 1]$ que valga 1 sobre K_1 y -1 sobre K_2 . Entonces

$$\|\mu\| = \int_A d\mu^+ + \int_B d\mu^- \leq \int_K f d\mu + \epsilon = T(f) + \epsilon \leq \|T\| + \epsilon,$$

luego $\|\mu\| \leq \|T\|$. En particular tenemos que la correspondencia $T \mapsto \mu$ es inyectiva (su núcleo es trivial) y obviamente es suprayectiva. ■

9.5 Derivación de medidas

Supongamos que μ representa la distribución de la masa en el espacio y m es la medida de Lebesgue. Si suponemos que $\mu \ll m$, es decir, si la materia ocupa un volumen, entonces tiene sentido hablar de la densidad de materia en un punto x del espacio, entendida como la cantidad de materia por unidad de volumen. Una aproximación a dicha densidad es el cociente

$$\frac{\mu(B)}{m(B)},$$

donde B es un entorno de x . Sin embargo, si la distribución de la materia no es uniforme, dicho cociente no es exactamente la densidad, pero se parecerá más a ella cuanto menor sea el entorno considerado. Estas ideas nos llevan a la definición siguiente:

Definición 9.38 Sea m la medida de Lebesgue en \mathbb{R}^n y μ una medida de Borel signada en \mathbb{R}^n . Definimos los cocientes

$$C_r \mu(x) = \frac{\mu(B_r(x))}{m(B_r(x))},$$

donde las bolas las tomamos respecto a la distancia euclídea. Definimos la *derivada* de μ en x como

$$\frac{d\mu}{dm}(x) = \lim_{r \rightarrow 0} C_r \mu(x).$$

Probaremos que si $\mu \ll m$ entonces la derivada existe p.c.t.p. Para ello nos apoyaremos en la *función maximal* M , definida por

$$M\mu(x) = \sup_{0 < r < +\infty} C_r |\mu|(x).$$

Veamos que es medible usando para ello el teorema 8.41. Sea E la antiimagen por $M\mu$ de un intervalo $]t, +\infty]$, es decir, $E = \{x \in X \mid M\mu(x) > t\}$. Veamos que es abierto. Dado $x \in E$, existe un $r > 0$ tal que $u = C_r|\mu|(x) > t$, luego $\mu(B_r(x)) = um(B_r(x))$. Tomemos $\delta > 0$ tal que $(r + \delta)^n < r^n u/t$. Así, si $|y - x| < \delta$ entonces $B_r(x) \subset B_{r+\delta}(y)$, con lo que

$$|\mu|(B_{r+\delta}(y)) \geq um(B_r(x)) = u \left(\frac{r}{r+\delta} \right)^n m(B_{r+\delta}(y)) > tm(B_{r+\delta}(y)).$$

Esto prueba que $y \in E$, es decir, tenemos que $B_\delta(x) \subset E$, luego E es abierto. ■

Ahora necesitamos un par de hechos técnicos:

Teorema 9.39 *Sea W la unión de una familia finita de bolas $B_{r_i}(x_i) \subset \mathbb{R}^n$, para $i = 1, \dots, N$. Entonces existe un conjunto $S \subset \{1, \dots, n\}$ tal que:*

- a) *Las bolas $B_{r_i}(x_i)$ con $i \in S$ son disjuntas,*
- b) $W \subset \bigcup_{i \in S} B_{3r_i}(x_i)$,
- c) $m(W) \leq 3^n \sum_{i \in S} m(B_{r_i}(x_i))$.

DEMOSTRACIÓN: Escribiremos $B_i = B_{r_i}(x_i)$. Ordenemos las bolas de modo que sus radios sean decrecientes. Sea $i_1 = 1$. Eliminemos todas las bolas que corten a B_{i_1} . Sea B_{i_2} la primera bola restante, si es que queda alguna, eliminemos las bolas que cortan a B_{i_2} y continuemos el proceso hasta que no queden bolas. Veamos que las bolas que hemos dejado cumplen el teorema. Ciertamente son disjuntas. Cada bola B_j de las que hemos eliminado está contenida en una bola $B_{3r_i}(x_i)$, para algún $i \in S$, pues si $r' \leq r$ y $B_{r'}(x')$ corta a $B_r(x)$ entonces $B_{r'}(x') \subset B_{3r}(x)$.

La parte c) es consecuencia inmediata de b). ■

Teorema 9.40 *Si μ es una medida signada de Borel en \mathbb{R}^n y $t > 0$, entonces*

$$m(\{x \in \mathbb{R}^n \mid M\mu(x) > t\}) \leq \frac{3^n}{t} |\mu|(\mathbb{R}^n).$$

DEMOSTRACIÓN: Sea K un subconjunto compacto del abierto que aparece en el miembro izquierdo. Cada $x \in K$ es el centro de una bola abierta B tal que $|\mu|(B) > tm(B)$.

Estas bolas forman un cubrimiento de K , del cual podemos extraer un subcubrimiento finito al que a su vez podemos aplicar el teorema anterior, digamos $\{B_1, \dots, B_k\}$ de modo que

$$m(K) \leq 3^n \sum_{i=1}^k m(B_i) \leq \frac{3^n}{t} \sum_{i=1}^k |\mu|(B_i) \leq \frac{3^n}{t} |\mu|(\mathbb{R}^n).$$

La medida μ es regular porque lo son sus variaciones positiva y negativa, luego tomando el supremo sobre todos los compactos K obtenemos la desigualdad del enunciado. ■

Si $f \in L^1(\mathbb{R}^n)$ podemos aplicar el teorema anterior a la medida definida por

$$\mu(E) = \int_E |f| dm.$$

En este caso $M\mu$ es la función

$$Mf(x) = \sup_{0 < r < +\infty} \frac{1}{m(B_r(x))} \int_{B_r(x)} |f| dm$$

y la tesis del teorema es

$$m(\{x \in \mathbb{R}^n \mid Mf(x) > t\}) \leq \frac{3^n}{t} \|f\|_1. \quad (9.4)$$

La existencia de la derivada de una medida se deducirá de un teorema de existencia de puntos de Lebesgue, que definiremos a continuación:

Definición 9.41 Sea $f \in L^1(\mathbb{R}^n)$ (representaremos así al espacio $L^1(m)$, para la medida de Lebesgue en \mathbb{R}^n). Un *punto de Lebesgue* de f es un punto $x \in \mathbb{R}^n$ tal que

$$\lim_{r \rightarrow 0} \frac{1}{m(B_r(x))} \int_{B_r(x)} |f(y) - f(x)| dm(y) = 0.$$

Notemos que si x es un punto de Lebesgue entonces, dado que

$$\begin{aligned} \left| \frac{1}{m(B_r(x))} \int_{B_r(x)} f dm - f(x) \right| &= \left| \frac{1}{m(B_r(x))} \int_{B_r(x)} (f(y) - f(x)) dm(y) \right| \\ &\leq \frac{1}{m(B_r(x))} \int_{B_r(x)} |f(y) - f(x)| dm(y), \end{aligned}$$

se cumple

$$f(x) = \lim_{r \rightarrow 0} \frac{1}{m(B_r(x))} \int_{B_r(x)} f dm.$$

Es claro que si f es continua en x entonces x es un punto de Lebesgue para f , pero necesitamos la existencia de puntos de Lebesgue de funciones integrables cualesquiera:

Teorema 9.42 Si $f \in L^1(\mathbb{R}^n)$, entonces casi todo $x \in \mathbb{R}^n$ es un punto de Lebesgue de f .

DEMOSTRACIÓN: Sea

$$T_r(f)(x) = \frac{1}{m(B_r(x))} \int_{B_r(x)} |f - f(x)| dm$$

y sea

$$T(f)(x) = \overline{\lim}_{r \rightarrow 0} T_r(f)(x).$$

Tenemos que probar que $Tf = 0$ p.c.t.p. Fijemos un número real $y > 0$ y un número natural k . Por el teorema 9.24 existe una función $g \in C_c(\mathbb{R}^n)$ tal que $\|f - g\| - 1 < 1/k$. Sea $h = f - g$. La continuidad de g implica que $T(g) = 0$. Como

$$T_r(h)(x) \leq \frac{1}{m(B_r(x))} \int_{B_r(x)} |h| dm + |h(x)|,$$

tenemos que $T(h) \leq Mh + |h|$. Por otra parte, dado que $T_r(f) \leq T_r(g) + T_r(h)$, vemos que $T(f) \leq Mh + |h|$. Así pues,

$$\{x \in \mathbb{R}^n \mid T(f)(x) > 2y\} \subset \{x \in \mathbb{R}^n \mid M(h)(x) > y\} \cup \{x \in \mathbb{R}^n \mid |h|(x) > y\}$$

Llamemos $E(y, k)$ al miembro derecho de la inclusión anterior. Por 9.4 tenemos que la medida del primero de los conjuntos de la unión es menor o igual que $3^n/(yk)$. Respecto al segundo, llamémoslo A , observamos que

$$ym(A) \leq \int_A |h| dm \leq \int_{\mathbb{R}^n} |h| dm = \|h\|_1 < \frac{1}{k},$$

luego en total, $m(E(y, k)) \leq (3^n + 1)/(yk)$.

El conjunto $\{x \in \mathbb{R}^n \mid T(f)(x) > 2y\}$ es independiente de k y está contenido en la intersección de los conjuntos $E(y, k)$ para todo k , que es nula. Por la completitud de la medida de Lebesgue concluimos que es medible Lebesgue y tiene medida nula. Como esto vale para todo $y > 0$ concluimos que $Tf = 0$ p.c.t.p. ■

Con esto llegamos al teorema principal de esta sección:

Teorema 9.43 *Sea μ una medida signada de Borel en \mathbb{R}^n tal que $\mu \ll m$ y sea f la derivada de Radon-Nikodým de μ respecto a m . Entonces $d\mu/dm = f$ p.c.t.p. y para todo conjunto de Borel $E \subset \mathbb{R}^n$ se cumple*

$$\mu(E) = \int_E \frac{d\mu}{dm} dm.$$

DEMOSTRACIÓN: El teorema de Radon-Nikodým afirma que se verifica la igualdad del enunciado con f en lugar de $d\mu/dm$. Para cada punto de Lebesgue x de f se cumple

$$f(x) = \lim_{r \rightarrow 0} \frac{1}{m(B_r(x))} \int_{B_r(x)} f dm = \lim_{r \rightarrow 0} \frac{\mu(B_r(x))}{m(B_r(x))} = \frac{d\mu}{dm}(x).$$

■

9.6 El teorema de cambio de variable

En esta sección probaremos un teorema fundamental para el cálculo de integrales, junto con el teorema de Fubini. Se trata de la generalización a funciones de varias variables de la regla de integración por sustitución. El planteamiento es el siguiente: Supongamos que $g : U \rightarrow V$ es un difeomorfismo entre dos abiertos en \mathbb{R}^n , el problema es relacionar la integral de una función $f : V \rightarrow \mathbb{R}$ con la de $g \circ f$. Según el teorema 8.36, si g fuera una aplicación lineal de determinante Δ se cumpliría que $m(g(A)) = |\Delta| m(A)$, para todo conjunto medible $A \subset U$. En este caso no es difícil deducir que

$$\int_V f \, dm = |\Delta| \int_U (g \circ f) \, dm.$$

En el caso general, sabemos que en un entorno de cada punto x la aplicación g se confunde con su diferencial $dg(x)$, que es una aplicación lineal de determinante $\Delta_g(x) = \det Jg(x)$ (el *determinante jacobiano* de g en x). Esto se traduce en que si A es un conjunto medible contenido en un entorno de x suficientemente pequeño, entonces $m(g(A)) \approx |\Delta_g(x)| m(A)$. Esto es suficiente para llegar a un resultado análogo al caso lineal:

$$\int_V f \, dm = \int_U (g \circ f) |\Delta_g| \, dm.$$

Éste es el contenido del teorema de cambio de variable. La prueba detallada no es trivial en absoluto, sino que depende de una gran parte de los resultados que hemos visto hasta ahora. Observemos que $dg(x)$ es de hecho un isomorfismo, luego $\Delta_g(x) \neq 0$. Comencemos probando la relación entre la medida de un conjunto y la de su imagen para el caso de bolas abiertas:

Teorema 9.44 *Sea $g : U \rightarrow V$ un difeomorfismo entre dos abiertos de \mathbb{R}^n y $x \in U$. Entonces*

$$\lim_{r \rightarrow 0} \frac{m(g[B_r(x)])}{m(B_r(x))} = |\Delta_g(x)|.$$

DEMOSTRACIÓN: Puesto que la medida de Lebesgue es invariante por traslaciones, no perdemos generalidad si suponemos $x = 0$ y $g(0) = 0$. Sea $\phi = dg(0)$ y $h = g \circ \phi^{-1}$. Probaremos que

$$\lim_{r \rightarrow 0} \frac{m(h[B_r(0)])}{m(B_r(0))} = 1.$$

El teorema 8.36 nos da que $m(h[B_r(0)]) = |\Delta_g(0)|^{-1} m(g[B_r(x)])$, luego la igualdad anterior implica la que figura en el enunciado. La aplicación h cumple $h(0) = 0$ y además $dh(0)$ es la aplicación identidad. Por definición de diferenciabilidad esto significa que

$$\lim_{x \rightarrow 0} \frac{h(x) - x}{\|x\|} = 0.$$

Así, dado $0 < \epsilon < 1$ existe un $\delta > 0$ tal que

$$\text{si } \|x\| < \delta \text{ entonces } \|h(x) - x\| < \epsilon\|x\|. \quad (9.5)$$

Como h es un difeomorfismo en particular es una aplicación abierta, luego podemos tomar δ de modo que además $B_\delta(0)$ está contenida en la imagen de h . Veamos que si $0 < r < \delta$ entonces

$$B_{(1-\epsilon)r}(0) \subset h[B_r(0)] \subset B_{(1+\epsilon)r}(0). \quad (9.6)$$

En efecto, si $y \in B_{(1-\epsilon)r}(0) \subset B_\delta(0)$ entonces existe un $x \in U$ tal que $h(x) = y$. La relación 9.5 implica $\|y - x\| < \epsilon\|x\|$. Entonces

$$\|x\| \leq \|x - y\| + \|y\| < \epsilon\|x\| + (1 - \epsilon)r,$$

luego $(1 - \epsilon)\|x\| < (1 - \epsilon)r$ y $\|x\| < r$. Así $y \in h[B_r(0)]$. Por otra parte, si $y \in h[B_r(0)]$, es decir, si $y = h(x)$ con $\|x\| < r$, la relación 9.5 implica que $\|y - x\| < \epsilon r$, luego $\|y\| < (1 + \epsilon)r$, lo que nos da la otra inclusión.

Tomando medidas en 9.6 resulta

$$(1 - \epsilon)^n \leq \frac{m(h[B_r(0)])}{m(B_r(0))} \leq (1 + \epsilon)^n, \quad \text{para todo } r < \delta,$$

y la conclusión es clara. ■

En las condiciones del teorema anterior la aplicación g biyecta claramente los conjuntos de Borel de U con los de V , por lo que podemos definir una medida de Borel en U mediante $\mu(A) = m(g[A])$. Si se trata de una medida finita el teorema afirma que

$$\frac{d\mu}{dm}(x) = \Delta_g(x).$$

En realidad no hay ningún problema en definir la derivada de una medida positiva, aunque no sea finita, pues se trata de un concepto local, y la igualdad anterior es cierta en cualquier caso.

Todo conjunto medible Lebesgue se puede expresar como unión de un conjunto de Borel y un conjunto nulo. El teorema 8.38 prueba que si $g : U \rightarrow V$ es un difeomorfismo entre dos abiertos de \mathbb{R}^n y $E \subset U$ es medible Lebesgue, entonces $g[E]$ es medible Lebesgue.

Veamos finalmente el teorema principal:

Teorema 9.45 (Teorema de cambio de variable) *Sea $g : U \rightarrow V$ un difeomorfismo entre dos abiertos de \mathbb{R}^n y sea $f : V \rightarrow \mathbb{R}$ una aplicación integrable Lebesgue. Entonces*

$$\int_V f \, dm = \int_U (g \circ f) |\Delta_g| \, dm.$$

DEMOSTRACIÓN: Para cada natural k , sea $U_k = \{x \in U \mid \|g(x)\| < k\}$. Claramente U_k es abierto. Para cada conjunto medible E definimos $\mu_k(E) = m(g[E \cap U_k])$. Claramente μ_k es una medida finita sobre la σ -álgebra de los conjuntos medibles Lebesgue en \mathbb{R}^n . El teorema 8.38 prueba que $\mu_k \ll m$.

Ahora podemos aplicar el teorema 9.43, según el cual existe $d\mu_k/dm$ p.c.t.p., es integrable Lebesgue y para todo conjunto medible E se cumple

$$\mu_k(E) = \int_E \frac{d\mu_k}{dm} dm.$$

En principio 9.43 prueba esto para conjuntos de Borel, pero la igualdad se extiende obviamente a conjuntos medibles arbitrarios. Del teorema 9.44 se sigue fácilmente que si $x \in U_k$ entonces

$$\frac{d\mu_k}{dm} = |\Delta_g(x)|.$$

En total hemos probado que si E es medible entonces

$$m(g[E \cap U_k]) = \int_{U_k} \chi_E |\Delta_g| dm.$$

Por el teorema de la convergencia monótona concluimos que

$$m(g[E \cap U]) = \int_U \chi_E |\Delta_g| dm. \quad (9.7)$$

Vamos a deducir de aquí que si A es un conjunto medible, entonces

$$\int_V \chi_A dm = \int_U (g \circ \chi_A) |\Delta_g| dm.$$

Basta tomar $E = g^{-1}(A) \subset U$. El comentario previo al teorema prueba que E es medible y claramente $\chi_E = g \circ \chi_A$. Además $g[E \cap U] = g[E] = A \cap V$, luego (9.7) se convierte en la igualdad buscada.

De aquí se sigue la fórmula del enunciado para el caso en que f es una función simple no negativa. Por el teorema de la convergencia monótona llegamos al mismo resultado para funciones medibles no negativas y a su vez se extiende a toda función integrable aplicándolo a f^+ y f^- . ■

Ejemplo Vamos a probar que

$$\int_{-\infty}^{+\infty} e^{-x^2} dx = \sqrt{\pi}.$$

Aparentemente se trata de un problema de análisis de una variable real, pero el cálculo es mucho más simple si nos apoyamos en una función de dos variables. Concretamente consideramos $f(x, y) = e^{-x^2 - y^2}$. Calculamos la integral de esta función en la bola de centro 0 y radio r mediante el cambio a coordenadas polares:

$$\int_{B_r(0)} e^{-x^2-y^2} dx dy = \int_0^{2\pi} \int_0^\infty e^{-\rho^2} \rho d\rho d\theta = 2\pi \left[-\frac{e^{-\rho^2}}{2} \right]_0^r = \pi(1 - e^{-r^2}).$$

El teorema de la convergencia monótona implica que f es integrable en \mathbb{R}^2 y además

$$\int_{\mathbb{R}^2} e^{-x^2-y^2} dx dy = \pi.$$

Por otro lado podemos aplicar el teorema de Fubini, que nos da

$$\left(\int_{-\infty}^{\infty} e^{-x^2} dx \right)^2 = \pi,$$

luego tenemos la igualdad que buscábamos. De aquí se deducen varias integrales de interés. En primer lugar

$$\int_0^{+\infty} e^{-x^2} dx = \frac{\sqrt{\pi}}{2},$$

y haciendo el cambio $x = \sqrt{t}$ resulta

$$\Gamma(-1/2) = \int_0^{+\infty} t^{-1/2} e^{-t} dt = \sqrt{\pi}.$$

Por la ecuación funcional de la función factorial concluimos que

$$\Gamma(1/2) = \frac{1}{2} \Gamma(-1/2) = \frac{\sqrt{\pi}}{2}.$$

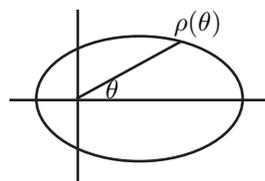
■

Ejemplo Con el cálculo que acabamos de hacer podemos dar una expresión explícita para las constantes v_n que calculamos en el ejemplo de la página 357. En efecto, ahora es inmediato que

$$v_n = \frac{\pi^{n/2}}{\Gamma(n/2)}, \quad (9.8)$$

pues esta función coincide con v_n para $n = 1$ y $n = 2$ y satisface la misma relación recurrente que v_n . En el capítulo siguiente daremos una prueba más elegante de esta fórmula. ■

Ejemplo Consideremos una curva cerrada en \mathbb{R}^2 que rodee a $(0, 0)$ y admita una expresión en coordenadas polares $\rho = \rho(\theta)$ (es decir, que corte a cada semirrecta de origen $(0, 0)$ en un único punto). El recinto S limitado por la curva estará formado por los puntos (ρ_0, θ_0) tales que $0 \leq \rho_0 \leq \rho(\theta_0)$. Para calcular el área de S



efectuamos el cambio de variables $x = \rho \cos \theta$, $y = \rho \operatorname{sen} \theta$, cuyo jacobiano es ρ . Así, el área se puede calcular como

$$\int_S dx dy = \int_0^{2\pi} \int_0^{\rho(\theta)} \rho d\rho d\theta = \int_0^{2\pi} \frac{\rho^2}{2} d\theta.$$

(Hemos aplicado el teorema anterior en el abierto

$$\{(\rho, \theta) \mid 0 < \rho < \rho(\theta), 0 < \theta < 2\pi\}.$$

El cambio de coordenadas lo transforma biyectivamente en S menos el radio $\theta = 0$, que tiene área nula, por lo que no importa despreciarlo.)

Por ejemplo, el área de la cardioide $\rho = (a/2)(1 + \cos \theta)$ viene dada por

$$\begin{aligned} \frac{a^2}{8} \int_0^{2\pi} (1 + \cos \theta)^2 d\theta &= \frac{a^2}{8} \int_0^{2\pi} \left(1 + 2 \cos \theta + \frac{1 + \cos 2\theta}{2} \right) d\theta \\ &= \frac{a^2}{8} \left[\theta + 2 \operatorname{sen} \theta + \frac{\theta}{2} + \frac{\operatorname{sen} 2\theta}{4} \right]_0^{2\pi} = \frac{3\pi}{8} a^2. \end{aligned}$$

■

Ejemplo En el capítulo VII demostramos que los cuerpos sometidos a la acción gravitatoria de una estrella o planeta describen trayectorias rectas o cónicas, pero no calculamos la posición del cuerpo en función del tiempo. Ahora probaremos la segunda ley de Kepler, que aporta información a este respecto. Se refiere a un cuerpo (un planeta, un cometa) que describe una trayectoria cónica alrededor (digamos) del Sol:

El radio que une el móvil con el Sol barre áreas iguales en tiempos iguales.

Tomemos como origen la posición del Sol y sea $\rho(\theta)$ la trayectoria del móvil. Sea A el sector de cónica que barre el radio que une al móvil con el Sol entre un ángulo θ_0 y un ángulo θ_1 . El área de A es

$$\int_A dx dy = \frac{1}{2} \int_{\theta_0}^{\theta_1} \rho^2 d\theta.$$

Hacemos el cambio $\theta = \theta(t)$, donde t es el tiempo. El resultado es

$$\frac{1}{2} \int_{t_0}^{t_1} \rho^2(\theta(t)) \theta'(t) dt = \frac{1}{2} \int_{t_0}^{t_1} \rho^2(t) \omega(t) dt = \frac{1}{2m} \int_{t_0}^{t_1} L dt = \frac{L}{2m} (t_1 - t_0).$$

Así pues, el área barrida es proporcional al tiempo recorrido.

A su vez de aquí se deduce la tercera ley de Kepler, válida para móviles que describen órbitas elípticas alrededor de un mismo cuerpo. El período de revolución de tal cuerpo es el tiempo que tarda en recorrer una órbita completa:

Los cuadrados de los períodos de revolución son proporcionales a los cubos de los semiejes de las órbitas.

Según vimos en el capítulo VII, la ecuación de la órbita es

$$\rho = \frac{L^2}{GMm^2} \frac{1}{1 + \epsilon \cos \theta}.$$

Los vértices mayores (los valores máximo y mínimo de ρ) se corresponden con los ángulos $\theta = 0, \pi$. Su semisuma es el semieje mayor:

$$a = \frac{L^2}{GMm^2} \frac{1}{2} \left(\frac{1}{1 + \epsilon} - \frac{1}{1 - \epsilon} \right) = \frac{L^2}{GMm^2} \frac{1}{1 - \epsilon^2}.$$

Puesto que ϵ es la excentricidad, el semieje menor es $b = a\sqrt{1 - \epsilon^2}$. El área de la elipse es

$$A = \pi ab = \pi a^2 \sqrt{1 - \epsilon^2} = \frac{\pi}{G^2 M^2} \frac{L^4}{m^4 (1 - \epsilon^2)^2} \sqrt{1 - \epsilon^2}.$$

Según hemos calculado, el período T cumple $A = LT/(2m)$, luego

$$T = \frac{2\pi}{G^2 M^2} \frac{L^3}{m^3} \frac{\sqrt{1 - \epsilon^2}}{(1 - \epsilon^2)^2}.$$

Reuniendo todo esto vemos que

$$\frac{T^2}{a^3} = \frac{4\pi^2}{GM},$$

luego tenemos la proporción buscada. ■

Centro de masas y tensor de inercia Supongamos que $V \subset \mathbb{R}^3$ es un sólido cuya masa está determinada por la función de densidad ρ (es decir, que la masa contenida en una región $U \subset V$ se obtiene integrando sobre U el diferencial de masa $dm = \rho dx dy dz$ (de modo que en este ejemplo dm no representará la medida de Lebesgue). En particular, la masa de V es

$$M = \int_V dm = \int_V \rho dx dy dz.$$

En Física se define el *centro de masas* del sólido V como

$$\text{CM} = \frac{1}{M} \int_V (x, y, z) dm,$$

donde la integral de una función vectorial se define como el vector formado por las integrales de las componentes. Similarmente, el *tensor de inercia* de V respecto de un punto O (que por simplicidad tomamos como el origen de

coordenadas) se define como la forma bilineal $J_O : \mathbb{R}^3 \rightarrow R$ determinada por la matriz simétrica I_O cuyas componentes son¹

$$I_{kl} = \int_V ((x_1^2 + x_2^2 + x_3^2)\delta_{kl} - x_k x_l) dm,$$

donde $\delta_{kl} = \begin{cases} 1 & \text{si } k = l, \\ 0 & \text{si } k \neq l. \end{cases}$

Vamos a estudiar estos conceptos para el caso de un sólido con simetría axial, es decir, parametrizable en la forma

$$V = \{(\lambda r(z) \cos \theta, \lambda r(z) \sin \theta, z) \mid 0 \leq \lambda \leq 1, 0 \leq \theta \leq 2\pi, a \leq z \leq b\},$$

cuya función de densidad tenga también simetría axial, con lo que será de la forma $\rho = \rho(\lambda, z)$. Vamos a probar que en este caso el centro de masas está en el eje de simetría (el eje Z) y que la matriz del tensor de inercia es diagonal con $I_{11} = I_{22}$.

En efecto, es fácil ver que el determinante jacobiano del cambio de coordenadas

$$(x, y, z) = (\lambda r(z) \cos \theta, \lambda r(z) \sin \theta, z)$$

es $\lambda r^2(z)$. Por lo tanto, la primera coordenada del centro de masas es

$$\begin{aligned} \frac{1}{M} \int_V x \rho dx dy dz &= \frac{1}{M} \int_0^1 \int_0^{2\pi} \int_a^b \lambda r(z) \cos \theta \rho(\lambda, z) \lambda r^2(z) dz d\theta d\lambda \\ &= \frac{1}{M} \int_0^1 \int_a^b \lambda^2 r^3(z) \rho(\lambda, z) dz d\lambda \int_0^{2\pi} \cos \theta d\theta = 0, \end{aligned}$$

pues la última integral vale 0. Igualmente se prueba que la segunda coordenada del centro de masas es nula, luego éste se encuentra sobre el eje Z . Para el tensor de inercia tenemos que

$$\begin{aligned} I_{12} &= - \int_V xy \rho dx dy dz = \\ &= - \int_0^1 \int_0^{2\pi} \int_a^b \lambda^2 r^2(z) \cos \theta \sin \theta \rho(\lambda, z) \lambda r^2(z) dz d\theta d\lambda = 0, \end{aligned}$$

porque al igual que antes podemos separar la parte en θ y su integral es 0.

El mismo razonamiento vale para I_{13} e I_{23} . Por último:

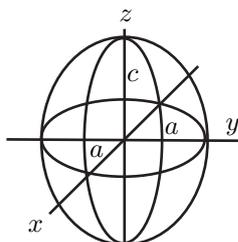
$$\begin{aligned} I_{11} &= \int_V (y^2 + z^2) \rho dx dy dz = \\ &= \int_0^1 \int_0^{2\pi} \int_a^b (\lambda^2 r^2(z) \sin^2 \theta + z^2) \rho(\lambda, z) \lambda r^2(z) dz d\theta d\lambda = \\ &= \int_0^1 \int_0^{2\pi} \int_a^b z^2 \rho(\lambda, z) \lambda r^2(z) dz d\theta d\lambda, \end{aligned}$$

porque la integral del primer sumando se anula por la misma razón de siempre,

¹Véase el apéndice C para una motivación de estas definiciones.

porque se puede separar la parte en θ y tiene integral nula. La expresión para I_{22} es igual salvo que tiene $\cos^2 \theta$ en lugar de $\sin^2 \theta$, por lo que se simplifica hasta la misma integral. Por lo tanto, $I_1 = I_2$. ■

Ejemplo Vamos a calcular el tensor de inercia de un elipsoide de revolución de densidad constante ρ respecto de su centro geométrico (que por simetría ha de ser su centro de masas). Un elipsoide es el volumen encerrado por la superficie que resulta de girar una elipse sobre sus ejes.



Si los ejes de la elipse miden $2a$ y $2c$, la ecuación (de la superficie) del elipsoide es

$$\frac{x^2}{a^2} + \frac{y^2}{a^2} + \frac{z^2}{c^2} = 1,$$

de modo que el eje de rotación es el eje Z y el origen de coordenadas es el centro del elipsoide. Se trata de un caso particular del elipsoide cuyo volumen hemos calculado en el ejemplo de la página 356, que resulta ser $V = \frac{4}{3}\pi a^2 c$. Por consiguiente, $M = \frac{4}{3}\pi \rho a^2 c$. Ahora calculamos

$$\begin{aligned} I_{11} &= \int_V (y^2 + z^2) dm = \rho \int_{-a}^a \int_{-\sqrt{a^2-x^2}}^{\sqrt{a^2-x^2}} \int_{-c\sqrt{1-x^2/a^2-y^2/a^2}}^{c\sqrt{1-x^2/a^2-y^2/a^2}} (y^2 + z^2) dz dy dx \\ &= \rho \int_{-a}^a \int_{-\sqrt{a^2-x^2}}^{\sqrt{a^2-x^2}} \left[y^2 z + \frac{z^3}{3} \right]_{-c\sqrt{1-x^2/a^2-y^2/a^2}}^{c\sqrt{1-x^2/a^2-y^2/a^2}} dy dx \\ &= 2\rho \int_{-a}^a \int_{-\sqrt{a^2-x^2}}^{\sqrt{a^2-x^2}} \left(cy^2 \sqrt{1 - \frac{x^2}{a^2} - \frac{y^2}{a^2}} + \frac{c^3}{3} \left(1 - \frac{x^2}{a^2} - \frac{y^2}{a^2} \right)^{3/2} \right) dy dx. \end{aligned}$$

Separamos las dos integrales. La primera, tras hacer un cambio a polares $x = r \cos \theta$, $y = r \sin \theta$ (con jacobiano r) queda

$$\frac{2c\rho}{a} \int_0^{2\pi} \int_0^a r^3 \sqrt{a^2 - r^2} \sin^2 \theta dr d\theta = \frac{2\pi c\rho}{a} \int_0^a r^3 \sqrt{a^2 - r^2} dr = \frac{4\pi a^4 c\rho}{15}.$$

(La última integral se hace con el cambio $r = a\sqrt{1-t^2}$.)

Falta calcular

$$\frac{2c^3\rho}{3a^3} \int_0^{2\pi} \int_0^a (a^2 - r^2)^{3/2} r dr d\theta = \frac{4\pi c^3\rho}{3a^3} \int_0^a (a^2 - r^2)^{3/2} r dr = \frac{4\pi a^2 c^3\rho}{15}.$$

Así pues, en total queda:

$$I_{11} = \frac{4\pi a^2 c \rho}{15} (a^2 + c^2) = \frac{M}{5} (a^2 + c^2),$$

donde M es la masa del sólido. El ejemplo anterior nos da que $I_{22} = I_{11}$. En cuanto a I_{33} , viene dado por

$$\begin{aligned} I_{33} &= \int_V (x^2 + y^2) dm = \rho \int_{-c}^c \int_0^{2\pi} \int_0^{a\sqrt{1-z^2/c^2}} r^3 dr d\theta dz \\ &= \frac{\pi a^4 \rho}{2} \int_{-c}^c \left(1 - \frac{z^2}{c^2}\right)^2 dz = \frac{8\pi a^4 c \rho}{15} = \frac{2Ma^2}{5}. \end{aligned}$$

En resumen: las componentes del tensor de inercia del elipsoide son

$$I_{11} = I_{22} = \frac{M}{5} (a^2 + c^2), \quad I_{33} = \frac{2Ma^2}{5}, \quad (9.9)$$

donde a y c son los semiejes y M es la masa. En particular, el tensor de inercia de una esfera homogénea de radio r y masa M es diagonal en cualquier sistema de referencia con origen en el centro de la esfera, y sus componentes diagonales son

$$I_{11} = I_{22} = I_{33} = \frac{2Mr^2}{5}.$$

■

Ejemplo: El tensor de inercia de la Tierra Vamos a aplicar los cálculos del ejemplo precedente al caso de la Tierra. Su masa es $M = 5.9736 \times 10^{24}$ kg y su radio $r = 6.371 \times 10^6$ m. Por lo tanto, las componentes diagonales de su tensor de inercia valen

$$I_{11} = I_{22} = I_{33} = 9.69865 \times 10^{37} \text{ kg}\cdot\text{m}^2.$$

Ahora bien, aquí hemos supuesto que la Tierra es esférica. En realidad está achatada por los polos, y sus radios ecuatorial y polar son, respectivamente

$$r_e = 6.3781 \times 10^6 \text{ m}, \quad r_p = 6.35675 \times 10^6 \text{ m}.$$

Si la consideramos un elipsoide de revolución, las componentes del tensor de inercia serán

$$I_{11} = I_{22} = 9.6878 \times 10^{37} \text{ kg}\cdot\text{m}^2, \quad I_{33} = 9.7203 \times 10^{37} \text{ kg}\cdot\text{m}^2.$$

Para estos cálculos hemos supuesto que la densidad de la Tierra es constante. Vamos a suponer ahora un núcleo esférico de radio $r_n = 3.5 \times 10^6$ m y densidad $d_n = 10\,000 \text{ kg/m}^3$. Esto deja para el *manto* (la capa que rodea al núcleo)² una

²En realidad la Tierra tiene una corteza externa menos densa que el manto, pero es tan fina que a estos efectos podemos despreciarla.

densidad $d_m = 4622.83 \text{ kg/m}^3$. Podemos calcular I_{11} (e igualmente I_{33}) como el valor correspondiente a un elipsoide del tamaño de la Tierra y densidad d_m , menos el valor correspondiente a una esfera del tamaño del núcleo y densidad d_m más el valor correspondiente al núcleo. Así se obtiene:

$$I_{11} = 8.59494 \times 10^{37} \text{ kg}\cdot\text{m}^2, \quad I_{33} = 8.62217 \times 10^{37} \text{ kg}\cdot\text{m}^2.$$

Por otra parte, observemos que de (9.9) se sigue también que

$$\frac{I_{33} - I_{11}}{I_{33}} = \frac{(a+c)/2}{a} \frac{a-c}{a} = \frac{r_m}{r_e} \epsilon, \quad (9.10)$$

donde r_m es la media de los radios y ϵ es el achatamiento de la Tierra. La ventaja de esta fórmula es que depende de magnitudes relativamente fáciles de medir con precisión. Para la Tierra (considerada como elipsoide homogéneo), $\epsilon = 0.00335$ y el cociente de momentos de inercia es casi el mismo. ■

9.7 Integración en variedades

Si queremos calcular el área de una esfera, lo primero que debemos preguntarnos es ¿qué es el área de una esfera? La pregunta es menos trivial de lo que parece, pues podemos calcular el área de cualquier figura plana razonable (es decir, medible Lebesgue), pero cualquier fragmento de esfera, por pequeño que sea, no es plano, y no es evidente cómo puede compararse su área con la de ninguna figura plana. Obviamente podemos definir aplicaciones de conjuntos planos en la esfera (difeomorfismos incluso, es decir, cartas de la esfera), pero nada nos garantiza que conserven el área, y difícilmente podríamos probar que así es sin tener de hecho una definición de área.

Para hacernos una idea de lo que vamos a hacer pensemos en una esfera del tamaño de la Tierra. Si situamos sobre ella una baldosa plana de un metro cuadrado, no descansaría perfectamente sobre el suelo esférico, pero, suponiendo que se apoyara en su centro, cada esquina se levantaría tan sólo 0.04 micras del suelo. Resulta, pues, razonable considerar que el fragmento de la Tierra cubierto por la baldosa (aunque sea imperfectamente en teoría) tiene una superficie de 1m^2 , salvo un error muy pequeño. Si cubrimos toda la Tierra con baldosas de 1m^2 , aunque sin duda no encajarán a la perfección, el número de baldosas empleadas será una buena aproximación de la superficie de la Tierra, y el mínimo error cometido se podrá reducir arbitrariamente a base de considerar baldosas más y más pequeñas.

Para formalizar (y generalizar) esta idea empezamos observando que si E^n es un espacio vectorial euclídeo de dimensión n (por ejemplo un subespacio de \mathbb{R}^m de dimensión n con el producto escalar inducido desde \mathbb{R}^m), entonces existe una isometría $\phi : E^n \rightarrow \mathbb{R}^n$ que es, de hecho, un homeomorfismo entre las topologías euclídeas. Definimos los conjuntos medibles de E^n como las antiimágenes por ϕ de los conjuntos medibles Lebesgue de \mathbb{R}^n y la *medida de Lebesgue* en E^n como la dada por $m(G) = m(\phi(G))$.

Teniendo en cuenta que las isometrías en \mathbb{R}^n conservan la medida de Lebesgue, es inmediato comprobar que m así definida es una medida en E^n que no depende de la elección de ϕ . Claramente se corresponde con la noción de longitud, área, volumen, etc. en E^n . El teorema 8.36 puede enunciarse ahora en este contexto general:

Teorema 9.46 *Sea $\phi : E^n \rightarrow F^n$ una aplicación lineal entre dos espacios vectoriales euclídeos de dimensión n y sea Δ_ϕ el determinante de la matriz de ϕ respecto a dos bases ortonormales cualesquiera. Entonces, para todo conjunto medible $A \subset E^n$ se cumple $m(\phi(A)) = |\Delta_\phi| m(A)$.*

DEMOSTRACIÓN: Sean $f : \mathbb{R}^n \rightarrow E^n$ y $g : \mathbb{R}^n \rightarrow F^n$ dos isometrías. Sea $h = f \circ \phi \circ g^{-1}$. Claramente el determinante de h es igual a Δ_ϕ y $\phi = f^{-1} \circ h \circ g$. Basta aplicar el teorema 8.36 y el hecho de que f y g conservan la medida. ■

Así pues, si $S \subset \mathbb{R}^m$ es una variedad diferenciable de dimensión n , para cada punto $p \in S$ tenemos definida la medida de Lebesgue en el espacio tangente $T_p S$ que es, según sabemos, una formalización adecuada del concepto de longitud, área, volumen, o una generalización natural de éste, según la dimensión n .

Consideremos una carta $X : U \rightarrow \mathbb{R}^m$, llamemos $V = X[U]$, fijemos un punto $p = X(x)$ y consideremos la diferencial $dX(x) : \mathbb{R}^n \rightarrow T_p S$. Para cada punto $t \in U$ tenemos que

$$X(t) \approx p + dX(x)(t - p),$$

de modo que, si t está suficientemente próximo a p , el punto $X(t) \in V$ “se confunde” con $p + dX(x)(t - p) \in p + T_p S$. Por lo tanto, si $B \subset V$ es un conjunto de Borel en un entorno suficientemente pequeño de p y $B_X = X^{-1}[B]$, el conjunto de Borel $B^t = dX(x)[B_X - x] \subset T_p S$ cumple que $p + B^t$ es prácticamente indistinguible de B .

Particularizando esto al ejemplo de la Tierra que hemos considerado antes, si B es la superficie cubierta por una baldosa que se apoya en el punto p , entonces $p + B^t$ es la baldosa tangente a la superficie terrestre.³ Por ello, lo que vamos a hacer es probar que existe una medida de Borel regular μ en V con la propiedad de que si $K \subset V$ es un subconjunto compacto y lo cubrimos por un conjunto finito de conjuntos de Borel B_i disjuntos dos a dos, entonces la suma de las medidas $m(B_i^t)$ se aproxima a $\mu(K)$, y la aproximación es mejor cuanto menores son los diámetros de los conjuntos B_i . Esta propiedad justifica que consideremos a $\mu(K)$ como el volumen de K .

Empezamos calculando las medidas $m(B^t)$:

³Al menos si elegimos la carta X de modo que la aplicación $t \mapsto p + dX(x)(t - p)$ sea la proyección ortogonal en el plano tangente. Puede probarse que podemos elegir así la carta, pero no será necesario, porque la medida que vamos a construir sobre S no dependerá de la elección de ninguna carta de S .

Teorema 9.47 Sea $X : U \rightarrow \mathbb{R}^m$ una carta de una variedad S . Sea $X(x) = p$. Para cada conjunto de Borel $B \subset U$ se cumple $m(dX(x)[B]) = \Delta_X(x) m(B)$, donde $\Delta_X = \sqrt{\det(g_{ij})}$ y las funciones g_{ij} son los coeficientes del tensor métrico de S .

DEMOSTRACIÓN: Sea $\phi : \mathbb{R}^m \rightarrow \mathbb{R}^m$ una isometría que transforme $T_p(S)$ en $\mathbb{R}^n \times \{0\}$ y sea $p_n : \mathbb{R}^m \rightarrow \mathbb{R}^n$ la proyección en las primeras componentes. De este modo la restricción a $T_p(S)$ de $\phi \circ p_n$ es una isometría del plano tangente en \mathbb{R}^n , luego $m(dX(x)[B]) = m((dX(x) \circ \phi \circ p_n)[B])$.

Sean A y P_n las matrices de ϕ y p_n en las bases canónicas. Entonces la matriz de $(dX(x) \circ \phi \circ p_n)$ es $J_X(x)AP_n$ y según el teorema 9.46 tenemos que $m(dX(x)[B_r]) = \Delta_X(x) m(B_r)$, donde $\Delta_X(x) = |\det(J_X(x)AP_n)|$. Ahora usamos que para toda matriz cuadrada M se cumple $|\det M| = \sqrt{\det(MM^t)}$, con lo que

$$\Delta_X(x) = \sqrt{\det(J_X(x)AP_nP_n^tA^tJ_X(x)^t)}.$$

Pero $J_X(x)AP_nP_n^tA^tJ_X(x)^t = J_X(x)AA^tJ_X(x)^t$. En efecto, el elemento (i, j) de esta matriz es el producto de $e_iJ_X(x)AP_n$ por $(e_jJ_X(x)AP_n)^t$, donde e_i y e_j son los vectores de la base canónica de \mathbb{R}^n , pero $e_iJ_X(x) \in T_p(S)$, luego $e_iJ_X(x)A \in \mathbb{R}^n \times \{0\}$, e igualmente $e_jJ_X(x)A \in \mathbb{R}^n \times \{0\}$, luego el producto es el mismo aunque suprimamos P_n . Como A es la matriz de una isometría, se cumple $AA^t = I$, luego concluimos que

$$\Delta_X(x) = \sqrt{\det(J_X(x)J_X(x)^t)} = \sqrt{\det(g_{ij}(x))}. \quad \blacksquare$$

Así, con la notación previa al teorema, tenemos que $m(B^t) = \Delta_X(x) m(B_X)$. Ahora vamos a probar que la medida en V dada por

$$\mu(B) = \int_{X^{-1}[B]} \Delta_X dm$$

cumple la propiedad requerida.

Claramente μ es una medida de Borel regular en V . Es fácil comprobar que si f es una función integrable en V , entonces

$$\int_V f d\mu = \int_U (X \circ f) \Delta_X dm$$

(se prueba primero para funciones simples y luego para no negativas).

Sea ahora $K \subset V$ un conjunto compacto, de modo que $K_X \subset U$ es también compacto. Dado $\epsilon > 0$, como $X^{-1} \circ \Delta_X$ es uniformemente continua sobre K , existe un $\delta > 0$ tal que si $p = X(x)$, $q = X(y) \in K$ cumplen $\|p - q\| < \delta$ entonces

$$|\Delta_X(x) - \Delta_X(y)| < \epsilon_0 = \frac{\epsilon}{1 + m(K_X)}.$$

A continuación consideramos cualquier cubrimiento finito de K por conjuntos de Borel B_i disjuntos dos a dos de diámetro menor que δ (es fácil probar que

existen tales cubrimientos) y tomemos puntos $x_i \in X^{-1}[B_i]$. Así, si $x \in X^{-1}[B_i]$ tenemos que $|\Delta_X(x) - \Delta_X(x_i)| < \epsilon_0$. Por consiguiente:

$$\begin{aligned} |\sum_i m(B_i^t) - \mu(K)| &= |\sum_i (m(B_i^t) - \mu(B_i))| \\ &= \left| \sum_i (\Delta_X(x_i) m(X^{-1}[B_i]) - \int_{X^{-1}[B_i]} \Delta_X dm) \right| \\ &= \left| \sum_i \int_{X^{-1}[B_i]} (\Delta_X(x_i) - \Delta_X) dm \right| \leq \sum_i \int_{X^{-1}[B_i]} |\Delta_X(x_i) - \Delta_X| dm \\ &\leq \epsilon_0 \sum_i m(X^{-1}[B_i]) = \epsilon_0 m(K_X) < \epsilon, \end{aligned}$$

como había que probar.

En principio tenemos definida una medida sobre la imagen de cada carta de la variedad S , pero a continuación probamos que todas estas medidas se extienden a una única medida en S que, en particular, no depende de ninguna carta en concreto:

Teorema 9.48 *Sea $S \subset \mathbb{R}^m$ una variedad diferenciable de dimensión n . Entonces existe una única medida de Borel regular m en S tal que si $X : U \rightarrow V$ es una carta y f es una función integrable en V , se tiene*

$$\int_V f dm = \int_U (X \circ f) \Delta_X dm, \quad \text{donde } \Delta_X = \sqrt{\det(g_{ij})}.$$

DEMOSTRACIÓN: Consideremos dos cartas $X : U_1 \rightarrow V_1$ e $Y : U_2 \rightarrow V_2$ y sea $B \subset V_1 \cap V_2$ un conjunto de Borel. Vamos a probar que

$$\int_{X^{-1}[B]} \Delta_X dm = \int_{Y^{-1}[B]} \Delta_Y dm.$$

Restringiendo las cartas podemos suponer que $V_1 = V_2$, y entonces basta aplicar el teorema de cambio de variable a la función $f = X \circ Y^{-1}$. En efecto, se cumple que $f : U_1 \rightarrow U_2$ es un difeomorfismo y, como $f \circ Y = X$, tenemos $J_f(x)J_Y(f(x)) = J_X(x)$, luego

$$\begin{aligned} \Delta_X^2(x) &= \det(J_X(x)J_X(x)^t) = \det(J_f(x)J_Y(f(x))J_Y(f(x))^t J_f(x)^t) \\ &= \det(J_f(x))^2 \Delta_Y^2(f(x)) \end{aligned}$$

y por consiguiente $\Delta_X(x) = \Delta_Y(f(x)) |\det J_f(x)| = \Delta_Y(f(x)) |\Delta_f(x)|$. Así:

$$\begin{aligned} \int_{X^{-1}[B]} \Delta_X dm &= \int_{U_1} \chi_{X^{-1}[B]} \Delta_X dm = \int_{U_1} (f \circ \chi_{Y^{-1}})(f \circ \Delta_Y) |\Delta_f| dm \\ &= \int_{U_2} \chi_{Y^{-1}[B]} \Delta_Y dm = \int_{Y^{-1}[B]} \Delta_Y dm. \end{aligned}$$

Con esto hemos probado que las distintas medidas que tenemos definidas sobre las imágenes de las cartas de S coinciden en su dominio común. El paso siguiente es “pegar” las medidas correspondientes a un número finito de cartas. Aunque intuitivamente es obvio que esto puede hacerse, formalmente conviene simplificar las comprobaciones usando el teorema de Riesz.

Supongamos que $X_i : U_i \rightarrow V_i$, para $i = 1, \dots, k$ son cartas de S con imágenes acotadas. Por el teorema 3.17 existe una partición de la unidad subordinada a los abiertos V_i , es decir, una familia de funciones $h_i \prec V_i$ tales que $h_1 + \dots + h_k = 1$. Sea $V = V_1 \cup \dots \cup V_k$. Para cada $f \in C_c(V)$ definimos

$$T(f) = \int_{V_1} h_1 f d\mu_1 + \dots + \int_{V_k} h_k f d\mu_k,$$

donde μ_i es la medida asociada a la carta X_i . Claramente T es un operador lineal y positivo, luego existe una medida μ en V tal que para toda $f \in C_c(V)$ se cumple

$$\int_V f d\mu = \int_{V_1} h_1 f d\mu_1 + \dots + \int_{V_k} h_k f d\mu_k.$$

Si en particular tomamos $f \in C_c(V_i)$ entonces $h_j f \in C_c(V_i \cap V_j)$, luego

$$\int_{V_j} h_j f d\mu_j = \int_{V_i \cap V_j} h_j f d\mu_j = \int_{V_i \cap V_j} h_j f d\mu_i = \int_{V_i} h_j f d\mu_i,$$

luego

$$\int_V f d\mu = \int_{V_i} (h_1 + \dots + h_k) f d\mu_i = \int_{V_i} f d\mu_i.$$

Por la unicidad del teorema de Riesz esto prueba que la restricción de μ a V_i es precisamente μ_i , y es claro que esta propiedad determina a μ . En particular la construcción de μ no depende de la partición de la unidad escogida.

Finalmente “pegamos” todas las medidas asociadas a todas las cartas en una única medida en S . Para ello definimos un operador $T : C_c(S) \rightarrow \mathbb{R}$. Para cada $f \in C_c(S)$ tomamos un número finito de cartas con imagen acotada cuya unión cubra el soporte de f . Sea V la unión de las imágenes y μ_V la medida sobre V que acabamos de construir. Definimos

$$T(f) = \int_V f d\mu_V.$$

Es claro que $T(f)$ no depende de las cartas con que cubrimos el soporte, pues si realizamos dos cubrimientos distintos $V = V_1 \cup \dots \cup V_k$ y $V' = V'_1 \cup \dots \cup V'_k$, entonces cada abierto $V_i \cap V'_j$ es la imagen de dos cartas que inducen la misma medida y $T(f)$ coincide con la integral de f en $V \cap V'$ respecto a la única medida que extiende a todas ellas. Teniendo esto en cuenta es fácil probar que T es lineal y positivo, con lo que existe una única medida de Borel regular m en S tal que

$$\int_S f dm = T(f).$$

Es claro que m extiende a la medida inducida por cualquier carta. El resto del teorema es ya inmediato. ■

La completión de la medida construida en el teorema anterior se llama a veces *medida de Lebesgue* en la variedad S . Observemos que si tomamos $S = \mathbb{R}^n$ con la carta identidad, la medida del teorema es precisamente la medida de Lebesgue en \mathbb{R}^n (pues $\Delta_X = 1$). Lo mismo es válido si S es un subespacio vectorial de \mathbb{R}^m de dimensión n .

En definitiva, hemos probado que si $K \subset S$ está contenido en la imagen de una carta X y cubrimos K por un número finito de conjuntos de Borel B_i disjuntos dos a dos, las sumas $\sum_i m(B_i^t)$, que en principio dependen de la carta X y de los puntos $p_i \in B_i$ respecto a los que se calculan las proyecciones, convergen, cuando los diámetros de los conjuntos B_i tienden a 0, a la medida de Lebesgue $m(K)$, que no depende de la carta ni de las elecciones de los puntos p_i . Podríamos generalizar esta propiedad para compactos no contenidos en el rango de una carta y para subconjuntos de Borel arbitrarios, pero no es necesario, puesto que el “volumen” que pretendemos definir sobre la variedad S debe extender a las medidas que hemos construido sobre las imágenes de las cartas (ya que cumplen la propiedad de aproximación que hemos tomado como condición necesaria para que la definición sea aceptable) concluimos que la única medida consistente con dicha condición es la medida de Lebesgue que acabamos de definir, luego es la única definición posible de volumen en una variedad.

Demostraremos únicamente la siguiente relación local entre la medida de Lebesgue y las medidas de las proyecciones en los espacios tangentes:

Teorema 9.49 *Sea $S \subset \mathbb{R}^m$ una variedad diferenciable de dimensión n y sea $X : U \rightarrow S \cap V$ una carta alrededor de un punto $p \in S$. Sea $x \in U$ tal que $X(x) = p$. Para cada conjunto de Borel $B \subset S \cap V$ sea $B^t = dX(x)[X^{-1}[B]]$. Entonces*

$$\lim_{B \rightarrow p} \frac{m(B)}{m(B^t)} = 1,$$

donde el límite ha de entenderse como sigue: Para todo $\epsilon > 0$ existe un entorno G de p en $S \cap V$ tal que si $B \subset G$ es un conjunto de Borel no nulo en S entonces $|m(B)/m(B^t) - 1| < \epsilon$.

DEMOSTRACIÓN: Dado $\epsilon > 0$, sea $\delta = (\epsilon/2)\Delta_X(x)$. Por la continuidad de X^{-1} y Δ_X en x existe un entorno G de p tal que si $y \in X^{-1}[G]$ entonces $|\Delta_X(y) - \Delta_X(x)| < \delta$. Si E es un conjunto de Borel no nulo contenido en G y $B_X = X^{-1}[B]$ tenemos que

$$m(B_X)(\Delta_X(x) - \delta) \leq m(B) = \int_{B_X} \Delta_X dm \leq m(B_X)(\Delta_X(x) + \delta),$$

luego

$$\left| \frac{m(B)}{m(B_X)} - \Delta_X(x) \right| \leq \delta < \Delta_X(x)\epsilon.$$

Por consiguiente:

$$\left| \frac{m(B)}{\Delta_X(x)m(B_X)} - 1 \right| < \epsilon,$$

pero por 9.47 tenemos que $m(B^t) = m(dX(x)[B_X]) = \Delta_X(x)m(B_X)$, con lo que

$$\left| \frac{m(B)}{m(B^t)} - 1 \right| < \epsilon. \quad \blacksquare$$

Ejemplo Si $\alpha :]a, b[\rightarrow \mathbb{R}^n$ es una curva parametrizada regular que no se corta a sí misma (de modo que su imagen es una variedad S de dimensión 1 con carta α) entonces $J_\alpha(t) = \alpha'(t)$, luego $\Delta_\alpha(t) = \|\alpha'(t)\|$ y por consiguiente

$$m(S) = \int_a^b \|\alpha'(t)\| dt$$

es la longitud de α tal y como la teníamos definida.

Si S es una superficie en \mathbb{R}^3 entonces el elemento de superficie se suele representar por $d\sigma = \sqrt{EG - F^2} dm$. Por consiguiente el área de una región C de S cubierta por la carta puede calcularse como

$$A = \int_{X^{-1}(C)} \sqrt{EG - F^2} dudv. \quad \blacksquare$$

Ejemplo Vamos a calcular el área la superficie de revolución determinada por

$$X = (r(u) \cos v, r(u) \sin v, z(u)).$$

Sabemos que

$$E = r'(u)^2 + z'(u)^2, \quad F = 0, \quad G = r(u)^2.$$

Por lo tanto

$$A = \int_0^{2\pi} \int_{u_0}^{u_1} r(u) \sqrt{r'(u)^2 + z'(u)^2} dudv = 2\pi \int_{u_0}^{u_1} r(u) \sqrt{r'(u)^2 + z'(u)^2} du.$$

Si en particular $z(u) = u$, la fórmula se reduce a

$$A = 2\pi \int_{u_0}^{u_1} r(u) |r'(u)| du.$$

Por ejemplo, el área de la esfera

$$g(\phi, \theta) = (R \sin \phi \cos \theta, R \sin \phi \sin \theta, R \cos \phi), \quad \phi \in]0, \pi[, \theta \in]0, 2\pi[.$$

es

$$A = 2\pi \int_0^\pi R^2 \sin \phi d\phi = 4\pi R^2.$$

Más detalladamente, si hacemos $\rho = R\phi$ entonces ρ es la distancia del punto (ρ, θ) al polo norte y las coordenadas

$$x = R \operatorname{sen} \frac{\rho}{R} \cos \theta, \quad y = R \operatorname{sen} \frac{\rho}{R} \operatorname{sen} \theta, \quad z = R \cos \frac{\rho}{R},$$

son el análogo esférico a las coordenadas polares en el plano. El área de un círculo esférico de radio (esférico) r es

$$A_r = 2\pi \int_0^r R^2 \operatorname{sen} \frac{\rho}{R} d\rho = 2\pi R^2 \left(1 - \cos \frac{r}{R}\right) = 4\pi R^2 \operatorname{sen}^2 \frac{r}{2R}.$$

Así, si $r = \pi R$ recuperamos el área de la esfera $4\pi R^2$, si r es pequeño con respecto a R entonces $\operatorname{sen}(r/R) \approx r/R$ y por consiguiente $A_r \approx \pi r^2$, el área del círculo plano del mismo radio. ■

El teorema siguiente nos conecta el teorema de Fubini con la integración en variedades:

Teorema 9.50 *Si $S_1 \subset \mathbb{R}^{m_1}$ y $S_2 \subset \mathbb{R}^{m_2}$ son variedades diferenciables, entonces la medida de Lebesgue en $S_1 \times S_2$ (restringida a los conjuntos de Borel) es el producto de las medidas de Lebesgue de S_1 y S_2 (sobre los conjuntos de Borel).*

DEMOSTRACIÓN: Sean m , m_1 y m_2 las medidas de Lebesgue en $S_1 \times S_2$, S_1 y S_2 respectivamente. Basta probar que si A_1 y A_2 son conjuntos de Borel en S_1 y S_2 entonces $m(A_1 \times A_2) = m_1(A_1)m_2(A_2)$. Es fácil ver que A_1 y A_2 se descomponen en una unión numerable disjunta de conjuntos de Borel, cada uno de los cuales está contenido en el rango de una carta. También es claro que si probamos la igualdad anterior para los productos de estos abiertos de ahí se sigue el caso general. En definitiva, podemos suponer que A_1 está contenido en el rango de una carta X_1 y A_2 está contenido en el rango de una carta X_2 . Una simple comprobación nos da que $\Delta_{X_1 \times X_2} = \Delta_{X_1} \Delta_{X_2}$, luego aplicando el teorema de Fubini concluimos que

$$\begin{aligned} m(A_1 \times A_2) &= \int_{X_1^{-1}[A_1] \times X_2^{-1}[A_2]} \Delta_{X_1} \Delta_{X_2} dx_1 \cdots dx_{n_1+n_2} \\ &= \left(\int_{X_1^{-1}[A_1]} \Delta_{X_1} dx_1 \cdots dx_{n_1} \right) \left(\int_{X_2^{-1}[A_2]} \Delta_{X_2} dx_{n_1+1} \cdots dx_{n_1+n_2} \right) \\ &= m_1(A_1) m_2(A_2). \end{aligned}$$

■

Terminamos la sección con una interpretación de la curvatura de Gauss de una superficie. De hecho se trata de la definición de curvatura que adoptó el propio Gauss.

Teorema 9.51 *Sea S una superficie y p un punto en el que la curvatura no sea nula. Sea n una determinación del vector normal en un entorno de p . Entonces*

$$|K(p)| = \lim_{E \rightarrow p} \frac{m(n[E])}{m(E)},$$

donde el límite se entiende en el mismo sentido que en el teorema 9.49.

DEMOSTRACIÓN: Por el teorema 6.29 sabemos que $K(p)$ es el determinante de $dn(p)$, luego ésta diferencial es un isomorfismo. Sea g una carta alrededor de p y sea $h = g \circ n$. Es fácil ver que h puede restringirse hasta una carta alrededor de $n(p)$ en la esfera unidad. Del teorema 9.49 se sigue que

$$\lim_{E \rightarrow p} \frac{m(E)}{m(E_t)} = 1, \quad \lim_{E \rightarrow p} \frac{m(n[E])}{m(n[E]_t)} = 1,$$

y por otra parte $n[E]_t = dn(p)[E_t]$, luego $m(n[E]_t) = |K(p)|m(E_t)$, de donde se sigue claramente el teorema. ■

Apéndice A

La completión de un espacio métrico

Vamos a probar que todo espacio métrico M se puede sumergir como subespacio denso de un espacio métrico completo \overline{M} , que es único salvo isometría. Si además M tiene estructura de cuerpo métrico o de cuerpo ordenado arquimediano, lo mismo valdrá para su completión \overline{M} . En particular, la completión de \mathbb{Q} será un cuerpo ordenado completo \mathbb{R} , que podemos tomar como cuerpo de los números reales.

Empezamos demostrando un par de resultados técnicos que nos van a hacer falta. El primero tiene que ver con lo que sucede al formar la sucesión de inversos de una sucesión de Cauchy, lo cual obliga a tomar precauciones con los términos que sean cero:

Teorema A.1 *Si $\{x_n\}_{n=0}^{\infty}$ es una sucesión de Cauchy en un cuerpo métrico y no converge a 0, entonces existe un $T \in \mathbb{R}$ y un $m \in \mathbb{N}$ tal que si $n \geq m$ entonces $0 < T \leq |x_n|$. Además, la sucesión $\{y_n\}_{n=0}^{\infty}$ dada por*

$$y_n = \begin{cases} 1/x_n & \text{si } x_n \neq 0, \\ 0 & \text{si } x_n = 0, \end{cases}$$

también es de Cauchy.

DEMOSTRACIÓN: Si no existe el T indicado, tomando $\epsilon = T$ obtenemos que para todo $\epsilon > 0$ y todo $m \in \mathbb{N}$ existe un $n \geq m$ tal que $|x_n| < \epsilon$.

Definimos $\{n_k\}_{k=0}^{\infty}$ por recurrencia estableciendo que n_k sea el menor número natural mayor que los ya definidos y tal que $|x_{n_k}| < 1/(k+1)$. Es claro entonces que $\{x_{n_k}\}_{k=0}^{\infty}$ es una subsucesión de $\{x_n\}_{n=0}^{\infty}$ convergente a 0, pero entonces, por 1.30 tenemos que la sucesión dada converge a 0, contradicción.

Para la segunda parte observamos que si $n, n' \geq m$ (el m dado por la primera parte) entonces $x_n, x_{n'} \neq 0$, luego

$$|y_{n'} - y_n| = \left| \frac{1}{x_{n'}} - \frac{1}{x_n} \right| = \frac{|x_n - x_{n'}|}{|x_{n'} x_n|} \leq \frac{|x_n - x_{n'}|}{T^2}.$$

Así, dado $\epsilon > 0$, podemos tomar un m mayor (si es preciso) que el de la primera parte de modo que si $n, n' \geq m$ se cumpla que $|x_n - x_{n'}| < T^2\epsilon$, y así concluimos que $|y_{n'} - y_n| < \epsilon$. ■

El segundo resultado previo nos simplifica considerablemente el problema que tenemos planteado:

Teorema A.2 *Sea M un espacio métrico y sea $D \subset M$ un conjunto denso tal que toda sucesión de Cauchy en D converge en M . Entonces M es un espacio métrico completo.*

DEMOSTRACIÓN: Sea $\{x_n\}_{n=0}^\infty$ una sucesión de Cauchy en M . Para cada $n \in \mathbb{N}$, sea $d_n \in D$ tal que $d(x_n, d_n) < 1/(n+1)$. La sucesión $\{d_n\}_{n \in \mathbb{N}}$ es de Cauchy, pues

$$|d_{n'} - d_n| \leq |d_{n'} - x_{n'}| + |x_{n'} - x_n| + |x_n - d_n|,$$

luego, dado $\epsilon > 0$, existe un $m > 3/\epsilon$ tal que si $n, n' \geq m$ se cumple que $|x_{n'} - x_n| < \epsilon/3$, y entonces

$$|d_{n'} - d_n| < \frac{1}{n'+1} + \frac{\epsilon}{3} + \frac{1}{n+1} < \frac{\epsilon}{3} + \frac{\epsilon}{3} + \frac{\epsilon}{3} = \epsilon.$$

Por hipótesis existe $\lim_n d_n = l$ y basta probar que $\lim_n x_n = l$. En efecto:

$$|x_n - l| \leq |x_n - d_n| + |d_n - l|.$$

Dado $\epsilon > 0$, basta tomar $m \geq 2/\epsilon$ tal que si $n \geq m$ se cumpla $|d_n - l| < \epsilon/2$, y entonces $|x_n - l| < \epsilon$. ■

Así pues, si partimos de un espacio métrico M y lo sumergimos como subconjunto denso de otro espacio métrico \overline{M} , para probar que \overline{M} es completo sólo tenemos que asegurar que las sucesiones de Cauchy de M son convergentes en \overline{M} , sin necesidad de preocuparnos de las nuevas sucesiones de Cauchy que pueda contener \overline{M} , que no estaban en M .

La estrategia para construir \overline{M} consiste en observar que queremos añadirle un límite a cada sucesión de Cauchy de M que no converja, pero no uno distinto a cada una, sino que varias sucesiones de Cauchy pueden estar obligadas a converger al mismo límite. Por ello, definiremos una relación de equivalencia en el conjunto de todas las sucesiones de Cauchy de M que se cumpla cuando dos de ellas “deberían” converger al mismo límite, y definiremos \overline{M} como el conjunto cociente determinado por la relación.

Definición A.3 Si M es un espacio métrico, llamaremos C_M al conjunto de todas las sucesiones de Cauchy en M . Definimos en C_M la relación de equivalencia dada por

$$\{x_n\}_{n=0}^\infty \sim \{y_n\}_{n=0}^\infty \leftrightarrow \lim_n d(x_n, y_n) = 0.$$

Es fácil ver que se trata, en efecto, de una relación de equivalencia. Por ejemplo, la transitividad se debe a que si

$$\lim_n d(x_n, y_n) = \lim_n d(y_n, z_n) = 0,$$

dado $\epsilon > 0$, existe un $m \in \mathbb{N}$ tal que si $n \geq m$ se cumple $d(x_n, y_n) < \epsilon/2$, $d(y_n, z_n) < \epsilon/2$, con lo que

$$d(x_n, z_n) \leq d(x_n, y_n) + d(y_n, z_n) < \epsilon/2 + \epsilon/2 = \epsilon,$$

luego $\lim_n d(x_n, z_n) = 0$.

Definimos la *compleción* de M como el conjunto cociente \overline{M} de C_M respecto de la relación de equivalencia que acabamos de definir.

Definimos $i : M \rightarrow \overline{M}$ como la aplicación que a cada $x \in M$ le asigna la clase de equivalencia de la sucesión constante $\{x\}_{n=0}^\infty$. Es una aplicación inyectiva, pues si $i(x) = i(y)$, entonces la sucesión constante $\{d(x, y)\}_{n=0}^\infty$ converge a 0, lo cual equivale a que $d(x, y) = 0$, luego $x = y$.

Así podemos identificar a M con un subconjunto de \overline{M} .

La compleción de un espacio métrico Si suponemos que el cuerpo R en el que la distancia toma sus valores es completo, entonces es fácil convertir a \overline{M} en espacio métrico. Basta observar que, por una parte, si $\{x_n\}_{n=0}^\infty, \{y_n\}_{n=0}^\infty$ son sucesiones de Cauchy en M , entonces $\{d(x_n, y_n)\}_{n=0}^\infty$ es una sucesión de Cauchy en R , pues

$$\begin{aligned} |d(x_n, y_n) - d(x_{n'}, y_{n'})| &\leq |d(x_n, y_n) - d(x_{n'}, y_n)| + |d(x_{n'}, y_n) - d(x_{n'}, y_{n'})| \\ &\leq d(x_n, x_{n'}) + d(y_n, y_{n'}), \end{aligned}$$

luego, dado $\epsilon > 0$, podemos tomar $m \in \mathbb{N}$ tal que, para $n, n' \geq m$ se cumpla $d(x_n, x_{n'}) < \epsilon/2$ y $d(y_n, y_{n'}) < \epsilon/2$, con lo que $|d(x_n, y_n) - d(x_{n'}, y_{n'})| < \epsilon$. Por lo tanto, si R es completo podemos definir

$$d([\{x_n\}_{n=0}^\infty], [\{y_n\}_{n=0}^\infty]) = \lim_n d(x_n, y_n).$$

En realidad, para que esta definición sea correcta hemos de comprobar que no depende de la elección de las sucesiones que tomamos en cada clase de equivalencia. Esto se debe a que si

$$[\{x_n\}_{n=0}^\infty] = [\{x'_n\}_{n=0}^\infty], \quad [\{y_n\}_{n=0}^\infty] = [\{y'_n\}_{n=0}^\infty],$$

entonces

$$\begin{aligned} |d(x_n, y_n) - d(x'_n, y'_n)| &\leq |d(x_n, y_n) - d(x_n, y'_n)| + |d(x_n, y'_n) - d(x'_n, y'_n)| \\ &\leq d(y_n, y'_n) + d(x_n, y'_n) \end{aligned}$$

y, teniendo en cuenta que los dos últimos sumandos tienden a 0, es fácil ver que

$$\lim_n (d(x_n, y_n) - d(x'_n, y'_n)) = 0,$$

de donde $\lim_n d(x_n, y_n) = \lim_n d(x'_n, y'_n)$.

Así pues, tenemos una aplicación $d : \overline{M} \times \overline{M} \rightarrow R$, y es fácil comprobar que es una distancia. Por ejemplo, para probar la desigualdad triangular, dados $\alpha = \{x_n\}_{n=0}^\infty$, $\beta = \{y_n\}_{n=0}^\infty$, $\gamma = \{z_n\}_{n=0}^\infty$ en \overline{M} , observamos que

$$0 \leq d(x_n, y_n) + d(y_n, z_n) - d(x_n, z_n),$$

luego por el teorema 1.16 (tomando como b cualquier cota de la sucesión, que es de Cauchy), obtenemos que $0 \leq d(\alpha, \beta) + d(\beta, \gamma) - d(\alpha, \gamma)$.

Además, la aplicación $i : M \rightarrow \overline{M}$ es una inmersión isométrica. En efecto, si $x, y \in M$, es claro que la sucesión constante $\{d(x, y)\}_{n=0}^\infty$ converge a $d(x, y)$, luego $d(i(x), i(y)) = d(x, y)$.

Por último, observamos que si $\{x_n\}_{n=0}^\infty$ es una sucesión de Cauchy en M y $\alpha = \{x_n\}_{n=0}^\infty$, entonces $\lim_n i(x_n) = \alpha$.

En efecto, por definición $d(i(x_n), \alpha) = \lim_{n'} d(x_n, x_{n'})$, luego, dado $\epsilon > 0$, existe un $m \in \mathbb{N}$ tal que si $n, n' \geq m$ entonces $0 \leq d(x_n, x_{n'}) < \epsilon/2$, luego por el teorema 1.16 tenemos que $0 \leq d(i(x_n), \alpha) \leq \epsilon/2 < \epsilon$, y esto es lo que había que probar.

En definitiva, lo que sucede es que cada sucesión de Cauchy en M , identificada con una sucesión de Cauchy en \overline{M} a través de i , converge a la clase de equivalencia que ella misma determina.

En particular vemos que todo elemento de \overline{M} es el límite de una sucesión en $i[M]$, luego $i[M]$ es denso en \overline{M} , y toda sucesión de Cauchy en $i[M]$ es de la forma $\{i(x_n)\}_{n=0}^\infty$, para una única sucesión de Cauchy $\{x_n\}_{n=0}^\infty$ en M , luego acabamos de ver que converge en \overline{M} , luego el teorema A.2 nos da que \overline{M} es un espacio métrico completo.

El único problema de esta demostración es que hemos supuesto que el cuerpo R es completo, luego en principio no vale cuando $R = \mathbb{Q}$, por ejemplo.

La completión de un cuerpo métrico Volvamos ahora al caso general en que R es un cuerpo ordenado arquimediano no necesariamente completo y veamos qué sucede en el caso en que partimos de un cuerpo métrico K .

Observamos entonces que el conjunto C_K de las sucesiones de Cauchy en K tiene estructura de anillo con la suma y el producto dadas por

$$\{x_n\}_{n=0}^\infty + \{y_n\}_{n=0}^\infty = \{x_n + y_n\}_{n=0}^\infty, \quad \{x_n\}_{n=0}^\infty \cdot \{y_n\}_{n=0}^\infty = \{x_n y_n\}_{n=0}^\infty.$$

Aquí usamos que la suma y el producto de sucesiones de Cauchy es una sucesión de Cauchy.

Por otra parte, el conjunto I_K de las sucesiones convergentes a 0 es un ideal de C_K , pues ciertamente contiene a la sucesión nula (que es el neutro de C_K), la suma de sucesiones convergentes a 0 converge a $0 + 0 = 0$ y si $\{x_n\}_{n=0}^\infty \in C_K$ e $\{y_n\}_{n=0}^\infty \in I_K$, entonces $\{x_n y_n\}_{n=0}^\infty \in I_K$, porque $\{x_n\}_{n=0}^\infty$ está acotada por ser de Cauchy, y basta aplicar el teorema 1.15.

Ahora observamos que la relación de equivalencia que hemos definido en C_K es la dada por

$$\{x_n\}_{n=0}^\infty \sim \{y_n\}_{n=0}^\infty \leftrightarrow \lim_n (x_n - y_n) = 0 \leftrightarrow \{x_n\}_{n=0}^\infty - \{y_n\}_{n=0}^\infty \in I_K,$$

es decir, se trata de la congruencia usual módulo el ideal I_K , luego el cociente $\bar{K} = C_K/I_K$ tiene estructura de anillo, de modo que si tenemos dos clases $\alpha = [\{x_n\}_{n=0}^\infty]$, $\beta = [\{y_n\}_{n=0}^\infty]$, entonces

$$\alpha + \beta = [\{x_n + y_n\}_{n=0}^\infty], \quad \alpha\beta = [\{x_n y_n\}_{n=0}^\infty].$$

Los neutros 0 y 1 de \bar{K} son las clases de las sucesiones constantes correspondientes, es decir, $i(0)$ e $i(1)$.

Veamos que \bar{K} es un cuerpo. Para ello tomamos un $\alpha \in \bar{K}$, $\alpha \neq 0$. Esto quiere decir que $\alpha = [\{x_n\}_{n=0}^\infty]$, donde la sucesión $\{x_n\}_{n=0}^\infty$ es de Cauchy en K , pero no converge a 0. Por el teorema A.1 sabemos que la sucesión $\{1/x_n\}_{n=0}^\infty$ (definida como 0 cuando $x_n = 0$, cosa que sólo puede suceder en un número finito de casos) es de Cauchy, luego define un $\beta \in K$, de modo que $\alpha\beta$ es la clase de equivalencia de una sucesión que vale 1 salvo a lo sumo en un número finito de casos. Es claro entonces que $\alpha\beta = 1$, luego α tiene inverso y \bar{K} es un cuerpo.

Se comprueba trivialmente que la aplicación $i : K \rightarrow \bar{K}$ es un monomorfismo de cuerpos.

Nuevamente, si suponemos que el cuerpo R es completo, podemos convertir a \bar{K} en un cuerpo métrico. Por el caso general en que M era un espacio métrico arbitrario, sabemos que si $\{x_n\}_{n=0}^\infty$ es una sucesión de Cauchy en M entonces la sucesión $\{|x_n|\}_{n=0}^\infty$ es de Cauchy en R (pues $|x_n| = d(x_n, 0)$), luego podemos definir $|\{x_n\}_{n=0}^\infty| = \lim_n |x_n|$, y ya hemos visto que la definición es correcta en el sentido de que no depende de la elección de la sucesión en la clase de equivalencia.

Ahora podemos probar que $|\cdot| : \bar{K} \rightarrow R$ es un valor absoluto en \bar{K} . Por ejemplo, si $\alpha = [\{x_n\}_{n=0}^\infty]$, $\beta = [\{y_n\}_{n=0}^\infty]$, entonces $|x_n + y_n| \leq |x_n| + |y_n|$, luego $|x_n| + |y_n| - |x_n + y_n| \geq 0$ y por 1.16 resulta que $|\alpha| + |\beta| - |\alpha + \beta| \geq 0$, luego tenemos la desigualdad triangular. Por otra parte,

$$|\alpha\beta| = \lim_n |x_n y_n| = \lim_n |x_n| |y_n| = (\lim_n |x_n|)(\lim_n |y_n|) = |\alpha||\beta|.$$

La distancia en \bar{K} definida a partir del valor absoluto es la misma que ya teníamos definida al considerar a K como un mero espacio métrico, luego ya sabemos que \bar{K} es un cuerpo métrico completo.

La completión de un cuerpo ordenado El problema sigue siendo que para definir el valor absoluto en \overline{K} hemos tenido que suponer que R es completo. Consideremos ahora el caso de un cuerpo ordenado K y vamos a probar que \overline{K} puede convertirse en un cuerpo ordenado completo sin suponer que R es completo. Sabemos que \overline{K} tiene estructura de cuerpo (pues en esta parte no hemos supuesto que R fuera completo).

Para definir un orden en \overline{K} diremos que $\alpha \in \overline{K}$ es *positivo* si $\alpha = [\{x_n\}_{n=0}^\infty]$ y existe un $c \in K$, $c > 0$ y un $m \in \mathbb{N}$ de modo que para todo $n \geq m$ se cumpla $x_n \geq c$.

Esta propiedad no depende de la sucesión elegida en α , pues si $\alpha = [\{y_n\}_{n=0}^\infty]$, entonces, para m es suficientemente grande, tenemos $x_n \geq c$ y $|x_n - y_n| < c/2$, y entonces tiene que ser $y_n \geq c/2$, pues si fuera $y_n < c/2 < c \leq x_n$, tendríamos que $|x_n - y_n| = x_n - y_n \geq c - c/2 = c/2$.

Es trivial que la suma y el producto de elementos positivos es positiva. Además, todo $\alpha \in \overline{K}$ se encuentra en uno y sólo uno de los tres casos siguientes: α es positivo, $\alpha = 0$ o bien $-\alpha$ es positivo.

En efecto, $\alpha = 0$ no es positivo, porque $\alpha = [\{0\}_{n=0}^\infty]$ y la sucesión nula no cumple la definición.

Si $\alpha \neq 0$, entonces $\alpha = [\{x_n\}_{n=0}^\infty]$ (y $-\alpha = [\{-x_n\}_{n=0}^\infty]$) donde, por A.1, existe un $c \in K$, $c > 0$ y un $m \in \mathbb{N}$, de modo que si $n \geq m$ se cumple $|x_n| \geq c$, es decir, $x_n \geq c > 0$, o bien $x_n \leq -c < 0$.

Si se da el primer caso para todo n suficientemente grande, entonces α es positivo y $-\alpha$ no lo es. Si se da el segundo caso para todo n suficientemente grande entonces $-\alpha$ es positivo y α no lo es. Sólo falta probar que no pueden darse los dos casos para n grande, es decir, que no puede ocurrir que para todo $m \in \mathbb{N}$ existan $n, n' \geq m$ tales que $x_n \geq c$ y $x_{n'} \leq -c$. Si ocurriera esto, entonces $|x_n - x_{n'}| \geq 2c$ y no se cumpliría la definición de sucesión de Cauchy para $\epsilon = 2c$.

Por lo tanto, si definimos en \overline{K} la relación dada por

$$\alpha < \beta \leftrightarrow \beta - \alpha \text{ es positivo,}$$

tenemos que se trata de una relación de orden estricto, pues no puede suceder que $\beta - \alpha$ y $\alpha - \beta$ sean ambos positivos y si $\beta - \alpha$ y $\gamma - \beta$ son positivos, también lo es la suma $\gamma - \alpha$, lo que nos da la transitividad. Además es una relación de orden total, pues o bien $\alpha - \beta$ es positivo, o bien lo es $\beta - \alpha$, o bien $\alpha - \beta = 0$, es decir, $\alpha < \beta \vee \beta < \alpha \vee \alpha = \beta$.

Más aún, \overline{K} cumple las dos propiedades de la definición de cuerpo ordenado, pues si $\alpha \leq \beta$, entonces $\beta - \alpha$ es positivo o nulo, luego lo mismo vale para $(\beta + \gamma) - (\alpha + \gamma)$, luego $\alpha + \gamma \leq \beta + \gamma$, y si $\alpha, \beta \geq 0$, ambos son positivos o nulos, luego lo mismo vale para $\alpha\beta \geq 0$.

Además, la aplicación $i : K \rightarrow \overline{K}$ conserva el orden, pues si $x < y$, es claro que la clase de la sucesión constante $\{y - x\}_{n=0}^\infty$ es positiva, es decir, que $i(y) - i(x)$ es positivo y, por consiguiente, $i(x) < i(y)$.

Así pues, \overline{K} es un cuerpo ordenado que tiene un subcuerpo $i[K]$ isomorfo a K como cuerpo ordenado. Veamos que si $\alpha = [\{x_n\}_{n=0}^\infty]$, entonces $\lim_n i(x_n) = \alpha$.

En efecto, dado $\epsilon > 0$ (ahora en \overline{K} , porque queremos probar una convergencia en \overline{K}), tenemos que $\epsilon = [\{e_n\}_{n=0}^\infty]$ de modo que existe un $c > 0$ en K y un $m \in \mathbb{N}$ de modo que para todo $n \geq m$ se cumpla $e_n \geq c$.

Cambiando m por otro mayor podemos suponer que si $n, n' \geq m$, entonces $|x_n - x_{n'}| < c/2$. Así $-c/2 < x_n - x_{n'} < c/2$, luego

$$e_{n'} - x_n + x_{n'} > c - c/2 = c/2.$$

Por lo tanto, $\epsilon - i(x_n) + \alpha = [\{e_{n'} - x_n + x_{n'}\}_{n'=0}^\infty]$ es positivo, luego concluimos que $i(x_n) - \alpha < \epsilon$.

Similarmente, $x_n - x_{n'} + e_{n'} > c - c/2 = c/2$, por lo que $i(x_n) - \alpha + \epsilon$ es positivo, luego $i(x_n) - \alpha > -\epsilon$. En total tenemos que $-\epsilon < i(x_n) - \alpha < \epsilon$, luego $|i(x_n) - \alpha| < \epsilon$ (siempre que $n \geq m$), luego se cumple la definición de límite.

La conclusión, como en el caso general para espacios métricos, es que $i[K]$ es denso en \overline{K} , que toda sucesión de Cauchy en $i[K]$ converge en \overline{K} y que \overline{K} es un cuerpo ordenado completo por el teorema A.2.

A partir de aquí ya podemos recoger los frutos de todo el desarrollo que hemos realizado:

Definición A.4 Llamaremos *cuerpo de los números reales* a la completión $\mathbb{R} = \overline{\mathbb{Q}}$ del cuerpo \mathbb{Q} de los números racionales. Hemos demostrado que es un cuerpo ordenado completo y, por 1.41 sabemos que es único salvo isomorfismo, es decir, que todo cuerpo ordenado completo es isomorfo a \mathbb{R} .

También sabemos que existe un monomorfismo de cuerpos $i : \mathbb{Q} \rightarrow \mathbb{R}$ que conserva el orden, por lo que en lo sucesivo identificaremos a \mathbb{Q} con $i[\mathbb{Q}]$, de modo que podemos considerar que $\mathbb{N} \subset \mathbb{Z} \subset \mathbb{Q} \subset \mathbb{R}$.

Si R es cualquier cuerpo ordenado arquimediano, hemos probado que su completión \overline{R} es un cuerpo ordenado completo, luego tiene que ser $\overline{R} \cong \mathbb{R}$, y entonces R resulta ser isomorfo a un subcuerpo de \mathbb{R} . Esto implica que no perdemos generalidad si consideramos únicamente distancias y valores absolutos de cuerpos métricos con valores en \mathbb{R} . Con este convenio adicional, los resultados que hemos obtenido nos permiten concluir:

Teorema A.5 Si M es un espacio métrico, entonces su completión \overline{M} es un espacio métrico completo tal que existe una inmersión isométrica $i : M \rightarrow \overline{M}$ de M en un subconjunto denso de \overline{M} (la cual nos permite considerar $M \subset \overline{M}$) y toda inmersión isométrica $f : M \rightarrow N$ de M en un espacio métrico completo N se extiende a una inmersión isométrica $F : \overline{M} \rightarrow N$, que será una isometría si $f[M]$ es denso en N . Lo mismo es válido para cuerpos métricos e inmersiones isométricas de cuerpos métricos.

Podemos parafrasear este teorema diciendo que la completión de un espacio métrico (resp. de un cuerpo métrico) es el menor espacio métrico (resp. cuerpo métrico) completo que lo contiene.

Apéndice B

Fracciones continuas

Todo número real está determinado por su parte entera y una sucesión de decimales (que no es única para ciertos números racionales). Ahora vamos a ver otra representación de características similares que resulta de interés en diversos contextos.

B.1 Propiedades básicas

Definición B.1 Partamos de una sucesión de enteros a_0, a_1, a_2, \dots , todos positivos salvo quizá el primero. Llamaremos

$$\begin{aligned} [a_0] &= a_0, \\ [a_0, a_1] &= a_0 + \frac{1}{a_1}, \\ [a_0, a_1, a_2] &= a_0 + \frac{1}{a_1 + \frac{1}{a_2}}, \\ [a_0, a_1, a_2, a_3] &= a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{a_3}}}, \end{aligned}$$

En general tenemos definido el número racional $[a_0, \dots, a_n]$ para todo n , que es no nulo si $n \geq 1$. Una definición formal se da por recurrencia de derecha a izquierda, es decir:

$$x_0 = a_n, \quad x_{i+1} = a_{n-(i+1)} + 1/x_i, \quad [a_0, \dots, a_n] = x_n.$$

Llamaremos $r_n = [a_0, \dots, a_n] = p_n/q_n$, donde p_n y q_n son enteros primos entre sí $q_n > 0$ (convenimos que si $a_0 = 0$, entonces $p_0 = 0$, $q_0 = 1$). Los números p_n, q_n pueden calcularse mediante una relación recurrente:

Teorema B.2 Con la notación anterior:

$$\begin{aligned} p_0 &= a_0, & q_0 &= 1, & p_1 &= a_0 a_1 + 1, & q_1 &= a_1, \\ p_n &= a_n p_{n-1} + p_{n-2}, & q_n &= a_n q_{n-1} + q_{n-2}. \end{aligned}$$

DEMOSTRACIÓN: Los casos $n = 0, 1, 2$ se comprueban directamente. Hay que probar que los valores dados por las fórmulas (en estos tres casos) son realmente primos entre sí, pero esto se ve fácilmente por los métodos usuales.

Supongámoslo cierto para $n - 1 \geq 2$ y probémoslo para n . Definimos los enteros primos entre sí

$$\frac{p'_j}{q'_j} = [a_1, \dots, a_{j+1}], \quad j = 0, 1, 2, \dots$$

Por la hipótesis de inducción aplicada a $n - 1$ se cumplen las fórmulas

$$p'_{n-1} = a_n p'_{n-2} + p'_{n-3}, \quad q'_{n-1} = a_n q'_{n-2} + q'_{n-3}. \quad (\text{B.1})$$

Por otra parte $\frac{p_j}{q_j} = a_0 + \frac{q'_{j-1}}{p'_{j-1}}$, luego

$$p_j = a_0 p'_{j-1} + q'_{j-1}, \quad q_j = p'_{j-1}, \quad (\text{B.2})$$

donde se ha usado que si $(p'_{j-1}, q'_{j-1}) = 1$, los valores que dan estas fórmulas también son primos entre sí. Haciendo $j = n$ en (B.2) y usando (B.1) obtenemos

$$\begin{aligned} p_n &= a_0(a_n p'_{n-2} + p'_{n-3}) + (a_n q'_{n-2} + q'_{n-3}) \\ &= a_n(a_0 p'_{n-2} + q'_{n-2}) + a_0 p'_{n-3} + q'_{n-3}, \\ q_n &= a_n q'_{n-2} + q'_{n-3}. \end{aligned}$$

Aplicando (B.2) con $j = n - 1$ y $n - 2$ se deduce

$$p_n = a_n p_{n-1} + p_{n-2}, \quad q_n = a_n q_{n-1} + q_{n-2}. \quad \blacksquare$$

De estas relaciones se sigue en particular que la sucesión q_n es creciente, y si $a_0 > 0$ entonces p_n también lo es. Veamos otra consecuencia sencilla:

Teorema B.3 *Con la notación anterior, $p_n q_{n+1} - p_{n+1} q_n = (-1)^{n+1}$ o, lo que es lo mismo: $r_n - r_{n+1} = (-1)^{n+1} / q_n q_{n+1}$.*

DEMOSTRACIÓN: Claramente

$$\begin{aligned} p_n q_{n+1} - p_{n+1} q_n &= p_n(a_{n+1} q_n + q_{n-1}) - (a_{n+1} p_n + p_{n-1}) q_n \\ &= p_n q_{n-1} - p_{n-1} q_n = -(p_{n-1} q_n - p_n q_{n-1}), \end{aligned}$$

y como $p_0 q_1 - p_1 q_0 = a_0 a_1 - (a_0 a_1 + 1) = -1$, se cumple el teorema. \blacksquare

Y así llegamos al resultado básico sobre fracciones continuas:

Teorema B.4 *Con la notación anterior, existe un único número real α tal que*

$$r_0 < r_2 < r_4 < r_6 < \dots < \alpha < \dots < r_7 < r_5 < r_3 < r_1,$$

de modo que $\alpha = \lim_n r_n$.

DEMOSTRACIÓN: Los r_n están ordenados como se indica, pues

$$r_{n+2} - r_n = r_{n+2} - r_{n+1} + r_{n+1} - r_n = (-1)^{n+1}/q_{n+1}q_{n+2} + (-1)^{n+1}/q_nq_{n+1},$$

luego la sucesión $\{r_{2n}\}$ es creciente y $\{r_{2m+1}\}$ decreciente. El teorema anterior nos da que cualquier r_{2n} es menor que cualquier r_{2m+1} , así como que sus distancias tienden a 0 (porque la sucesión q_nq_{n+1} tiende a infinito), luego r_n converge a un número α , que es el supremo de los r_{2n} y el ínfimo de los r_{2m+1} . ■

Ahora ya podemos definir lo que entendemos exactamente por “fracción continua”:

Definición B.5 La *fracción continua* asociada a la sucesión finita de números naturales no nulos a_0, \dots, a_n (salvo a_0 , que puede ser un entero arbitrario) es el número racional $[a_0, \dots, a_n]$. La *fracción continua* asociada a una sucesión infinita $\{a_n\}_{n=0}^{\infty}$ en las mismas condiciones es el número real

$$[a_0, a_1, a_2, \dots] = a_0 + \frac{1}{a_1 + \frac{1}{a_2 + \frac{1}{\ddots}}} = \lim_n [a_0, \dots, a_n]$$

dado por el teorema anterior.

Las fracciones continuas finitas $r_n = [a_0, \dots, a_n]$ se llaman *convergentes* de la fracción continua infinita $[a_0, a_1, a_2, \dots]$.

Teorema B.6 *Las fracciones continuas infinitas son números irracionales.*

DEMOSTRACIÓN: Con la notación anterior, supongamos que una fracción continua infinita es un número racional $\alpha = p/q$ (con p y q primos entre sí).

Como la sucesión q_n es creciente, existe un n tal que $q < q_{n+1}$. Puesto que α está entre r_n y r_{n+1} , se cumple que $|\alpha - r_n| \leq |r_n - r_{n+1}| = 1/q_nq_{n+1} < 1/q_nq$.

Pero por otro lado $|\alpha - r_n| = |p/q - p_n/q_n| = |pq_n - qp_n|/q_nq \geq 1/q_nq$, puesto que $p/q \neq p_n/q_n$, luego $|pq_n - qp_n| \geq 1$, contradicción. ■

El resultado que da importancia a las fracciones continuas es el que garantiza que todo número irracional positivo admite un único desarrollo en fracción continua. En efecto:

Teorema B.7 *Sea α un número real cualquiera.*

- a) *Si α es racional admite un desarrollo $\alpha = [a_0, \dots, a_n]$ en fracción continua finita.*
- b) *Si α es irracional admite un único desarrollo $\alpha = [a_0, a_1, a_2, a_3, \dots]$ en fracción continua infinita.*

DEMOSTRACIÓN: Definimos $a_0 = E(\alpha)$ (la parte entera de α). Si $\alpha \neq [a_0]$, entonces podemos escribir $\alpha = a_0 + 1/\alpha_1$ para un cierto número real positivo α_1 . Tomamos $a_1 = E(\alpha_1)$. Si $a_1 = \alpha_1$ entonces $\alpha = [a_0, a_1]$. En otro caso $\alpha_1 = a_1 + 1/\alpha_2$ para cierto número real positivo α_2 .

Si el proceso termina es que α es un número racional. Veamos que si no termina obtenemos una fracción continua que converge a α .

Por construcción se tiene que $\alpha = [a_0, \dots, a_n, \alpha_{n+1}]$ (notemos que el último término no es un número natural, pero la definición vale igualmente).

Es fácil ver que la función $[a_0, \dots, a_n, x]$ es monótona creciente cuando n es impar y monótona decreciente cuando n es par. Como $a_{n+1} = E(\alpha_{n+1}) < \alpha_{n+1}$, se cumple que α es mayor que todos los convergentes pares y menor que todos los impares. Esto prueba que la fracción continua converge a α .

Para probar la unicidad supongamos que dos sucesiones definen la misma fracción continua $[a_0, a_1, \dots] = [b_0, b_1, \dots]$. Entonces $a_0 \leq [a_0, a_1, \dots] \leq a_0 + 1$ e igualmente con la otra fracción. Como el límite es irracional no se dan las igualdades, luego $a_0 = E([a_0, a_1, \dots]) = E([b_0, b_1, \dots]) = b_0$.

Restando a_0 de ambas y tomando inversos resulta $[a_1, a_2, \dots] = [b_1, b_2, \dots]$. Siguiendo así llegamos a que todos los coeficientes coinciden. ■

La unicidad del desarrollo en fracción continua no vale para fracciones continuas finitas (o, equivalentemente, para números racionales). La razón esencial es la identidad

$$a + \frac{1}{1} = a + 1$$

o, equivalentemente, $[a, 1] = [a + 1]$. De aquí se sigue en general que

$$[a_0, \dots, a_n, 1] = [a_0, \dots, a_{n-1}, a_n + 1].$$

Una leve adaptación del argumento de la unicidad en el caso irracional prueba que cada número racional admite exactamente dos desarrollos en fracción continua de la forma precedente. (Por ejemplo, es claro que todo número racional admite un desarrollo terminado en 1 y, por inducción sobre su longitud, se puede probar que la única alternativa es la que proporciona la igualdad anterior.)

El teorema B.3 afirma que $|r_n - r_{n+1}| = 1/q_n q_{n+1}$ para cualquier par de convergentes consecutivos de una fracción continua. Puesto que su límite α se halla entre ambos, tenemos que

$$|\alpha - r_n| < 1/q_n a_{n+1} < 1/q_n^2.$$

Esto significa que los convergentes son buenas aproximaciones de sus límites. Podemos mejorar ligeramente este hecho observando que

$$|\alpha - r_n| + |\alpha - r_{n+1}| = |r_n - r_{n+1}| = 1/q_n q_{n+1}.$$

Cualquier par de números reales distintos cumple $xy < (x^2 + y^2)/2$, concluimos que

$$|\alpha - r_n| + |\alpha - r_{n+1}| < \frac{1}{2q_n^2} + \frac{1}{2q_{n+1}^2}.$$

Esto prueba que de cada dos convergentes consecutivos de un número irracional α , uno de ellos, p/q cumple $|\alpha - p/q| < 1/2q^2$. El resultado principal que necesitamos es el recíproco de este hecho.

Teorema B.8 Si α es un número irracional y p, q son naturales primos entre sí tales que $|\alpha - p/q| < 1/2q^2$, entonces p/q es un convergente de α .

DEMOSTRACIÓN: Vamos a probar que si p y q son enteros cualesquiera tales que $0 < q < q_{n+1}$, entonces $|q\alpha - p| \geq |q_n\alpha - p_n|$. Esto significa que el convergente n -simo es la mejor aproximación racional de α con denominador menor que q_{n+1} .

En efecto, la matriz de los coeficientes del sistema de ecuaciones

$$\begin{aligned} p &= up_n + vp_{n+1} \\ q &= uq_n + vq_{n+1} \end{aligned}$$

tiene determinante ± 1 , luego tiene una solución entera (u, v) . Por la hipótesis se ha de cumplir $u \neq 0$ y en el caso en que $v \neq 0$ entonces u y v tienen signos opuestos, y así

$$\begin{aligned} |q\alpha - p| &= |(uq_n + vq_{n+1})\alpha - (up_n + vp_{n+1})| \\ &= |u(q_n\alpha - p_n) + v(q_{n+1}\alpha - p_{n+1})| \geq |q_n\alpha - p_n|. \end{aligned}$$

Ahora, en las hipótesis del teorema, tomamos un n tal que $q_n \leq q < q_{n+1}$. Entonces

$$\left| \frac{p}{q} - \frac{p_n}{q_n} \right| \leq \left| \alpha - \frac{p}{q} \right| + \left| \alpha - \frac{p_n}{q_n} \right| = \frac{|\alpha q - p|}{q} + \frac{|\alpha q_n - p_n|}{q_n} \leq \left(\frac{1}{q} + \frac{1}{q_n} \right) |\alpha q - p|.$$

Como $q \geq q_n$ y $|\alpha q - p| < 1/2q$, concluimos que

$$\frac{|pq_n - qp_n|}{qq_n} < \frac{1}{qq_n},$$

y como el numerador es entero, ha de ser 0, o sea, p/q es el convergente n -simo. ■

Ejemplo El teorema B.7 nos permite calcular los convergentes de cualquier número real. Por ejemplo, para el caso de $\alpha_0 = \pi$ tenemos que $a_0 = E[\alpha_0] = 3$. Entonces calculamos $\alpha_1 = 1/(\alpha_0 - 3)$ y $a_1 = E[\alpha_1] = 7$, luego $\alpha_2 = 1/(\alpha_1 - 7)$ y $a_2 = E[\alpha_2] = 15$, igualmente $\alpha_3 = 1/(\alpha_2 - 15)$ y $a_3 = E[\alpha_3] = 1$, etc. Así podemos ir obteniendo:

$$\pi = [3, 7, 15, 1, 292, 1, 1, 1, 2, 1, 3, 1, 14, \dots].$$

Al igual que sucede con las cifras decimales de π , las cifras de su desarrollo en fracción continua no siguen ningún patrón conocido. El aproximante

$$[3, 7, 15, 1] = [3, 7, 16] = 3 + \frac{1}{7 + \frac{1}{16}} = \frac{355}{113} = 3.1415929 \dots$$

proporciona 6 cifras decimales exactas de π y además es muy fácil de recordar, pues formalmente se obtiene de la sucesión 113355 partiéndola por la mitad, usando la segunda mitad como numerador y la primera como denominador. Esta aproximación era conocida por el astrónomo chino Zu Chongzhi (del siglo V d.C.), a la que denominó *Milü* (aproximación fina), por contraposición a $[3, 7] = 22/7 = 3.142 \dots$, que era la *Yuelü* (aproximación bruta). ■

Terminamos la sección con un resultado sencillo, pero útil en la manipulación de fracciones continuas.

Teorema B.9 *Sea $\alpha = [a_0, a_1, a_2, \dots]$ y sea $\beta = [a_{n+1}, a_{n+2}, a_{n+3}, \dots]$, para $n \geq 1$. Entonces se cumple que*

$$\alpha = \frac{\beta p_n + p_{n-1}}{\beta q_n + q_{n-1}}.$$

DEMOSTRACIÓN: La prueba consiste simplemente en observar que en la demostración del teorema B.2 no se ha usado que los coeficientes a_n sean enteros salvo para probar que $(p_n, q_n) = 1$. Por lo tanto podemos aplicarlo a $\alpha = [a_0, \dots, a_n, \beta]$ y concluir que, aunque ahora p_{n+1} y q_{n+1} no sean números racionales,

$$\alpha = \frac{p_{n+1}}{q_{n+1}} = \frac{\beta p_n + p_{n-1}}{\beta q_n + q_{n-1}}$$

■

B.2 Desarrollos de irracionales cuadráticos

Los números irracionales con los desarrollos más simples en fracción continua son los irracionales cuadráticos, pues vamos a ver que se caracterizan por que sus desarrollos en fracción continua son finalmente periódicos:

Teorema B.10 *Un número irracional α es cuadrático si y sólo si los coeficientes de su fracción continua se repiten periódicamente a partir de un cierto término.*

DEMOSTRACIÓN: Supongamos que los coeficientes de la fracción continua de α se repiten a partir de un cierto término.

Puesto que $[a_0, a_1, a_2, \dots] = a_0 + 1/[a_1, a_2, \dots]$, es claro que uno es cuadrático si y sólo si lo es el otro, luego podemos suponer que los coeficientes de α se repiten desde el primero (sin anteperiodo), o sea,

$$\alpha = [a_0, \dots, a_n, a_0, \dots, a_n, a_0, \dots, a_n, \dots].$$

El teorema anterior nos da entonces que

$$\alpha = \frac{\alpha p_n + p_{n-1}}{\alpha q_n + q_{n-1}}.$$

Operando obtenemos un polinomio de segundo grado del cual es raíz α . Observemos que la fórmula anterior no vale si el periodo tiene longitud 1, pero en tal caso también podemos considerar que el periodo tiene longitud 2.

Supongamos ahora que α es un irracional cuadrático. Digamos que α es raíz del polinomio $ax^2 + bx + c$, donde a, b, c son enteros, $a > 0$ y $d = b^2 - 4ac > 0$.

Consideremos la forma cuadrática

$$f(x, y) = ax^2 + bxy + cy^2 = (x, y) \begin{pmatrix} a & b/2 \\ b/2 & c \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}.$$

Así $f(\alpha, 1) = 0$. Para cada n consideramos también la forma

$$f_n(x, y) = f(p_n x + p_{n-1} y, q_n x + q_{n-1} y) = a_n x^2 + b_n xy + c_n y^2.$$

Equivalentemente,

$$f_n(x, y) = (x, y) \begin{pmatrix} p_n & q_n \\ p_{n-1} & q_{n-1} \end{pmatrix} \begin{pmatrix} a & b/2 \\ b/2 & c \end{pmatrix} \begin{pmatrix} p_n & p_{n-1} \\ q_n & q_{n-1} \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}.$$

Teniendo en cuenta que la matriz de cambio tiene determinante ± 1 , concluimos que el determinante $d = b_n^2/4 - a_n c_n$ es el mismo para todo n , concretamente, $d = b^2/4 - ac$.

Si llamamos $\alpha_n = [a_n, a_{n+1}, \dots]$, el teorema B.9 nos da que

$$\alpha = \frac{p_n \alpha_{n+1} + p_{n-1}}{q_n \alpha_{n+1} + q_{n-1}},$$

luego

$$\begin{aligned} 0 &= f(\alpha, 1) = \frac{1}{(q_n \alpha_{n+1} + q_{n-1})^2} f(p_n \alpha_{n+1} + p_{n-1}, q_n \alpha_{n+1} + q_{n-1}) \\ &= \frac{1}{(q_n \alpha_{n+1} + q_{n-1})^2} f_n(\alpha_{n+1}, 1), \end{aligned}$$

o sea, $f_n(\alpha_{n+1}, 1) = 0$. También se cumple que $a_n = f_n(1, 0) = f(p_n, q_n)$, $c_n = f_n(0, 1) = f(p_{n-1}, q_{n-1}) = a_{n-1}$ y $b_n^2 - 4a_n c_n = d$.

De $f(\alpha, 1) = 0$ se sigue

$$\frac{a_n}{q_n^2} = f\left(\frac{p_n}{q_n}, \frac{q_n}{q_n}\right) - f(\alpha, 1) = a \left(\left(\frac{p_n}{q_n}\right)^2 - \alpha^2 \right) + b \left(\frac{p_n}{q_n} - \alpha \right).$$

Sabemos que $|\alpha - p_n/q_n| < 1/q_n^2$, luego

$$|\alpha^2 - (p_n/q_n)^2| < \frac{|\alpha + p_n/q_n|}{q_n^2} < \frac{2|\alpha| + 1}{q_n^2}.$$

Todo esto implica que $|a_n| < |a|(2|\alpha| + 1) + |b|$, o sea, $|a_n|$ satisface una cota independiente de n . Las relaciones que hemos obtenido prueban que $|b_n|$ y $|c_n|$ también están acotadas.

Por lo tanto los polinomios $f_n(x, 1)$ varían en un conjunto finito, al igual que sus raíces, entre las que se encuentran los números α_n . En consecuencia existen naturales n y k tales que $\alpha_n = \alpha_{n+k}$, y es claro que esto implica que $a_{m+k} = a_m$ para todo $m \geq n$, o sea, los coeficientes de α se repiten periódicamente. ■

Ejemplo Consideremos el número áureo $\alpha = (1 + \sqrt{5})/2$, que es raíz del polinomio $x^2 - x - 1$. Puesto que $\alpha^2 = \alpha + 1$, resulta que $\alpha = 1 + 1/\alpha$, lo que implica claramente que su fracción continua es la más sencilla de todas:

$$\frac{1 + \sqrt{5}}{2} = [1, 1, 1, \dots]. \quad \blacksquare$$

En general, para calcular el desarrollo de un irracional cuadrático α vamos calculando sus coeficientes a_n al mismo tiempo que los restos α_n . Concretamente a_n es la parte entera de α_n y $\alpha_{n+1} = 1/(\alpha_n - a_n)$. Si tenemos la precaución de expresar siempre α_n en forma canónica, $a + b\sqrt{d}$, detectaremos cuándo α_n coincide con otro resto anterior, con lo que terminará el periodo.

Ejemplo Desarrollemos $\sqrt{19}$:

$$\begin{aligned} \alpha_0 &= \sqrt{19}, & a_0 &= 4, & \alpha_1 &= \frac{4+\sqrt{19}}{3}, & a_1 &= 2, & \alpha_2 &= \frac{2+\sqrt{19}}{5}, & a_2 &= 1, \\ \alpha_3 &= \frac{3+\sqrt{19}}{2}, & a_3 &= 3, & \alpha_4 &= \frac{3+\sqrt{19}}{5}, & a_4 &= 1, & \alpha_5 &= \frac{2+\sqrt{19}}{3}, & a_5 &= 2, \\ \alpha_6 &= 4 + \sqrt{19}, & a_6 &= 8, & \alpha_7 &= \frac{4+\sqrt{19}}{3}, & a_7 &= 2. \end{aligned}$$

Así pues, $\sqrt{19} = [4, \overline{2, 1, 3, 1, 2, 8}]$, donde la barra indica el periodo que se repite. Este número tiene un anteperiodo de longitud 1. Enseguida veremos que esto no es casual. \blacksquare

Una fracción continua es *periódica pura* si no tiene anteperiodo.

Teorema B.11 *Un irracional cuadrático α tiene fracción continua periódica pura si y sólo si $\alpha > 1$ y su conjugado $\bar{\alpha}$ (es decir, la otra raíz de polmín α) cumple $-1 < \bar{\alpha} < 0$.*

DEMOSTRACIÓN: Recordemos que el desarrollo en fracción continua se calcula partiendo de $\alpha_0 = \alpha$ y de aquí $a_n = E(\alpha_n)$, $\alpha_{n+1} = 1/(\alpha_n - a_n)$.

Por inducción es claro que $-1 < \bar{\alpha}_n < 0$. En efecto, $\bar{\alpha}_{n+1} = 1/(\bar{\alpha}_n - a_n)$ y admitiendo $-1 < \bar{\alpha}_n < 0$, tenemos $-1 - a_n < \bar{\alpha}_n - a_n < -a_n$, con lo que $-1 < -1/(a_n + 1) < \bar{\alpha}_{n+1} < -1/a_n < 0$.

Ahora, despejando en $\alpha_{n+1} = 1/(\alpha_n - a_n)$, tenemos que $-1/\bar{\alpha}_{n+1} = a_n - \bar{\alpha}_n$, y como $0 < -\bar{\alpha}_n < 1$, concluimos que $a_n = E(a_n - \bar{\alpha}_n) = E(-1/\bar{\alpha}_{n+1})$.

Por el teorema anterior sabemos que $\alpha_m = \alpha_n$ para ciertos $m < n$, luego también $1/\bar{\alpha}_m = 1/\bar{\alpha}_n$, y así $a_{m-1} = a_{n-1}$. Por lo tanto

$$\alpha_{m-1} = a_{m-1} + 1/\alpha_m = a_{n-1} + 1/\alpha_n = \alpha_{n-1}.$$

Repitiendo el argumento llegamos a que $\alpha_0 = \alpha_{n-m}$, luego la fracción es periódica pura.

Ahora supongamos que la fracción es periódica pura. Entonces a_0 coincide con un coeficiente posterior, luego $\alpha \geq a_0 \geq 1$. Por el teorema B.9 resulta que

$$\alpha = \frac{p_n \alpha + p_{n-1}}{q_n \alpha + q_{n-1}},$$

luego α es raíz del polinomio $f(x) = q_n x^2 + (q_{n-1} - p_n)x - p_{n-1}$.

Ahora bien, $\bar{\alpha}$ también es raíz de este polinomio, y $f(0) = -p_{n-1} < 0$, $f(-1) = p_n - p_{n-1} + q_n - q_{n-1} > 0$, por el teorema B.2, luego $-1 < \bar{\alpha} < 0$. ■

Si d no es un cuadrado perfecto, entonces el conjugado de $E(\sqrt{d}) + \sqrt{d}$ es $E(\sqrt{d}) - \sqrt{d}$, que claramente está entre -1 y 0 , luego $E(\sqrt{d}) + \sqrt{d}$ tiene un desarrollo periódico puro. Por lo tanto el desarrollo de \sqrt{d} tiene exactamente una cifra de anteperiodo.

B.3 Transformaciones modulares

Seguidamente investigamos cuándo dos irracionales tienen fracciones continuas finalmente iguales. Veremos que esto sucede cuando son equivalentes en el sentido siguiente:

Definición B.12 Dos números α y β son *equivalentes* si existen enteros a, b, c, d tales que

$$\alpha = \frac{a\beta + b}{c\beta + d}, \quad ad - bc = \pm 1. \quad (\text{B.3})$$

Se comprueba enseguida que dos números racionales cualesquiera son equivalentes, y que un número racional nunca es equivalente a uno irracional, por lo que podemos limitarnos a considerar números irracionales.

También es fácil ver que la fórmula anterior define una biyección sobre los números irracionales. Las biyecciones de este tipo se llaman *transformaciones modulares*. Las inversas y la composición de transformaciones modulares son de nuevo transformaciones modulares, por lo que la equivalencia de números irracionales (y en general la de números reales) es una relación de equivalencia.

Los teoremas B.3 y B.9 nos dan que la transformación $\alpha = [a_0, \dots, a_n, \beta]$ es modular, dada concretamente por

$$\alpha = \frac{\beta p_n + p_{n-1}}{\beta q_n + q_{n-1}}.$$

El teorema siguiente caracteriza las transformaciones modulares que se pueden expresar de esta forma.

Teorema B.13 Si una transformación modular (B.3) cumple $c > d > 0$ entonces se puede expresar de la forma $\alpha = [a_0, \dots, a_n, \beta]$ para ciertos enteros a_0, \dots, a_n , todos positivos salvo quizá el primero.

DEMOSTRACIÓN: Hay que probar que existen a_0, \dots, a_n tales que

$$p_n = a, \quad p_{n-1} = b, \quad q_n = c, \quad q_{n-1} = d. \quad (\text{B.4})$$

Lo probaremos por inducción sobre d .

Si $d = 1$ tenemos que $a = bc \pm 1$. En el caso $a = bc + 1$ sirve $\alpha = [b, c, \beta]$. Si se cumple $a = bc - 1$, entonces $\alpha = [b - 1, 1, c - 1, \beta]$.

Supongamos ahora que $d > 1$. Aplicando el teorema B.2, las ecuaciones (B.4) equivalen a

$$p_{n-1} = b, \quad p_{n-2} = a - a_n b, \quad q_{n-1} = d, \quad q_{n-2} = c - a_n d. \quad (\text{B.5})$$

Se sigue cumpliendo $b(c - a_n d) - (a - a_n b)d = \pm 1$ para cualquier a_n , y por hipótesis de inducción (B.5) tendrá solución si garantizamos $d > c - a_n d > 0$, o equivalentemente, si $c/d > a_n > (c - d)/d$.

Notemos que c/d no puede ser entero, pues si $c = kd$ entonces $d \mid 1$. Como $c/d - (c - d)/d = 1$, podemos tomar un número natural a_n en estas condiciones y así se cumple el teorema. ■

Teorema B.14 *Dos números irracionales α y β son equivalentes si y sólo si sus desarrollos en fracción continua son finalmente iguales, es decir, si*

$$\alpha = [a_0, \dots, a_m, c_0, c_1, \dots], \quad \beta = [b_0, \dots, b_n, c_0, c_1, \dots].$$

DEMOSTRACIÓN: El teorema B.9 nos da que en estas condiciones tanto α como β son equivalentes al número $[c_0, c_1, \dots]$, luego son equivalentes entre sí. Supongamos ahora que α y β son equivalentes. Digamos que

$$\alpha = \frac{a\beta + b}{c\beta + d}, \quad ad - bc = \pm 1.$$

Podemos suponer que $c\beta + d > 0$. Sea $\beta = [b_0, \dots, b_{k-1}, \beta_k]$, donde $\beta_k = [b_k, b_{k+1}, \dots]$. Entonces:

$$\beta = \frac{\beta'_k p_{k-1} + p_{k-2}}{\beta'_k q_{k-1} + q_{k-2}}.$$

Componiendo las transformaciones modulares obtenemos que

$$\alpha = \frac{P\beta'_k + R}{Q\beta'_k + S},$$

donde

$$\begin{aligned} P &= ap_{k-1} + bq_{k-1}, \\ R &= ap_{k-2} + bq_{k-2}, \\ Q &= cp_{k-1} + dq_{k-1}, \\ S &= cp_{k-2} + dq_{k-2}, \end{aligned}$$

que son enteros y cumplen $PS - QR = \pm 1$.

Por B.3, y puesto que β se encuentra entre dos convergentes consecutivos cualesquiera, $|p_{k-1}/q_{k-1} - \beta| < 1/q_{k-1}q_k$, o sea, $|p_{k-1} - \beta q_{k-1}| < 1/q_k$. Por lo tanto $p_{k-1} = \beta q_{k-1} + \delta/q_{k-1}$, e igualmente $p_{k-2} = \beta q_{k-2} + \delta'/q_{k-2}$, con $|\delta|, |\delta'| < 1$.

De aquí resulta que

$$Q = (c\beta + d)q_{k-1} + \frac{c\delta}{q_{k-1}}, \quad S = (c\beta + d)q_{k-2} + \frac{c\delta'}{q_{k-2}}.$$

Teniendo en cuenta que $c\beta + d > 0$, es claro que haciendo k suficientemente grande podemos conseguir $Q > S > 0$. Aplicando el teorema anterior resulta que $\alpha = [a_0, \dots, a_m, \beta_k]$, de donde se sigue el teorema. ■

B.4 El espacio de Baire

El espacio de Baire es el espacio topológico $\mathcal{N} = \prod_{i \in \mathbb{N}} \mathbb{N}$, donde consideramos a \mathbb{N} como espacio topológico con la topología discreta y en \mathcal{N} consideramos la topología producto. Observemos que los elementos de \mathbb{N} son simplemente las sucesiones de números naturales o, equivalentemente, las funciones $x : \mathbb{N} \rightarrow \mathbb{N}$.

La topología de \mathcal{N} no es discreta. Como una base de \mathbb{N} está formada por los abiertos de la forma $\{n\}$, con $n \in \mathbb{N}$, el teorema 2.17 nos da que una base del espacio de Baire la forman los conjuntos de la forma

$$\{x \in \mathcal{N} \mid x(i_1) = n_1, \dots, x(i_k) = n_k\},$$

para ciertos $i_1 < \dots < i_k$ en \mathbb{N} y $n_i \in \mathbb{N}$. Ahora bien, podemos tomar una base más simple: para cada $x \in \mathcal{N}$ y cada $n \in \mathbb{N}$, llamamos $x|_n = (x(0), \dots, x(n))$ y para cada sucesión finita de números naturales $s = (s_0, \dots, s_n)$, definimos

$$B_s = \{x \in \mathcal{N} \mid x|_n = s\}.$$

Sucede que estos conjuntos son una base de \mathcal{N} . En efecto, por una parte, un conjunto de esta forma es uno de los abiertos básicos dados por el teorema 2.17 y, dado uno cualquiera de ellos y

$$x \in \{x \in \mathcal{N} \mid x(i_1) = n_1, \dots, x(i_k) = n_k\},$$

basta tomar como n el máximo de los i_k y definir $s = (s_0, \dots, s_n)$ de modo que $s_{i_j} = n_j$, y así

$$x \in B_s \subset \{x \in \mathcal{N} \mid x(i_1) = n_1, \dots, x(i_k) = n_k\},$$

lo que prueba que los conjuntos B_s forman una base por sí solos. Más aún, para cada $x \in \mathcal{N}$, la familia $\{B_{x|_n} \mid n \in \mathbb{N}\}$ es una base de entornos de x .

Esto se interpreta como que, respecto de la topología de \mathcal{N} , dos sucesiones x y y están más próximas cuanto mayor es el n tal que $x|_n = y|_n$.

El propósito de esta sección es demostrar el teorema siguiente:

Teorema B.15 *La aplicación $\phi : \mathcal{N} \rightarrow \mathbb{R}^+ \setminus \mathbb{Q}$ dada por*

$$\phi(x) = [x(0), x(1) + 1, x(2) + 1, \dots]$$

es un homeomorfismo.

DEMOSTRACIÓN: Sabemos que ϕ es biyectiva. Dado $x \in \mathcal{N}$, sea $\alpha = \phi(x)$ y sea $r_n = [x(0), x(1) + 1, \dots, x(n) + 1]$ el convergente n -ésimo de α . Fijado $\epsilon > 0$, existe un $n_0 \in \mathbb{N}$ (que podemos tomar par) tal que si $n \geq n_0$ entonces $|r_n - \alpha| < \epsilon/2$. En particular, según el teorema B.4, tenemos que

$$r_{n_0} = [x(0), \dots, x(n_0) + 1] < \phi(x) < [x(0), \dots, x(n_0 + 1) + 1] = r_{n_0+1}.$$

Entonces, si $y \in B_{x|_{n_0+1}}$, los convergentes de y hasta el correspondiente a $n_0 + 1$ son los mismos, luego

$$r_{n_0} = [y(0), \dots, y(n_0) + 1] < \phi(y) < [y(0), \dots, y(n_0 + 1) + 1] = r_{n_0+1},$$

luego $|\phi(y) - \phi(x)| < r_{n_0+1} - r_{n_0} < \epsilon$. Esto prueba que ϕ es continua en x .

Por otra parte, dado un $n_0 \in \mathbb{N}$ (que podemos suponer par sin pérdida de generalidad), sea $\delta = \min\{\alpha - r_{n_0}, r_{n_0+1} - \alpha\} > 0$. Basta probar que si $y \in \mathcal{N}$ cumple que $|\phi(y) - \alpha| < \delta$, entonces $y \in B_{x|_{n_0}}$, pues esto significa que ϕ^{-1} es continua en α . Tenemos que

$$r_{n_0} = [x(0), \dots, x(n_0) + 1] < \phi(y) < [x(0), \dots, x(n_0 + 1) + 1] = r_{n_0+1}.$$

Tomando partes enteras concluimos que $x(0) = y(0)$, de donde a su vez

$$[x(1) + 1, \dots, x(n_0 + 1) + 1] < [y(1) + 1, y(2) + 1, \dots] < [x(1) + 1, \dots, x(n_0) + 1].$$

Tomando de nuevo partes enteras llegamos a que $x(1) = y(1)$ y, tras un número finito de pasos, llegamos a que

$$[x(n_0) + 1] < [y(n_0) + 1, y(n_0 + 1) + 1, \dots] < [x(n_0) + 1, x(n_0 + 1) + 1],$$

de donde concluimos que $y(n_0) = x(n_0)$ y, por consiguiente, $y \in B_{x|_{n_0}}$. ■

Así pues, a través de los desarrollos en fracción continua, la topología de los números irracionales positivos es la misma que la del espacio de Baire. No es difícil retocar ϕ para abarcar todos los números irracionales. En efecto, basta observar que si

$$\mathcal{N}_i = \{x \in \mathcal{N} \mid x(0) \equiv i \pmod{2}\}, \quad i = 0, 1,$$

entonces $\mathcal{N} = \mathcal{N}_0 \cup \mathcal{N}_1$ es una descomposición de \mathcal{N} en abiertos disjuntos, y ambos son homeomorfos a \mathcal{N} , pues la aplicación $\phi_i : \mathcal{N}_i \rightarrow \mathcal{N}$ dada por

$$\phi_i(x) = \left[\frac{x(0) - i}{2}, x(1), x(2), \dots \right]$$

es un homeomorfismo. Por otra parte, $\mathbb{R} \setminus \mathbb{Q} = (\mathbb{R}^+ \setminus \mathbb{Q}) \cup (\mathbb{R}^- \setminus \mathbb{Q})$ es también una descomposición en abiertos disjuntos. La aplicación $\psi_0 : \mathcal{N}_0 \rightarrow \mathbb{R}^+ \setminus \mathbb{Q}$ que resulta de componer ϕ_0 con ϕ es un homeomorfismo, al igual que la aplicación $\psi_1 : \mathcal{N}_1 \rightarrow \mathbb{R}^- \setminus \mathbb{Q}$ que resulta de componer ϕ_1 con ϕ y con $x \mapsto -x$. Finalmente, los dos homeomorfismos ψ_0 y ψ_1 determinan un homeomorfismo $\psi : \mathcal{N} \rightarrow \mathbb{R} \setminus \mathbb{Q}$. Explícitamente:

Teorema B.16 La aplicación $\psi : \mathcal{N} \rightarrow \mathbb{R} \setminus \mathbb{Q}$ dada por

$$\psi(x) = \begin{cases} [\frac{x(0)}{2}, x(1) + 1, x(2) + 1, \dots] & \text{si } x(0) \text{ es par,} \\ -[\frac{x(0)-1}{2}, x(1) + 1, x(2) + 1, \dots] & \text{si } x(0) \text{ es impar,} \end{cases}$$

es un homeomorfismo.

B.5 La fracción continua de e

Como hemos indicado, el desarrollo en fracción continua de π no sigue ningún patrón reconocible. Cabría esperar que lo mismo sucede con el número e , pero no es así. En esta sección demostraremos que

$$e = [2, 1, 2, 1, 1, 4, 1, 1, 6, 1, 1, 8, \dots]$$

Fijemos un número natural m no nulo y para cada $n \geq 0$ definamos

$$\psi_n = \sum_{r=0}^{\infty} \frac{2r+2n+1}{1 \cdot 3 \cdot 5 \cdots (2r+2n+1)} \frac{2r+2}{2 \cdot 4 \cdot 6 \cdots (2r+2)} \frac{1}{m^{2r}}.$$

En primer lugar observamos que

$$\begin{aligned} \psi_0 &= \sum_{r=0}^{\infty} \frac{1}{(2r)!} \frac{1}{m^{2r}} = \frac{1}{2}(e^{1/m} + e^{-1/m}), \\ \psi_1 &= \sum_{r=0}^{\infty} \frac{1}{(2r+1)!} \frac{1}{m^{2r}} = \frac{m}{2}(e^{1/m} - e^{-1/m}). \end{aligned}$$

Comprobemos además que se cumple la relación

$$m^2\psi_n = (2n+1)m^2\psi_{n+1} + \psi_{n+2}, \quad n = 0, 1, 2, \dots \quad (\text{B.6})$$

de donde se sigue en particular que todas las series convergen. En efecto:

$$m^2\psi_n - (2n+1)m^2\psi_{n+1} = \sum_{r=0}^{\infty} \frac{(2r+2n+3)m^2 2r}{1 \cdot 3 \cdot 5 \cdots (2r+2n+3)} \frac{2r+2}{2 \cdot 4 \cdot 6 \cdots (2r+2)} \frac{1}{m^{2r}}.$$

Si eliminamos el primer sumando, que es nulo, y cambiamos el índice r por $r+1$ obtenemos la expresión que define a ψ_{n+2} .

Es claro que $\psi_n > 0$ para todo número natural n . Por lo tanto podemos definir

$$\omega_n = \frac{m\psi_n}{\psi_{n+1}}, \quad n = 0, 1, 2, \dots$$

Dividiendo entre $m\psi_{n+1}$ en (B.6) llegamos a la fórmula siguiente:

$$\omega_n = (2n+1)m + \frac{1}{\omega_{n+1}}, \quad n = 0, 1, 2, \dots$$

de donde se sigue que $\omega_n > 1$ para todo n , y que el desarrollo en fracción continua de ω_0 es

$$\omega_0 = [m, 3m, 5m, \dots].$$

Ahora bien,

$$\omega_0 = \frac{m\psi_0}{\psi_1} = \frac{e^{1/m} + e^{-1/m}}{e^{1/m} - e^{-1/m}} = \frac{e^{2/m} + 1}{e^{2/m} - 1},$$

con lo cual obtenemos en particular que

$$\frac{e + 1}{e - 1} = [2, 6, 10, 14, \dots].$$

Puesto que las fracciones continuas (infinitas) representan números irracionales, esto prueba que el número e no es racional. Más aún, que no es un irracional cuadrático, pues la fracción continua que nos ha aparecido no es periódica.

Sea ahora

$$\xi = \frac{e^{2/m} + 1}{2} = 1 + \frac{1}{\omega_0 - 1}.$$

Es inmediato que $\xi = [1, m - 1, 3m, 5m, \dots]$.

Para obtener el desarrollo en fracción continua de e necesitamos eliminar el 2 del denominador de ξ . Llamemos $\eta = e^{2/m} = 2\xi - 1$. Vamos a exponer un método general que permite calcular en muchos casos la fracción continua de un número η a partir de la fracción continua de un número ξ cuando entre ellos se da una relación del tipo

$$\eta = \frac{u\xi + v}{w},$$

donde u y w son números naturales no nulos y v es un número entero.

Antes de enunciar el resultado principal observamos que las fórmulas del teorema B.2 son válidas para $n = 0, 1$ si convenimos en que $p_{-1} = 1$, $q_{-1} = 0$, $p_{-2} = 0$, $q_{-2} = 1$.

Teorema B.17 *Sea $\xi = [a_0, a_1, a_2, \dots]$ el desarrollo en fracción continua de un irracional ξ . Sea p_n/q_n el convergente n -ésimo y $\xi_n = [a_n, a_{n+1}, a_{n+2}, \dots]$. Sea $\eta = (u\xi + v)/w$, donde u, v, w son números enteros, $u > 0$, $w > 0$, $uw = D > 1$. Para un índice cualquiera $n \geq 1$ desarrollamos el número racional*

$$\frac{u[a_0, a_1, \dots, a_{n-1}] + v}{w} = \frac{up_{n-1} + vq_{n-1}}{wq_{n-1}} = [b_0, b_1, \dots, b_{m-1}]$$

eligiendo el final de modo que $m \equiv n \pmod{2}$. Sea r_j/s_j el convergente j -ésimo de este desarrollo, de modo que en particular se tiene

$$\frac{up_{n-1} + vq_{n-1}}{wq_{n-1}} = \frac{r_{m-1}}{s_{m-1}}. \quad (\text{B.7})$$

Entonces existen números enteros u' , v' , w' tales que

$$\begin{pmatrix} u & v \\ 0 & w \end{pmatrix} \begin{pmatrix} p_{n-1} & p_{n-2} \\ q_{n-1} & q_{n-2} \end{pmatrix} = \begin{pmatrix} r_{m-1} & r_{m-2} \\ s_{m-1} & r_{m-2} \end{pmatrix} \begin{pmatrix} u' & v' \\ 0 & w' \end{pmatrix},$$

$u' > 0$, $w' > 0$, $u'w' = D$, $-w' \leq v' \leq u'$, y $\eta = [b_0, b_1, \dots, b_{m-1}, \eta_m]$, donde $\eta_m = (u'\xi_n + v')/w'$.

DEMOSTRACIÓN: La ecuación matricial equivale al siguiente sistema de ecuaciones:

$$up_{n-1} + vq_{n-1} = r_{m-1}u', \quad (\text{B.8})$$

$$wq_{n-1} = s_{m-1}u', \quad (\text{B.9})$$

$$up_{n-2} + vq_{n-2} = r_{m-1}v' + r_{m-2}w', \quad (\text{B.10})$$

$$wq_{n-2} = s_{m-1}v' + s_{m-2}w'. \quad (\text{B.11})$$

Como r_{m-1} y s_{m-1} son enteros primos entre sí, de (B.7) se sigue que los cocientes

$$\frac{up_{n-1} + vq_{n-1}}{r_{m-1}} = \frac{wq_{n-1}}{s_{m-1}}$$

son un mismo número entero u' que satisface (B.8) y (B.9). Considerando el segundo cociente concluimos que $u' > 0$.

Las ecuaciones (B.10) y (B.11) forman un sistema de ecuaciones lineales de determinante ± 1 , luego tiene solución entera v' , w' .

Tomando determinantes en la ecuación matricial llegamos a que

$$uw(-1)^{n-1} = (-1)^{m-1}u'w',$$

y puesto que $m \equiv n \pmod{2}$, podemos concluir que $D = uw = u'w'$. De aquí se deduce además que $w' > 0$. De (B.11) se sigue que

$$v' = \frac{wq_{n-1} - s_{m-2}w'}{s_{m-1}} \geq -\frac{s_{m-2}}{s_{m-1}}w' \geq -w',$$

y usando además (B.9)

$$v' = \frac{wq_{n-2} - s_{m-2}w'}{s_{m-1}} \leq \frac{w}{s_{m-1}}q_{n-2} = \frac{u'}{q_{n-1}}q_{n-2} \leq u'.$$

Por el teorema B.9 tenemos

$$\xi = \frac{p_{n-1}\xi_n + p_{n-2}}{q_{n-1}\xi_n + q_{n-2}}.$$

Haciendo uso de esto y de las ecuaciones que definen a u' , v' , w' llegamos a que

$$\begin{aligned} \eta &= \frac{u\xi + v}{w} = \frac{(up_{n-1} + vq_{n-1})\xi_n + (up_{n-2} + vq_{n-2})}{w(q_{n-1}\xi_n + q_{n-2})} \\ &= \frac{r_{m-1}u'\xi_n + r_{m-1}v' + r_{m-2}w'}{s_{m-1}u'\xi_n + s_{m-1}v' + s_{m-2}w'}, \end{aligned}$$

de donde, de acuerdo con la definición $\eta_m = (u'\xi_n + v')/w'$, se concluye

$$\eta = \frac{r_{m-1}\eta_m + r_{m-2}}{s_{m-1}\eta_m + s_{m-2}}.$$

Consecuentemente $\eta = [b_0, b_1, \dots, b_{m-1}, \eta_m]$. ■

Ahora observamos que en las hipótesis del teorema anterior se cumple

$$\eta_m = (u'\xi_n + v')/w' > v'/w' \geq -1.$$

Más aún, si $a_n \geq D$, teniendo en cuenta que a_n es la parte entera de ξ_n , de hecho

$$\eta_m = (u'\xi_n + v')/w' > (u'D + v')/w' \geq (u'^2 w' - w')/w' = u'^2 - 1 \geq 0,$$

y si $a_n \geq 2D$ entonces

$$\eta_m = (u'\xi_n + v')/w' > (u'2D + v')/w' \geq 2u'^2 - 1 \geq 1.$$

Esto es importante porque cuando $\eta_m > 1$, la relación

$$\eta = [b_0, b_1, \dots, b_{m-1}, \eta_m]$$

indica que los coeficientes de la fracción continua de η_m son la prolongación del desarrollo de η en fracción continua, que comienza con $[b_0, b_1, \dots, b_{m-1}, \dots]$.

Es fácil ver que esto sigue siendo cierto cuando $\eta_m \geq 0$ si convenimos en que

$$[\dots, a, 0, b, c, \dots] = [\dots, a + b, c, \dots].$$

Nuestra intención es partir de un número irracional ξ_0 y dividir su fracción continua en secciones

$$\xi_0 = [a_0, \dots, a_{n_1-1} \mid a_{n_1}, \dots, a_{n_2-1} \mid a_{n_2}, \dots, a_{n_3-1} \mid a_{n_3}, \dots],$$

a las que aplicar sucesivamente el teorema anterior.

Dado $\eta_0 = (u_0\xi_0 + v_0)/w_0$ tal que $u_0, w_0 > 0$ y $D = u_0w_0 > 1$, el teorema nos da números u_1, v_1, w_1 en las mismas condiciones (con el mismo D) y b_0, \dots, b_{m_1-1} tales que

$$\eta_0 = [b_0, \dots, b_{m_1-1}, \eta_{m_1}] \quad \text{con} \quad \eta_{m_1} = (u_1\xi_{n_1} + v_1)/w_1.$$

Ahora aplicamos el teorema a $\xi_{n_1} = [a_{n_1}, \dots, a_{n_2} - 1 \mid a_{n_2}, \dots, a_{n_3-1} \mid a_{n_3}, \dots]$ y obtenemos números u_2, v_2, w_2 con el mismo D y $b_{m_1}, \dots, b_{m_2-1}$ tales que

$$\eta_{m_1} = [b_{m_1}, \dots, b_{m_2-1}, \eta_{m_2}] \quad \text{con} \quad \eta_{m_2} = (u_2\xi_{n_1} + v_2)/w_2.$$

Suponiendo que $b_{m_1} \geq 0$ podemos enlazar ambos pasos y escribir

$$\eta_0 = [b_0, \dots, b_{m_1-1}, \eta_{m_1}] = [b_0, \dots, b_{m_1-1} \mid b_{m_1}, \dots, b_{m_2-1}, \eta_{m_2}].$$

A continuación aplicamos el teorema a ξ_{n_2} , y así sucesivamente. De este modo vamos obteniendo el desarrollo en fracción continua de η_0 , suponiendo que los sucesivos b_{m_i} que vamos obteniendo no sean negativos. Una forma de garantizarlo es partir la fracción original de modo que cada $a_{n_i} \geq D$, aunque no es necesario.

Con la ayuda del teorema siguiente podremos garantizar que, con las hipótesis adecuadas, al cabo de un número finito de pasos entraremos en un ciclo que nos dará una fórmula general para el desarrollo completo de η_0 . Al mismo tiempo nos dará una técnica útil para simplificar los cálculos.

Teorema B.18 *En las hipótesis del teorema B.17, si sustituimos a_0 por otro número congruente módulo D , digamos $a_0 + Dg$ (pero mantenemos los mismos a_1, \dots, a_{n-1}) entonces se obtienen los mismos números u', v', w' , así como los mismos m y b_1, \dots, b_{m-1} . El número b_0 se transforma en $b_0 + u^2g$.*

DEMOSTRACIÓN: Claramente

$$\begin{aligned} \frac{u[a_0 + Dg, a_1, \dots, a_{n-1}] + v}{w} &= \frac{u[a_0, a_1, \dots, a_{n-1}] + v}{w} + \frac{uDg}{w} \\ &= \frac{u[a_0, a_1, \dots, a_{n-1}] + v}{w} + u^2g. \end{aligned}$$

Según el teorema B.17 el desarrollo de este número es $[b_0, b_1, \dots, b_{m-1}]$, luego es inmediato que con el cambio todos los coeficientes quedan igual salvo el primero que se incrementa en u^2g .

Las relaciones recurrentes que determinan los denominadores de los convergentes no dependen del primer término de la fracción continua, luego los números q_i y s_i permanecen invariantes.

La fórmula (B.9) nos da que u' tampoco varía. Como $u'w' = D$, también w' permanece inalterado. Por último, la ecuación (B.11) garantiza la conservación de v' . ■

Con esto tenemos en realidad un método general para calcular las fracciones continuas de números η_0 a partir de números ξ_0 , pero explicaremos mejor este método aplicándolo al caso que nos interesa. Digamos sólo en general que si aplicamos sucesivamente el teorema B.17, las ternas (u_i, v_i, w_i) que vamos obteniendo varían en un conjunto finito (a causa de las restricciones que impone el teorema), luego después de un número finito de pasos volveremos a la misma terna.

Recordemos que si $\xi_0 = (e^{2/m} + 1)/2$ habíamos calculado

$$\xi_0 = [1, m-1, 3m, 5m, \dots]$$

y que $\eta_0 = e^{2/m} = 2\xi_0 - 1$. En este caso $u = 2$, $v = -1$, $w = 1$. Como $D = 2$, para obtener congruencias módulo 2 haremos $m = 2t$ (y después estudiaremos el caso $m = 2t + 1$). Dividimos la fracción de este modo:

$$\xi_0 = [1 \mid 2t-1 \mid 6t \mid 10t \mid 14t \mid \dots].$$

Vamos a aplicar el teorema B.17 a cada segmento. El teorema B.18 nos dice que podemos sustituir cada coeficiente por otro congruente módulo 2. Por ejemplo podemos considerar

$$\xi_0^* = [1 | 1 | 0 | 0 | 0 | \dots].$$

Ciertamente esto no tiene sentido como fracción continua, pero los cálculos a realizar sí lo tienen porque cada uno de ellos sólo involucra a un segmento, es decir a una fracción $[1]$ o $[0]$ que sí es correcta. Al hacer los cálculos obtendremos para cada segmento unos coeficientes $|b_{m_i}, \dots, b_{m_i+1} - 1|$, que serán los que buscamos salvo el primero. A estos primeros coeficientes tendremos que sumarles las cantidades $0, u_1^2(t-1), u_2^2 3t, u_3^2 5t, \dots$

Aplicamos el teorema B.17 al primer segmento:

$$\frac{1[1] - 1}{1} = 1 = [1] = [b_0], \quad m = 1.$$

$$\begin{pmatrix} p_0 & p_{-1} \\ q_0 & q_{-1} \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix}, \quad \begin{pmatrix} r_0 & r_{-1} \\ s_0 & s_{-1} \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix},$$

$$\begin{pmatrix} u_0 & v_0 \\ 0 & w_0 \end{pmatrix} = \begin{pmatrix} 2 & -1 \\ 0 & 1 \end{pmatrix}.$$

La ecuación matricial es

$$\begin{pmatrix} 2 & -1 \\ 0 & 1 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 2 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} u_1 & v_1 \\ 0 & w_1 \end{pmatrix},$$

y la solución:

$$\begin{pmatrix} u_1 & v_1 \\ 0 & w_1 \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix}.$$

Ahora aplicamos el teorema al segundo segmento $[1]$:

$$\frac{1[1] + 0}{1} = \frac{1}{2} = [0, 1, 1] = [b_2, b_3, b_4],$$

donde hemos tomado el desarrollo con tres cifras para que la longitud sea impar, como la de $[1]$. Ahora

$$\begin{pmatrix} r_2 & r_1 \\ s_2 & s_1 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 2 & 1 \end{pmatrix},$$

$$\begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 1 & 1 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 2 & 0 \end{pmatrix} = \begin{pmatrix} 1 & 1 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} u_2 & v_2 \\ 0 & w_2 \end{pmatrix},$$

de donde

$$\begin{pmatrix} u_2 & v_2 \\ 0 & w_2 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ 0 & 2 \end{pmatrix}$$

Sólo hay que rectificar el valor de b_2 , que en realidad es $u_1^2(t-1) = t-1 \geq 0$, luego por ahora tenemos que $\eta_0 = [1 | t-1, 1, 1 | \dots]$.

La siguiente aplicación del teorema es al segmento $[0]$:

$$\frac{1[0] - 1}{2} = -\frac{1}{2} = [-1, 1, 1] = [b_5, b_6, b_7].$$

$$\begin{pmatrix} 1 & -1 \\ 0 & 2 \end{pmatrix} \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix} = \begin{pmatrix} -1 & 1 \\ 2 & 0 \end{pmatrix} \begin{pmatrix} u_3 & v_3 \\ 0 & w_3 \end{pmatrix},$$

y esta vez llegamos a que

$$\begin{pmatrix} u_3 & v_3 \\ 0 & w_3 \end{pmatrix} = \begin{pmatrix} 1 & -1 \\ 0 & 2 \end{pmatrix} = \begin{pmatrix} u_2 & v_2 \\ 0 & w_2 \end{pmatrix},$$

El valor corregido de b_5 es $b_5 = -1 + u_2^2 3t = 3t - 1 \geq 0$.

Tenemos, pues, que $\eta_0 = [1 \mid t - 1, 1, 1 \mid 3t - 1, 1, 1 \mid \dots]$.

Ahora bien, para los cálculos relativos al cuarto segmento partimos exactamente de los mismos datos que para el tercero (la fracción $[0]$ y la terna $(u_3, v_3, w_3) = (1, -1, 2)$), luego llegaremos exactamente a los mismos coeficientes $[-1, 1, 1]$, y otra vez a la misma terna. Lo único que cambiará será la corrección del primer coeficiente, que ahora será $5t$, y después $7t$, etc., dando lugar siempre a coeficientes mayores que 0.

Consecuentemente tenemos la fracción continua de η_0 , que no es sino

$$\eta_0 = [1, t - 1, 1, 1, 3t - 1, 1, 1, 5t - 1, 1, 1, 7t - 1, 1, 1, \dots],$$

o más brevemente:

$$\sqrt[t]{e} = \eta_0 = [1, (2k + 1)t - 1, 1]_{k=0}^{\infty}.$$

En el caso $t = 1$ aparece un cero que debe ser cancelado:

$$e = [1, 0, 1, 1, 2, 1, 1, 4, 1, 1, 6, \dots] = [2, 1, 2, 1, 1, 4, 1, 1, 6, \dots],$$

así, $e = [2, \overline{1, 2k, 1}]_{k=0}^{\infty}$.

En general, este método puede ser aplicado siempre que la fracción continua de ξ_0 pueda ser dividida en segmentos que (por lo menos desde uno dado en adelante) tengan todos la misma longitud y los mismos términos, salvo quizá el primero, y de modo que los primeros términos de cada segmento sean mayores o iguales que D (para que los coeficientes que obtenemos puedan ser enlazados) y congruentes módulo D (para que podamos reducirlos a constantes por el teorema B.18 y así llegar a un ciclo como ha ocurrido en el ejemplo).

Otra aplicación la tenemos cuando hacemos $m = 2t + 1$ en la expresión original. Entonces queda

$$\xi_0 = [1 \mid 2t \mid 6t + 3 \mid 10t + 5 \mid 14t + 7 \mid \dots],$$

y con este método podemos calcular la fracción continua de $e^{2/(2t+1)}$. Para ello reducimos módulo 2 a la fracción $\xi_0^* = [1 \mid 0 \mid 1 \mid 1 \mid 1 \mid \dots]$.

Esta vez se obtienen las ternas

$$(2, -1, 1), \quad (1, 0, 2), \quad (2, 0, 1), \quad (1, 0, 2), \quad (1, -1, 2), \quad (2, 0, 1).$$

La primera repetición $(u_1, v_1, w_1) = (u_3, v_3, w_3)$ no es significativa, pues los primeros (y únicos) coeficientes de los segmentos primero y tercero son $[0]$ y $[1]$ respectivamente, luego no son congruentes y por lo tanto no podemos garantizar que comience un ciclo (y de hecho no comienza).

En cambio la repetición $(u_5, v_5, w_5) = (u_2, v_2, w_2)$ sí cierra el proceso. La fracción que se obtiene es

$$\eta_0^* = [1 \mid 0 \mid 2 \mid 0, 1, 1 \mid \dots]$$

Para corregir los primeros coeficientes observamos que al pasar de ξ_0 a ξ_0^* hemos restado $2 \cdot 0, 2t, 2(t+1), 2(5t+2), 2(7t+3), \dots$ así como que los valores de u_i son $2, 1, 2, 1, 1, 2, 1, 1, 2, 1, 1, 2, \dots$

Por lo tanto ahora hemos de sumar

$$0, t, 4(t+1), 5t+2, 7t+3, 4(9t+5), 11t+7, 13t+9, 4(15t+11), \dots$$

Omitimos los detalles, pero no es difícil llegar a que la expresión final es

$$\begin{aligned} e^{2/(2t+1)} &= [1, \overline{(1+6k)t+3k}, \overline{(12+24k)t+6+12k}, \overline{(5+6k)t+2+3k}, 1, 1] \\ &= [1, \overline{(1+6k)t+3k}, \overline{(12+24k)t+6+12k}, \overline{(5+6k)t+2+3k}, 1]_{k=0}^{\infty}. \end{aligned}$$

La fórmula se simplifica bastante en el caso $t=0$, que nos da

$$\begin{aligned} e^2 &= [1, \overline{3k}, \overline{6+12k}, \overline{2+3k}, 1]_{k=0}^{\infty} = [1, 0, \overline{6}, \overline{2+3k}, 1, 1, \overline{3+3k}, \overline{18+12k}]_{k=0}^{\infty} \\ &= [7, \overline{2+3k}, 1, 1, \overline{3+3k}, \overline{18+12k}]_{k=0}^{\infty} \end{aligned}$$

Explícitamente:

$$e^2 = [7, 2, 1, 1, 3, 18, 5, 1, 1, 6, 30, 8, 1, 1, 9, 42, 11, 1, 1, 12, 54, \dots].$$

■

Apéndice C

Resumen de dinámica clásica

Discutimos aquí los elementos básicos de la dinámica newtoniana que se han usado en muchos de los ejemplos y aplicaciones de este libro, junto con algunas aplicaciones más.

C.1 El espacio y el tiempo

El *movimiento* puede definirse como el cambio de posición de un cuerpo con el paso del tiempo. Lo primero que hemos de hacer para poder estudiar matemáticamente el movimiento es expresar numéricamente el espacio y el tiempo o, más precisamente, determinar mediante números adecuados cada posición posible en el espacio y cada instante de tiempo. Nos ocupamos primero del tiempo, porque es más sencillo.

Sistemas de referencia temporales Para que tenga sentido afirmar, por ejemplo, que un determinado suceso ocurrió en el instante $t = 10$ necesitamos establecer dos convenios arbitrarios:

- a) Hemos de establecer un *origen* de tiempos, es decir, hemos de indicar un suceso que determine el instante que vamos a etiquetar con el número $t = 0$.
- b) Hemos de establecer una *unidad de tiempo*, lo cual supone especificar un fenómeno físico que pueda repetirse cuantas veces se quiera en cualquier momento con evidencias de que va a durar siempre lo mismo.

Un aparato capaz de repetir indefinidamente el fenómeno físico que define la unidad de tiempo y registrar el número de repeticiones que se han producido hasta un instante dado es lo que comúnmente se llama un *reloj*. El número real que se asigna a cada instante es precisamente el número de veces que el reloj ha

repetido la unidad de tiempo¹ desde el instante seleccionado como $t = 0$. Los instantes anteriores al origen de tiempos se identifican con números negativos.

Aquí usaremos siempre las unidades especificadas por el *Sistema Internacional de Unidades* o *Sistema Internacional de Medidas*, el cual, como unidad de tiempo, establece el *segundo*, (abreviado s) cuya definición es la siguiente:

Un *segundo* es la duración de 9 192 631 770 oscilaciones de la radiación emitida en la transición entre los dos niveles hiperfinos del estado fundamental del isótopo 133 del átomo de cesio a una temperatura de 0 K”.

Los relojes cuyo funcionamiento se basa en esta definición de segundo —que el lector no necesita entender en absoluto— se llaman *relojes atómicos*, y son, por definición, los más precisos posibles. En principio, permiten medir el tiempo con una precisión de un 9 192 631 770-avo de segundo. El uso de cualquier reloj basado en fenómenos físicos de otro género introducirá necesariamente errores en las mediciones, si bien éstos pueden considerarse despreciables en la mayoría de los casos (aunque, por otra parte, cualquier reloj es exacto si tomamos como unidad de tiempo la que él mismo determina, al menos si podemos confiar en que es capaz de repetirla sin alterar significativamente su duración). El problema aparece entonces a la hora de comparar las mediciones realizadas con relojes diferentes.

Al fijar un origen de tiempos y una unidad de tiempo diremos que hemos fijado un *sistema de referencia temporal*, en el sentido de que afirmar que un suceso ha ocurrido en el instante $t = 10$ sólo tiene un significado físico *en referencia* a dicha elección del origen y de la unidad.

En realidad, al asignar un valor numérico a un instante dado con respecto a un sistema de referencia temporal estamos adoptando un tercer convenio arbitrario (aparte de la elección del origen y de la unidad de tiempo), y es que estamos conviniendo en representar con números positivos los instantes posteriores al instante $t = 0$, y con números negativos los anteriores. Nada impediría hacerlo al revés, e incluso podríamos considerar diferentes sistemas de referencia temporales cada uno con un criterio diferente, pero las cosas se simplifican notablemente si acordamos considerar únicamente sistemas de referencia en los que los instantes posteriores tengan asignados números reales mayores que los anteriores, y así lo haremos siempre de forma tácita.

Sistemas de referencia espaciales El espacio físico, en el contexto de la mecánica clásica, se identifica con un espacio afín tridimensional, de modo que, fijado un origen O , todo punto P del espacio está determinado por su *vector de posición* $\vec{r} = \overrightarrow{OP}$ y, si además fijamos una base ortonormal, $(O; \vec{u}_1, \vec{u}_2, \vec{u}_3)$, entonces cada punto P está determinado por sus tres coordenadas cartesianas (x, y, z) , de modo que $P = O + x\vec{u}_1 + y\vec{u}_2 + z\vec{u}_3$.

¹En realidad, si queremos medidas de mayor precisión, deberemos contar con un reloj capaz de repetir submúltiplos de la unidad de tiempo, y la precisión de la medida será a lo sumo la del menor submúltiplo de la unidad que es capaz de distinguir el reloj.

Sin embargo, para que estas afirmaciones tengan pleno sentido físico, tanto el origen de coordenadas como los vectores de la base tienen que definirse en términos de objetos físicos concretos que los determinen y además la unidad de longitud que nos permite hablar de bases ortonormales tiene que definirse en términos de algún proceso físico, al igual que hemos hecho con la unidad de tiempo. La unidad de longitud en Sistema Internacional de Unidades es el *metro*, (abreviado m) definido así:

Un *metro* es la distancia que recorre la luz en el vacío durante un tiempo de $1/299\,792\,458$ s.

Al igual que hemos convenido que el semieje temporal positivo corresponde a los sucesos futuros respecto del origen de tiempos, también es útil adoptar el convenio de considerar exclusivamente sistemas de referencia espaciales positivamente orientados, donde la orientación positiva es la determinada por la “regla de la mano derecha” (véase [G, Sección 7.6]).

A menudo se usa la palabra *observador* como sinónimo de “sistema de referencia espacial y temporal”, lo que nos dará pie a usar un lenguaje formalmente más “humano”. Así, cuando digamos que “un observador ve a tal objeto en la posición \vec{r} ” querremos decir simplemente que \vec{r} es el vector de posición del objeto en el sistema de referencia considerado. Cuando hablemos de “un observador situado en tal punto” querremos decir simplemente que el punto indicado es el origen de coordenadas del observador, etc.

Un hecho fundamental que el lector nunca debe olvidar es que el número t con que un observador identifica un instante o el vector de posición \vec{r} con el que identifica una posición en el espacio son *relativos* a dicho observador, en el sentido de que dos observadores distintos O y O' pueden identificar el mismo instante o la misma posición con números completamente distintos.

A la hora de comparar determinaciones espaciales o temporales correspondientes a observadores distintos supondremos siempre que ambos utilizan el mismo sistema de unidades. De no ser así, tendríamos que convertir unas unidades en otras antes de proceder a ninguna clase de comparación. Con este convenio, relacionar las coordenadas temporales de dos observadores O y O' es muy simple: sólo necesitamos conocer el valor t_0 que O asigna al origen temporal de O' . Entonces, la relación entre la coordenada temporal t que O asigna a un suceso y la coordenada t' que O' asigna al mismo suceso es simplemente:

$$t = t_0 + t'.$$

En la práctica, cuando tengamos que considerar varios observadores a un tiempo, siempre podremos suponer que todos ellos fijan el mismo origen de tiempos, con lo que nunca tendremos necesidad de aplicar la fórmula precedente.

La relación entre las coordenadas espaciales es más delicada. Para expresarla necesitamos conocer el vector de posición \vec{O}' que O asigna al origen de coordenadas de O' , así como los vectores de posición $\vec{P}_1, \vec{P}_2, \vec{P}_3$ de los puntos que O' identifica con los vectores de la base canónica de \mathbb{R}^3 , es decir, $(1, 0, 0)$, $(0, 1, 0)$,

$(0, 0, 1)$. Aquí hay que tener presente que tanto \vec{O}' como $\vec{P}_1, \vec{P}_2, \vec{P}_3$ pueden ser funciones de t , pues el observador O' puede estar en movimiento respecto de O . Si identificamos a \vec{O}' con sus coordenadas respecto de O y llamamos M a la matriz cuyas filas son las coordenadas de los vectores \vec{P}_i , entonces las coordenadas X y X' de un mismo punto respecto de cada observador están relacionadas en la forma

$$X = \vec{O}' + X'M.$$

Los límites de la mecánica clásica La física moderna nos enseña que todo el planteamiento que hemos establecido aquí sólo puede entenderse como una aproximación a la geometría del espacio y del tiempo que será razonable en la medida en que no pretendamos usarlo como base para analizar situaciones excesivamente alejadas de la experiencia cotidiana. Éstos son algunos de los aspectos a tener en cuenta para no incurrir en errores por el mero hecho de aplicar la mecánica clásica cuando es impropio hacerlo:

- Estamos suponiendo que el conjunto de todos los instantes puede biyectarse con la recta real, de modo que, en particular es posible avanzar y retroceder indefinidamente en el tiempo. Sin embargo, la cosmología nos enseña que el universo tiene un primer instante, de modo que es absurdo hablar de instantes anteriores a éste. Ahora bien, dicho primer instante hay que situarlo hace unos 13.000 millones de años, por lo que no es algo que deba preocuparnos a la hora de estudiar, por ejemplo, la forma en que la Tierra se mueve alrededor del Sol, dado que, por esa época, ni siquiera existía el Sol. Similarmente, algunas teorías cosmológicas predicen la existencia de un último instante tras el cual no existe más tiempo.
- También hemos supuesto que los puntos del espacio se pueden identificar con los puntos de \mathbb{R}^3 , pero la teoría de la relatividad permite que la geometría global del espacio pueda ser muy diferente de la de \mathbb{R}^3 , de modo que un cuerpo que se moviera durante suficiente tiempo siempre en la misma dirección podría acabar en el punto de partida.
- También estamos suponiendo que tiene sentido considerar distancias e intervalos de tiempo arbitrariamente pequeños, mientras que la mecánica cuántica pone límites teóricos (no meramente prácticos, debidos al margen de imprecisión que necesariamente ha de tener cualquier instrumento de medida) a la posibilidad de dividir el espacio y el tiempo. Ahora bien, para encontrarnos con estas dificultades tendríamos que ocuparnos de partículas subatómicas, cosa que en ningún momento vamos a hacer.

Todo esto es así en lo tocante a las cuestiones puramente geométricas que hemos analizado hasta ahora. Cuando entremos en consideraciones propiamente de naturaleza física, podríamos señalar otras limitaciones a su aplicabilidad, como que los principios de la mecánica clásica no son válidos para estudiar objetos que se muevan a velocidades comparables a la velocidad de la luz, pero estos hechos no deben preocuparnos porque todas las situaciones que vamos a analizar quedarán muy lejos de tales límites de aplicación de la mecánica clásica.

C.2 Cinemática de una partícula puntual

La *cinemática* es el estudio del movimiento de un cuerpo desde un punto de vista puramente descriptivo (en términos matemáticos), sin entrar en las causas que provocan o modifican el movimiento de los cuerpos. Cuando entramos en el análisis de dichas causas pasamos al dominio de la *dinámica*.

Aquí vamos a considerar objetos puntuales, es decir, objetos cuya extensión en el espacio puede considerarse despreciable y se pueden identificar con puntos matemáticos. Por ejemplo, a la hora de estudiar el movimiento de la Tierra alrededor del Sol o el movimiento de un coche por una carretera, la forma física de los objetos involucrados es irrelevante, y podemos permitirnos identificarlos con puntos.

Reposo y movimiento Fijado un observador O , el vector de posición de una partícula puntual P no tiene por qué ser el mismo en todo instante, sino que en general será una función $\vec{r}(t)$ del tiempo t . Si se trata de una función constante (al menos durante un intervalo de tiempo), diremos que P está en *reposo* (respecto del observador O y en el intervalo de tiempo considerado), mientras que si no es constante diremos que P está en *movimiento* (siempre respecto del observador fijado).

Vemos así que los conceptos de reposo y movimiento son relativos al observador, pues un mismo objeto puede estar en reposo para un observador y en movimiento para otro. Por ejemplo, los libros de mi estantería están en reposo respecto de un sistema de referencia basado en mi casa (por ejemplo, uno que tenga por origen de coordenadas una esquina de mi habitación y por ejes las tres aristas que confluyen en ella), pero están en movimiento respecto a un sistema de referencia que tenga por origen de coordenadas al Sol y cuyos ejes X - Y sean paralelos a los semiejes de la órbita elíptica que la Tierra describe alrededor del Sol. Respecto a dicho sistema de referencia, mis libros acompañan a la Tierra en su rápido movimiento de traslación.

Esto es algo más profundo de lo que a alguno le podría parecer a primera vista, ya que no hay modo alguno de dar sentido a una afirmación como que mis libros “parecen” estar en reposo cuando considero mi habitación como referencia, pero que “en realidad” se mueven porque la Tierra lo hace. En principio, sería posible postular un “sistema de referencia absoluto” cuyo origen de coordenadas y sus ejes estuvieran en reposo “de verdad”, de modo que un cuerpo está en reposo “de verdad” si lo está respecto de dicho sistema de referencia, pero cuando los físicos trataron de medir el “movimiento absoluto” de la Tierra en el espacio los experimentos demostraron que dicho concepto era inconsistente, y dicha conclusión fue el primer paso que llevó al descubrimiento de la teoría de la relatividad. Diremos algo más sobre esto en la sección siguiente. Aquí no es algo que deba preocuparnos puesto que nos limitaremos a estudiar el movimiento respecto de un sistema de referencia arbitrario.

Supondremos que la función $\vec{r}(t)$ es una función continua (lo que es tanto como afirmar que un objeto no puede desaparecer en un punto del espacio y aparecer en otro sin haber recorrido un camino que los conecte). Su imagen

puede ser un punto (si el objeto está en reposo) o bien una curva en \mathbb{R}^3 a la que llamaremos *trayectoria* del móvil. (En la práctica llamaremos también trayectoria a la propia función $\vec{r}(t)$, considerada como una curva paramétrica en \mathbb{R}^3 .)

Velocidad Si la función $\vec{r}(t)$ es derivable podemos definir la *velocidad* de un móvil como la derivada

$$\vec{v}(t) = \frac{d\vec{r}}{dt}.$$

De este modo, la velocidad de un móvil en un instante t es un vector que mide la variación instantánea de su posición. La interpretación geométrica de la derivada implica en particular que (si es no nula) la velocidad es tangente a la trayectoria.

La velocidad se mide en metros por segundo. Notemos que la velocidad, al igual que la posición, es una magnitud vectorial, de modo que afirmar que, en un momento dado, la velocidad de un móvil es $\vec{v} = (2, 1, -3)$ m/s significa que su coordenada x (para un observador dado) aumenta a razón de 2 metros por segundo, su coordenada y aumenta a razón de 1 metro por segundo y su coordenada z disminuye a razón de 3 metros por segundo.

Si suponemos que la velocidad nunca se anula y llamamos $v(t) = \|\vec{v}(t)\|$, la geometría diferencial nos da (definición 5.23) que la longitud de un arco de trayectoria, es decir, el espacio recorrido por el móvil, entre los instantes t_0 y t viene dada por

$$L(t) = \int_{t_0}^t v(t) dt,$$

de donde se sigue que

$$v(t) = \frac{dL}{dt},$$

es decir, que la *velocidad escalar* $v(t)$ representa la distancia recorrida por unidad de tiempo.

Ejemplo Si lanzamos una piedra al aire en un ángulo de 60 grados y tomamos un sistema de referencia cuyo origen esté en el punto de partida y cuyo eje X apunte en la dirección del lanzamiento (y el eje Z sea vertical), la trayectoria que describe la piedra será

$$\vec{r}(t) = (v_x t, 0, v_z t - 5t^2),$$

donde $\vec{v}_0 = (v_x, 0, v_z)$ es la velocidad con que la lanzamos.² Supongamos además que la piedra tarda 4 s en caer.

En primer lugar podemos averiguar la velocidad inicial. Que tarde 4 s en caer significa que

$$4v_z - 5 \cdot 4^2 = 0,$$

²Esto será inmediato dentro de poco. De momento nos limitaremos a estudiar la cinemática del movimiento dando por cierta esta expresión.

luego $v_z = 20$ m/s. Que el ángulo de lanzamiento sea de 60° significa que

$$\frac{v_z}{v_x} = \tan 60^\circ = \sqrt{3},$$

luego $\vec{v}_0 = (11.55, 0, 20)$ m/s y $v_0 = \|\vec{v}_0\| = 23.1$ m/s.

Derivando la posición obtenemos la velocidad:

$$\vec{v} = (11.55, 0, 20 - 10t) \text{ m/s.}$$

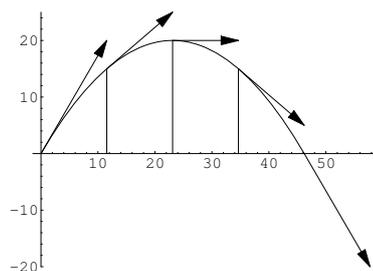
El punto donde cae la piedra es

$$\vec{r}(4) = (46.2, 0, 0) \text{ m,}$$

de modo que la piedra ha recorrido 46.2 m en horizontal, si bien la longitud de su trayectoria ha sido

$$e = \int_0^4 \sqrt{133.3 + (20 - 10t)^2} dt = \int_0^4 \sqrt{533.3 - 400t + 100t^2} dt = 63.75 \text{ m.}$$

La figura siguiente muestra la trayectoria de la piedra junto con su velocidad en los instantes $t = 0, 1, 2, 3, 4$. ■



La posición de un móvil se puede reconstruir a partir de su velocidad y de su posición en un instante dado $\vec{r}_0 = \vec{r}(t_0)$ sin más que integrar:

$$\vec{r}(t) = \vec{r}_0 + \int_{t_0}^t \vec{v}(t) dt.$$

En particular, si un móvil se mueve con velocidad constante \vec{v} , su trayectoria es $\vec{r}(t) = \vec{r}_0 + \vec{v}t$, es decir, una línea recta recorrida a velocidad constante.

Aceleración Si la función $\vec{v}(t)$ es a su vez derivable, podemos definir la *aceleración* de un móvil como

$$\vec{a}(t) = \frac{d\vec{v}}{dt} = \frac{d^2\vec{r}}{dt^2}.$$

La aceleración se mide en m/s^2 . Así, si la aceleración de un móvil en un instante dado es $\vec{a} = (2, 1, -1)$, esto significa que la primera coordenada v_x de su velocidad está aumentando a razón de 2 m/s cada segundo, etc.

Es claro que la posición $\vec{r}(t)$ puede reconstruirse a partir de $\vec{a}(t)$ sin más que conocer la posición \vec{r}_0 y la velocidad \vec{v}_0 en un instante t_0 , pues la velocidad es

$$\vec{v}(t) = \vec{v}_0 + \int_{t_0}^t \vec{a}(t) dt$$

y a partir de $\vec{v}(t)$ y \vec{r}_0 se reconstruye $\vec{r}(t)$. Por ejemplo, si un objeto se mueve con aceleración constante \vec{a} , entonces su velocidad en un instante dado viene dada por $\vec{v}(t) = \vec{v}_0 + \vec{a}t$ y, volviendo a integrar, su posición en un instante dado es $\vec{r}(t) = \vec{r}_0 + \vec{v}_0t + \frac{1}{2}\vec{a}t^2$, donde \vec{r}_0 es su posición en el instante $t = 0$.

Definimos

$$\vec{\tau}(t) = \frac{\vec{v}}{\|\vec{v}\|},$$

que es el vector unitario tangente a la trayectoria del móvil (y cuyo sentido indica la dirección de avance). Entonces $\vec{v} = v\vec{\tau}$ y, derivando:

$$\vec{a} = \frac{dv}{dt} \vec{\tau} + v \frac{d\vec{\tau}}{dt}.$$

Ahora observamos que, como $\vec{\tau} \cdot \vec{\tau} = 1$, al derivar resulta que

$$\frac{d\vec{\tau}}{dt} \cdot \vec{\tau} = 0,$$

lo que se interpreta como que la derivada de $\vec{\tau}$ es un vector perpendicular a $\vec{\tau}$. Si llamamos

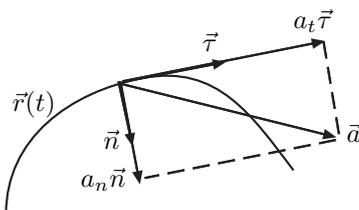
$$\vec{n} = \frac{\frac{d\vec{\tau}}{dt}}{\left\| \frac{d\vec{\tau}}{dt} \right\|},$$

tenemos entonces que

$$\vec{a} = a_t \vec{\tau} + a_n \vec{n},$$

donde

$$a_t = \frac{dv}{dt} \quad \text{y} \quad a_n = v \left\| \frac{d\vec{\tau}}{dt} \right\| \geq 0.$$



Las funciones $a_t(t)$ y $a_n(t)$ se llaman respectivamente *aceleración tangencial* y *aceleración normal*. Geométricamente son las proyecciones del vector aceleración sobre la tangente a la trayectoria y sobre su perpendicular en el plano determinado por los vectores \vec{v} y \vec{a} ; físicamente a_t es la derivada de la velocidad escalar, mientras que a_n determina la curvatura de la trayectoria.

En efecto, por una parte observamos que si la función a_n es idénticamente nula entonces la trayectoria es rectilínea, pues en tal caso

$$\frac{d\vec{\tau}}{dt} = 0$$

y, por consiguiente, $\vec{\tau}$ es constante. Así pues, $\vec{v} = v\vec{\tau}$ y

$$\vec{r}(t) = \vec{r}_0 + \left(\int_{t_0}^t v dt \right) \vec{\tau}$$

es una recta.

Por otra parte, consideremos una trayectoria circular de radio R . Su forma más general es

$$\vec{r}(t) = \vec{r}_0 - R \cos \alpha(t) \vec{i} + R \sin \alpha(t) \vec{j},$$

donde \vec{i}, \vec{j} son dos vectores unitarios ortogonales y $\alpha(t)$ es una función arbitraria cuya derivada no se anule. Entonces

$$\vec{v}(t) = R\alpha'(t) \sin \alpha(t) \vec{i} + R\alpha'(t) \cos \alpha(t) \vec{j} = R\alpha'(t) \vec{\tau}(t),$$

$$\vec{a}(t) = R\alpha''(t) \vec{\tau}(t) + R\alpha'^2(t) \vec{n}(t).$$

Observamos que $v(t) = R\alpha'(t)$ y $a_n(t) = R\alpha'^2(t) = v^2(t)/R$. Así pues, la aceleración normal de un movimiento circular viene dada por

$$a_n = v^2/R.$$

Más aún, fijado un tiempo t_0 y tres números reales $v_0 \neq 0$, a_{t_0} , a_{n_0} , con $a_{n_0} > 0$, podemos elegir $R = v_0^2/a_{n_0}$ y una función α tal que

$$\alpha(t_0) = 0, \quad \alpha'(t_0) = a_{n_0}/v_0, \quad \alpha''(t_0) = a_{t_0}/R,$$

de modo que

$$\vec{\tau} = \vec{j}, \quad \vec{n} = \vec{i}, \quad \vec{v}(t_0) = v_0 \vec{\tau}, \quad \vec{a}(t_0) = a_{t_0} \vec{\tau} + a_{n_0} \vec{n}.$$

Además, eligiendo adecuadamente \vec{r}_0 podemos hacer que $\vec{r}(t_0)$ sea un punto arbitrario.

En definitiva, concluimos que si $\vec{r}(t)$ es una trayectoria arbitraria con velocidad y aceleración normal no nulas en un instante t_0 , podemos construir otra trayectoria circular que coincida con $\vec{r}(t)$ en posición, velocidad y aceleración en el instante t_0 , y el radio de esa trayectoria circular será $R = v^2(t_0)/a_n(t_0)$. En términos analíticos, dicha trayectoria circular tiene un contacto de segundo orden con la trayectoria dada en el punto $\vec{r}(t_0)$, es lo que se llama la *circunferencia oscultriz* de la curva dada en el punto $\vec{r}(t_0)$. En resumen:

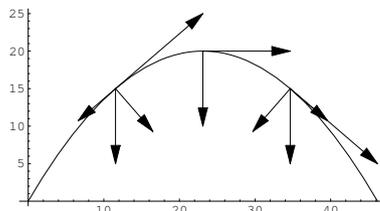
Si $\vec{r}(t)$ es una trayectoria con velocidad y aceleración normal no nulas en el instante t , entonces el valor $R(t) = v^2(t)/a_n(t)$ es el radio de la circunferencia oscultriz a la trayectoria en el punto $\vec{r}(t)$, es decir, de la circunferencia que más se parece a la trayectoria en un entorno de $\vec{r}(t)$.

En particular, cuanto mayor sea la aceleración normal, menor será el radio de la circunferencia oscultriz y más curvada estará la trayectoria en el punto considerado.

Ejemplo Continuamos nuestro análisis de la piedra que lanzábamos en un ángulo de 60° . Derivando su velocidad obtenemos su aceleración:

$$\vec{a} = (0, 0, -10) \text{ m/s}^2.$$

Vemos que es constante y está dirigida siempre hacia abajo.



La figura muestra los vectores \vec{a} , \vec{a}_t y \vec{a}_n para los tiempos $t = 1, 2, 3$. Vemos que, mientras la piedra sube, la aceleración tangencial es opuesta a la velocidad, por lo que la velocidad escalar va disminuyendo, mientras que, cuando baja, la aceleración tangencial tiene el mismo sentido que la velocidad, por lo que la velocidad escalar va aumentando. ■

Velocidad angular Vamos a introducir más elementos para describir la curvatura de la trayectoria de un móvil. Fijamos un observador situado en un punto O por el que no pase la trayectoria del móvil. Llamemos $A(t, u) \in [0, \pi]$ al ángulo que forman los vectores de posición $\vec{r}(t)$ y $\vec{r}(u)$ del móvil en dos instantes t y u (con el convenio de que es negativo cuando $u < t$ y positivo cuando $u > t$). Definimos la *velocidad angular (escalar)* del móvil en un instante t como

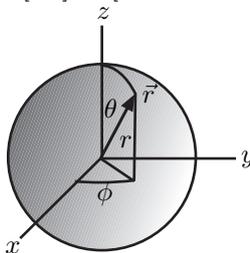
$$\omega(t) = \lim_{\Delta t \rightarrow 0} \frac{A(t, t + \Delta t)}{\Delta t}.$$

Observemos que la definición depende de la elección del origen O del sistema de referencia, pero no de la elección de los ejes coordenados. (En efecto, los vectores de posición dependen de la elección de O , pero el ángulo que forman dos de ellos es independiente de la base ortonormal respecto a la que se expresan.)

Vamos a obtener una expresión fácil de calcular para esta velocidad angular. Por la observación precedente podemos calcular $\omega(t_0)$ tomando unos ejes tales que $\vec{r}(t_0) = (0, 0, z(t_0))$. Ahora consideramos coordenadas esféricas, es decir,

$$\vec{r} = (r \cos \phi \sen \theta, r \sen \phi \sen \theta, r \cos \theta),$$

para $(r, \phi, \theta) \in]0, +\infty[\times]-\pi, \pi[\times]0, \pi[$.



Aquí θ es el ángulo que forma \vec{r} con el eje z . Notemos que podemos definir $\theta(t_0) = 0$ y $\rho(t_0) = z(t_0)$, pero el ángulo $\phi(t_0)$ no está definido.

Para tiempos $t > t_0$, es claro que $A(t_0, t) = \theta(t)$, luego, admitiendo que existe el límite que define a $\omega(t_0)$, éste tiene que ser $\omega(t_0) = \theta'_+(t_0)$. Supondremos además que $\theta'(t)$ es continua en $[0, t_1[$, para cierto tiempo t_1 . Observemos que si θ es idénticamente nula en un intervalo $[t_0, t_1[$, entonces $\omega = 0$ y no hay nada que calcular. En caso contrario vamos a considerar únicamente tiempos $t > t_0$ donde $\theta(t) \neq 0$. Para cualquiera de ellos tenemos que

$$\vec{v} = (r \cos \phi \cos \theta \theta' + \dots, r \sin \phi \cos \theta \theta' + \dots, r' \cos \theta - r \sin \theta \theta'),$$

donde los puntos suspensivos representan términos con un factor $\sin \theta$. Ahora calculamos

$$\vec{r} \times \vec{v} = (-r^2 \sin \phi \cos^2 \theta \theta' + \dots, r^2 \cos \phi \cos^2 \theta \theta', \dots),$$

de donde

$$\|\vec{r} \times \vec{v}\| = r^2 \cos^2 \theta \theta' + \dots$$

Finalmente hacemos tender $t \rightarrow t_0$ en ambos miembros. Notemos que esto tiene sentido, puesto que podemos acercarnos a t_0 por puntos donde $\theta(t) \neq 0$, que son los únicos que estamos considerando. Los puntos suspensivos representan productos de potencias de r (que tiende a $z(t_0)$) por potencias de $\sin \theta(t)$ (que tiende a 0) por otras funciones trigonométricas (acotadas), luego el límite de los puntos suspensivos es cero. Concluimos que en t_0 se cumple la igualdad

$$\|\vec{r} \times \vec{v}\| = r^2 \omega.$$

Esta igualdad vale igualmente en el caso en que θ es idénticamente nula en un entorno de t_0 , que hemos tratado aparte, pues en tal caso \vec{v} tiene la misma dirección que \vec{r} y el término izquierdo se anula, al igual que ω . Ahora bien, es claro que las expresiones de ambos miembros son independientes de la elección de los ejes coordenados, luego la igualdad no depende de la elección del sistema de referencia (para un origen O fijo) que hemos usado para demostrarla, luego también es válida para todo tiempo t . Concluimos, pues que

$$\omega = \frac{1}{r^2} \|\vec{r} \times \vec{v}\|,$$

y esto nos lleva a definir la *velocidad angular* $\vec{\omega}$ como

$$\vec{\omega} = \frac{1}{r^2} \vec{r} \times \vec{v},$$

que es un vector de módulo ω perpendicular tanto a \vec{v} como a \vec{r} .

Por ejemplo, si el movimiento se produce en un plano y consideramos coordenadas polares

$$\vec{r} = (\rho \cos \theta, \rho \sin \theta, 0), \quad \vec{v} = (\rho' \cos \theta - \rho \sin \theta \theta', \rho' \sin \theta + \rho \cos \theta \theta', 0),$$

se comprueba fácilmente que $\omega = \theta'$, es decir, la velocidad angular (escalar) representa el ángulo que gira el móvil por unidad de tiempo respecto del origen de coordenadas, y $\vec{\omega}$ señala el eje de giro (es perpendicular al plano de giro).

En el caso general podemos considerar la recta que pasa por O con dirección $\vec{\omega}$ como el eje de giro instantáneo del móvil.

Cambio de sistema de referencia Recordamos aquí las fórmulas que hemos obtenido en la página 229 y siguientes que relacionan la velocidad y la aceleración de un móvil en dos sistemas de referencia. Modificando ligeramente la notación, si tenemos dos observadores O y O' , existe un vector $\vec{\omega}$ determinado por ambos de modo que, para cualquier partícula puntual cuya posición respecto a cada uno de ellos venga dada por los vectores \vec{r} y \vec{r}' , respectivamente, las velocidades \vec{v} y \vec{v}' cumplen

$$\vec{v}' = \vec{v} - \vec{v}_{O'} - \vec{\omega} \times \vec{r}', \quad (\text{C.1})$$

donde $\vec{v}_{O'}$ es la velocidad de O' respecto de O , y las aceleraciones cumplen

$$\vec{a}' = \vec{a} - \vec{a}_{O'} - \vec{\omega} \times (\vec{\omega} \times \vec{r}') - 2\vec{\omega} \times \vec{v}' - \dot{\vec{\omega}} \times \vec{r}', \quad (\text{C.2})$$

donde $\dot{\vec{\omega}}$ es la derivada de $\vec{\omega}$.

C.3 Fuerzas

En la sección anterior hemos visto cómo el movimiento de un cuerpo puede ser descrito en términos de su velocidad y su aceleración, pero no hemos dicho nada sobre qué hace que un cuerpo dado se mueva con una determinada velocidad o una determinada aceleración en cada instante. No hemos dicho nada, en suma, sobre las causas que producen y determinan el movimiento de un cuerpo. Al abordar esta cuestión pasamos de la cinemática a la dinámica.

El principio de inercia El primer hecho fundamental que podemos tener en consideración a la hora de determinar el estado de movimiento que cabe esperar en un cuerpo es el *principio de inercia* de Galileo, o *primera ley de Newton*:

Si un cuerpo está libre de toda influencia externa, su velocidad permanece constante, es decir, o permanece en reposo, o se mueve en línea recta con velocidad escalar constante.

Galileo invirtió muchos esfuerzos en tratar de convencer a sus coetáneos de que un cuerpo en movimiento no termina parándose al cabo de un tiempo si no es por la acción de otros cuerpos. El hecho de que al darle una patada a un balón éste termina parándose no contradice el principio de inercia, pues el balón es frenado por el rozamiento con el suelo y con el aire. Sin tal rozamiento, nunca dejaría de moverse, y no necesitaría ninguna clase de “motor” que lo mantuviera en movimiento. La finalidad del motor de un coche no es mover el coche, sino contrarrestar el efecto del rozamiento con la carretera y con el aire, que tiende a detenerlo.

Sin embargo, hemos de advertir inmediatamente que el principio de inercia no puede ser cierto en todos los sistemas de referencia: Si una piedra está parada y alguien la empuja, la piedra ha alterado su velocidad debido a una causa externa (el empujón), lo cual no contradice al principio de inercia, pero si la piedra está en el asiento de un coche y éste frena bruscamente, la piedra puede experimentar exactamente el mismo efecto que si alguien la hubiera empujado, con la diferencia de que en este caso nadie lo ha hecho. Respecto a un observador para el que la carretera esté en reposo no hay ningún misterio: el coche (y la piedra con él) se están moviendo a una velocidad constante, de repente el coche frena (por efecto de los frenos, que detienen las ruedas y aumentan el rozamiento con el suelo) y la piedra sigue moviéndose con la misma velocidad. Esto hace que, para un observador situado en el coche, “parece” que la piedra es empujada hacia delante, aunque nadie la empuje de hecho.

Diremos que un sistema de referencia es *inercial* si en él se cumple el principio de inercia, es decir, si, respecto a él, cualquier aceleración que experimente un objeto puede explicarse como el efecto de algún fenómeno físico que ha actuado sobre él.

En estos términos, observador que viaja en un coche que frena no es un observador inercial, porque verá cómo una piedra situada en un asiento experimenta una aceleración que no podrá imputar a nada que haya actuado sobre ella (y esto sólo puede explicarse en términos del frenado del coche, lo cual obliga a considerar otro sistema de referencia, ya que, respecto al propio coche, éste se encuentra siempre en reposo y no tiene sentido decir que frena).

Si O es un observador inercial, cualquier otro observador O' que se mueva a velocidad constante respecto de O y cuyos ejes coordenados mantengan siempre la misma dirección es también un sistema de referencia inercial. Esto es consecuencia de la fórmula (C.2), que en este caso se reduce a $\vec{a}' = \vec{a}$ (pues $\vec{\omega} = \vec{0}$), de modo que un objeto experimenta una aceleración respecto de O si y sólo si la experimenta respecto de O' , si y sólo si alguna causa física está actuando sobre él. No es difícil convencerse de que el recíproco también es cierto, de modo que:

Si O es un observador inercial, los demás observadores inerciales son exactamente aquellos observadores O' que se mueven respecto de O con velocidad constante y cuyos ejes coordenados mantienen siempre la misma dirección.

Además, hemos visto que la aceleración que experimenta un móvil es la misma para todos los observadores inerciales, en el sentido de que la aceleración que observa uno de ellos se transforma en la que observa el otro sin más que aplicar el correspondiente cambio de coordenadas algebraico

$$(a_1, a_2, a_3) = \sum_i a'_i \vec{u}_i,$$

pero sin tener que añadir otros términos como los que aparecen en la fórmula general (C.2).

Fuerza En principio, una *fuerza* es cualquier causa física capaz de alterar la velocidad de un móvil, es decir, de comunicarle una aceleración. Matemáticamente, una fuerza se cuantifica mediante un vector. En principio podría pensarse que, desde un punto de vista matemático, “fuerza” es sinónimo de “aceleración”, pues hay una fuerza allí donde hay una aceleración, pero esto es falso por dos motivos:

En primer lugar, dos o más fuerzas pueden compensarse unas a otras. Por ejemplo, no es cierto que un libro que permanece en reposo sobre mi mesa no esté afectado por fuerza alguna. Al contrario, sobre él actúan dos fuerzas que se compensan mutuamente: la fuerza gravitatoria causada por la Tierra, que lo haría caer si fuera la única presente, y la fuerza eléctrica con la que los electrones de la superficie de mi mesa repelen a los electrones de la superficie del libro. Tiene sentido decir que hay dos fuerzas que se cancelan mutuamente, pero no que hay dos aceleraciones que se cancelan. Dado que el libro está en reposo, su aceleración es cero y ya está.

En segundo lugar, es razonable considerar que una misma causa produce siempre la misma fuerza. Por ejemplo, si le damos una patada a un balón en reposo, de modo que éste pasa a moverse por el suelo hasta que se detiene, nuestro pie ha ejercido una fuerza sobre el balón durante el breve periodo de tiempo de contacto. Esta fuerza le ha comunicado una aceleración al balón, y la velocidad que tenía en el momento en que se ha separado del pie ha empezado a disminuir paulatinamente por efecto de la fuerza de rozamiento (de naturaleza eléctrica, como la de la patada) que ejerce sobre él el aire y el suelo. Sin embargo, si le damos una patada exactamente igual a una piedra, su velocidad cuando se separe de nuestro pie —y, por consiguiente— la aceleración que ha experimentado durante el contacto, será significativamente menor que la que había adquirido el balón. Es razonable considerar que ambas patadas han aplicado la misma fuerza, tanto al balón como a la piedra, pero las aceleraciones producidas han sido distintas.

La segunda ley de Newton afirma que las diferencias de aceleración que una misma fuerza provoca a objetos diferentes dependen únicamente de una magnitud de cada objeto, que es la que llamamos su “masa”. La masa, como la longitud y la duración, es una magnitud fundamental, y su unidad de medida en el sistema internacional de unidades es el *kilogramo*, que se define simplemente como la masa de un patrón de platino e iridio conservado en la *Oficina Internacional de Pesos y Medidas* en Sèvres (cerca de París).³

Ahora podemos enunciar la *segunda ley de Newton*:

La suma de todas las fuerzas que actúan sobre un móvil en un instante dado es igual al producto de su masa por su aceleración en dicho instante: $\vec{F} = m\vec{a}$.

³La idea es que un kilogramo es la masa de un decímetro cúbico de agua en determinadas condiciones de presión, temperatura, etc., pero, por razones prácticas, oficialmente se toma como definición la masa del patrón. Es la única unidad del Sistema Internacional que no se define mediante un fenómeno físico susceptible de ser reproducido sin necesidad de viajar a Francia.

Observemos que el contenido físico de la segunda ley no está en la fórmula $\vec{F} = m\vec{a}$, que puede considerarse como la definición matemática del concepto de fuerza, sino en la afirmación de que las interacciones físicas entre distintos cuerpos se pueden cuantificar mediante vectores “sumables” y de modo que el efecto (la aceleración) que produce una misma causa sólo depende de la fuerza que la describe a través del factor de proporcionalidad dado por la masa.

La fuerza se mide en *Newtons*, donde un Newton no es más que un $\text{kg}\cdot\text{m}/\text{s}^2$ o, en otros términos, que una fuerza tiene una intensidad de 1 N si produce una aceleración de $1\text{m}/\text{s}^2$ sobre un cuerpo de 1kg de masa.

La segunda ley de Newton incluye a la primera como un caso particular: si un objeto está libre de toda influencia externa, sobre él no actúa ninguna fuerza, luego su aceleración será nula y su velocidad permanecerá constante. En particular, la segunda ley de Newton sólo puede ser cierta en los sistemas de referencia en los que sea válida la primera, es decir, en los sistemas de referencia inerciales.

Si dos observadores inerciales ven cómo un pie da una patada a un balón, sabemos que la aceleración final del balón será la misma para ambos, luego, si convenimos en que “la misma patada” ha de ser descrita matemáticamente por “la misma fuerza”, la segunda ley de Newton nos da que la masa del balón ha de ser la misma para ambos observadores. En realidad, la teoría de la relatividad afirma que la masa de un cuerpo varía de un sistema inercial a otro, pero esto sólo se aprecia en móviles que viajan a velocidades comparables a la velocidad de la luz. La mecánica clásica supone, por el contrario, que la masa de un objeto es independiente del observador (en principio, inercial) que la mide. El carácter absoluto (es decir, no relativo al sistema de referencia) de la masa y la aceleración implican inmediatamente que las fuerzas también son absolutas.

Ejemplo: La gravedad terrestre Un cuerpo situado sobre la superficie de la Tierra experimenta la fuerza de gravedad terrestre, dirigida hacia el centro de la Tierra y cuya intensidad es aproximadamente $g = 9.8\text{N}/\text{kg}$. Este valor puede oscilar levemente en función de la latitud y la altitud.

Teniendo en cuenta que dos vectores que apunten hacia el centro de la Tierra son prácticamente paralelos si los consideramos aplicados a puntos no excesivamente lejanos, resulta que la fuerza gravitatoria que actúa en la superficie terrestre sobre un cuerpo de masa m puede considerarse de la forma

$$\vec{P} = (0, 0, -mg),$$

respecto a cualquier sistema de referencia en el que el eje Z sea vertical. Dicha fuerza es lo que comúnmente se llama el *peso* del objeto.⁴

⁴Es frecuente expresar el peso de un objeto en kg, con lo cual no estamos dando su peso, sino su masa. Restringidos a la superficie terrestre ambos conceptos están relacionados por la proporción $P = mg$, pero en otros contextos no son equivalentes. Por ejemplo, un objeto de 1 kg de masa sigue teniendo la misma masa sobre la superficie de la Luna, pero allí pesará mucho menos que sobre la superficie de la Tierra.

La segunda ley de Newton nos asegura entonces que un cuerpo sometido únicamente a la fuerza de la gravedad experimenta una aceleración constante \vec{a} determinada por la ecuación $\vec{P} = m\vec{a}$, con lo que.⁵ $\vec{a} = (0, 0, -g) \text{ m/s}^2$.

Esto significa que un cuerpo en caída libre experimenta una aceleración hacia abajo de 9.8 m/s^2 independientemente de cuál sea su masa. En particular, si dejamos caer dos objetos desde una misma altura, en la medida en que podamos despreciar las fuerzas de rozamiento, ambos llegarán al suelo al mismo tiempo, independientemente de cuáles sean sus masas. ■

Para completar las leyes fundamentales de la dinámica hemos de considerar la *tercera ley de Newton*:

Si un cuerpo ejerce una fuerza sobre otro, entonces éste ejerce sobre el primero una fuerza igual en módulo y dirección, pero de sentido contrario.

Así, por ejemplo, si estoy sentado en una silla con ruedas y aplico una fuerza a una pared, yo me veré impulsado hacia atrás, porque la pared está ejerciendo sobre mí la misma fuerza que yo ejerzo sobre ella, pero en sentido contrario. Si no estoy sobre una silla con ruedas, tal vez no me mueva, pero ello se deberá a que la fuerza que la pared ejerce sobre mí se cancela con la fuerza de rozamiento del suelo.

Ejemplo Si dejamos caer un objeto de 1 kg de masa desde una altura de 1 m, y analizamos el movimiento desde un sistema de referencia inercial, teóricamente no es cierto que el objeto corra hacia la Tierra mientras ésta permanece inmóvil. Por el contrario, la Tierra está siendo atraída por el objeto con la misma fuerza, luego también experimenta una aceleración y ambos —Tierra y objeto— corren mutuamente hasta encontrarse. Eso sí, las aceleraciones que experimentan son inversamente proporcionales a sus masas respectivas, luego los espacios recorridos son muy desiguales. Concretamente, la fuerza con la que el objeto es atraída es de 9.8 N (la masa del objeto por la aceleración que le imprime la Tierra), luego la Tierra también es atraída con una fuerza de 9.8 N por el objeto. Como la masa de la Tierra es de $M_T = 5.9736 \cdot 10^{24} \text{ kg}$ y tiene que cumplirse que $9.8 = M_T a_T$, la aceleración que experimenta es $a_T = 1.6405 \cdot 10^{-24} \text{ m/s}^2$.

Si tomamos un sistema de referencia en el que la superficie de la Tierra está en $x = 0$ y el cuerpo en $x = 1$ cuando $t = 0$, ambos se encontrarán cuando

$$\frac{1}{2} a_T t^2 = 1 - \frac{1}{2} 9.8 t^2.$$

Esto nos da que $t = \sqrt{2/(9.8 + a_T)}$, luego el espacio recorrido por la Tierra es

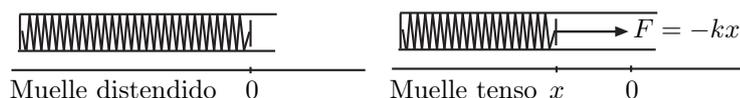
$$r = \frac{a_T}{9.8 + a_T} = 1.674 \cdot 10^{-25} \text{ m}.$$

Concluimos que, en la práctica, podemos afirmar que la Tierra permanece en reposo mientras el objeto cae sobre ella. ■

⁵En el ejemplo de la página 442 hemos redondeado $g \approx 10 \text{ m/s}^2$.

Ejemplo: Muelles y dinamómetros Los muelles proporcionan una técnica para medir fuerzas gracias a que verifican la llamada ley de Hooke. Un muelle tiene la propiedad de que se puede estirar o contraer bajo la acción de una fuerza, pero es capaz de regresar a su posición inicial una vez ha cesado la fuerza causante de la deformación. Por simplificar la situación, supongamos que metemos un muelle en un tubo, de modo que no pueda realizar ninguna clase de desplazamiento más que en sentido longitudinal, y fijamos uno de sus extremos dejando libre el otro. Tomamos un sistema de referencia de modo que el muelle esté sobre el eje x y su extremo libre ocupe la posición $x = 0$ cuando no actúa ninguna fuerza sobre él.

La *ley de Hooke* afirma que el muelle ejerce sobre su extremo libre una fuerza proporcional a la distancia en que ha sido desplazado de la posición que ocuparía en ausencia de fuerzas. Como esta fuerza está dirigida en sentido contrario al desplazamiento, respecto al sistema de referencia que hemos fijado se expresará como $F = -kx$, donde k es una constante que depende del muelle. Así, por ejemplo, para mantener el muelle en la posición de la figura de la derecha, es necesario mantener aplicada en su extremo una fuerza $F' = kx$ (hacia la izquierda) que compense la fuerza hacia la derecha $-kx$ que ejerce el muelle.



La constante k puede determinarse sin más que aplicar al muelle una fuerza de intensidad conocida (por ejemplo, colgando de él un cuerpo de masa conocida). Una vez determinada, el muelle se convierte en un *dinamómetro*, es decir, en un aparato de medida de fuerzas. Para medir una fuerza se la hace actuar sobre el extremo del muelle, se mide el desplazamiento que produce y se multiplica por la constante k . En particular, un muelle puede usarse como báscula, ya sea poniéndolo debajo de un plato para que cualquier objeto depositado en él oprima el muelle con la fuerza de su peso, ya sea colgando de él un plato para que el peso lo estire.

Es claro que la ley de Hooke sólo es válida para desplazamientos que no sobrepasen un cierto umbral, a partir del cual el muelle sufrirá una deformación permanente. Para desplazamientos aún mayores el muelle se romperá. ■

Momento El *momento lineal* o *cantidad de movimiento* de una partícula puntual de masa m se define como el vector

$$\vec{p} = m\vec{v},$$

donde \vec{v} es su velocidad.

A la hora de estudiar el movimiento de una única partícula este concepto no aporta nada que no esté reflejado —más precisamente incluso— al considerar por separado la masa m y la velocidad \vec{v} . Sin embargo, su importancia se pone de manifiesto cuando consideramos simultáneamente varias partículas. Así, al

sumar el momento lineal de todas las partículas de un sistema, estaremos diciendo que una partícula que se mueva con la misma velocidad que otra pero con el doble de masa ha de contar como dos partículas. La idea es que, por ejemplo, a la hora de evaluar el efecto de una bala de cañón que golpea a una pared, no sólo influye la velocidad de la bala, sino también su masa, pues una bala de 2 kg producirá el mismo efecto que dos balas de 1 kg que llegaran simultáneamente a la misma velocidad.

En términos del momento lineal, la segunda ley de Newton puede escribirse en la forma

$$\vec{F} = \frac{d\vec{p}}{dt},$$

donde \vec{F} es la fuerza total que actúa sobre la partícula en un instante dado.

Definimos el *momento angular* de una partícula como el vector

$$\vec{L} = \vec{r} \times \vec{p} = mr^2 \vec{\omega}, \quad (\text{C.3})$$

donde $\vec{\omega}$ es su velocidad angular, que depende del origen del sistema de referencia considerado. Si además definimos el *momento* de una fuerza \vec{F} que actúa sobre la partícula como

$$\vec{M} = \vec{r} \times \vec{F},$$

se cumple la versión angular de la segunda ley de Newton, pues

$$\frac{d\vec{L}}{dt} = \vec{v} \times \vec{p} + \vec{r} \times \vec{F} = \vec{r} \times \vec{F} = \vec{M},$$

ya que \vec{v} y \vec{p} tienen el mismo sentido, luego su producto vectorial es nulo.

C.4 Trabajo y energía

Trabajo Sea $\vec{r}(t)$ la posición de un móvil respecto de un observador O y sea $\vec{F}(t)$ una fuerza que está actuando sobre él (no necesariamente la fuerza total). Se define el *trabajo* que la fuerza \vec{F} realiza sobre el móvil entre un tiempo t_0 y un tiempo t_1 como

$$W = \int_{\vec{r}(t_0)}^{\vec{r}(t_1)} \vec{F} \cdot d\vec{r}.$$

Se trata de una integral sobre una curva (la trayectoria del móvil) de una forma diferencial. Si el lector no está familiarizado con este tipo de integrales puede considerar la siguiente definición alternativa:

$$W = \int_{t_0}^{t_1} \vec{F} \cdot \frac{d\vec{r}}{dt} dt = \int_{t_0}^{t_1} (F_x v_x + F_y v_y + F_z v_z) dt.$$

Algunas consecuencias elementales de la definición son las siguientes:

- El trabajo es lineal en la fuerza: el trabajo realizado por una suma de fuerzas es igual a la suma de los trabajos realizados por cada una de ellas.
- Si una fuerza actúa sobre un objeto y éste no se mueve (porque la fuerza total que actúa sobre él es nula) entonces $\vec{v} = 0$ y el trabajo realizado por la fuerza es nulo.
- Si descomponemos la fuerza como $\vec{F} = \vec{F}_t + \vec{F}_n$, donde \vec{F}_t es tangente a la trayectoria del móvil y \vec{F}_n es ortogonal a \vec{F}_t , entonces $\vec{F} \cdot \vec{v} = \vec{F}_t \cdot \vec{v}$, luego el trabajo realizado por \vec{F} es el mismo que el realizado por \vec{F}_t . En otras palabras: el trabajo sólo depende de la componente tangencial de \vec{F} .
- Si la fuerza es constante, entonces

$$W = \vec{F} \cdot \int_{t_0}^{t_1} \vec{v} dt = \vec{F} \cdot \Delta\vec{r},$$

donde $\Delta\vec{r} = \vec{r}(t_1) - \vec{r}(t_0)$.

Observemos que el trabajo es una magnitud escalar y no vectorial. Su unidad en el internacional de unidades es el *julio*, que no es sino un $N \cdot m$. Más explícitamente: *un julio es el trabajo que realiza una fuerza constante de un Newton cuando se aplica a un cuerpo en un desplazamiento rectilíneo de un metro en la misma dirección del desplazamiento.*

Energía cinética Si \vec{F} es la fuerza total que actúa sobre un móvil, entonces la segunda ley de Newton afirma que $\vec{F} = m\vec{a}$, por lo que el trabajo es

$$\begin{aligned} W &= \int_{t_0}^{t_1} m\vec{a} \cdot \vec{v} dt = \sum_{i=1}^3 \int_{t_0}^{t_1} m \frac{dv_i}{dt} v_i dt = \sum_{i=1}^3 \int_{t_0}^{t_1} \frac{d}{dt} \left(\frac{1}{2} m v_i^2 \right) dt \\ &= \int_{t_0}^{t_1} \frac{d}{dt} \left(\frac{1}{2} m v^2 \right) dt = \Delta E_c, \end{aligned}$$

donde hemos llamado E_c a la *energía cinética* del móvil, definida como

$$E_c = \frac{1}{2} m v^2,$$

y $\Delta E_c = E_c(t_1) - E_c(t_0)$. Con esto hemos demostrado el teorema siguiente:

El trabajo que realiza la fuerza total que actúa sobre un móvil durante un intervalo de tiempo es igual al incremento de energía cinética que experimenta el móvil en dicho intervalo.

La energía cinética se mide también en julios, lo cual es coherente, pues

$$1 \text{ J} = 1 \text{ N} \cdot \text{m} = 1 \text{ kg} \cdot \text{m/s}^2 \cdot \text{m} = 1 \text{ kg} \cdot (\text{m/s})^2.$$

En particular vemos que, si sobre un cuerpo no actúan fuerzas en la dirección de su movimiento, su energía cinética permanece constante, pues las fuerzas que actúan sobre él no están realizando ningún trabajo.

Energía potencial En muchas circunstancias, la fuerza que actúa sobre un objeto depende únicamente de la posición de éste en el espacio a través de una relación de la forma $\vec{F} = -\nabla V$, donde V es una función escalar definida sobre los puntos del espacio.

Por ejemplo, sobre la superficie terrestre, la fuerza gravitatoria que actúa sobre un cuerpo de masa m es

$$\vec{F} = (0, 0, -mg) = -\nabla V, \quad \text{donde} \quad V = mgz.$$

Otro ejemplo es la fuerza que un muelle ejerce sobre una partícula de masa m , que es

$$\vec{F} = -kx = -\nabla V, \quad \text{donde} \quad V = \frac{1}{2}kx^2.$$

Las fuerzas que pueden expresarse de este modo se llaman *fuerzas conservativas* y la función V correspondiente se llama *función potencial* asociada a la fuerza. Observemos que las unidades de V son N·m, es decir, julios. Para comprender el significado de este hecho observamos que el trabajo que realiza una fuerza conservativa sobre un cuerpo que se mueve por una trayectoria $\vec{r}(t)$ es

$$\begin{aligned} W &= \int_{t_0}^{t_1} \vec{F}(\vec{r}(t)) \cdot \vec{v}(t) dt = - \int_{t_0}^{t_1} \nabla V(\vec{r}(t)) \cdot \vec{v}(t) dt \\ &= - \int_{t_0}^{t_1} \frac{dV(\vec{r}(t))}{dt} dt = -(V(\vec{r}(t_1)) - V(\vec{r}(t_0))) = -\Delta V, \end{aligned}$$

donde en la última expresión V es $V(\vec{r}(t))$. En particular, vemos que este trabajo sólo depende de la posición inicial y final del objeto, pero no de la trayectoria que ha seguido.

Esto nos lleva a definir la *energía potencial* de un objeto sometido a una fuerza conservativa \vec{F} como el valor de la función potencial de \vec{F} en el punto del espacio en el que se encuentra. Si llamamos \vec{G} a la resultante de las demás fuerzas que actúan sobre el cuerpo y W_{ext} es el trabajo realizado por G , el trabajo total realizado sobre el cuerpo es $W_{\text{ext}} - \Delta V = \Delta E_c$ o, equivalentemente, si llamamos *energía mecánica* del cuerpo a la suma $E_m = E_c + V$, tenemos la relación $W_{\text{ext}} = \Delta E_m$, es decir:

La variación de la energía mecánica de un cuerpo es igual al trabajo realizado sobre él por las fuerzas exteriores no incluidas en el cálculo de la energía mecánica.

Observemos que si sobre un cuerpo actúan varias fuerzas conservativas, la suma de todas ellas es una fuerza conservativa, cuya función potencial es la suma de las funciones potenciales de cada fuerza, luego en el cálculo de la energía mecánica de un cuerpo podemos incluir todas las fuerzas conservativas que actúan sobre él. En particular, si sobre un cuerpo sólo actúan fuerzas conservativas, concluimos que su energía mecánica no varía con el tiempo.

Por ejemplo, si tomo un libro de mi mesa y lo elevo hasta una estantería, he aumentado la energía potencial del libro en mgh julios, donde m es la masa del libro, g la intensidad del campo gravitatorio y h la distancia que he elevado el libro. Eso se traduce en que si dejo caer el libro desde esa altura, su velocidad al llegar a la mesa será la necesaria para que $\frac{1}{2}mv^2 = mgh$. En efecto, cuando haya perdido al caer la energía potencial que ha ganado, ésta se habrá transformado en energía cinética, de modo que la energía cinética final coincidirá con la energía potencial inicial.

Energía interna Un ejemplo de fuerza no conservativa es la fuerza de rozamiento. Si pegamos una patada a un balón en reposo, durante el breve tiempo que nuestro pie está en contacto con el balón estamos realizando un trabajo que le comunica a éste una cierta energía cinética. Si el balón se mueve a ras del suelo, al cabo de un tiempo se detendrá, lo que significa que habrá perdido toda la energía cinética que le hemos dado sin haber ganado ninguna clase de energía potencial.

Sin embargo, no es correcto decir que la energía cinética “ha desaparecido”. Sólo ha desaparecido en el mismo sentido en que el humo de una hoguera que asciende por el aire termina disipándose, lo cual significa que se ha expandido tanto por la atmósfera que se vuelve imperceptible, pero sigue ahí. En el caso de la energía cinética del balón, lo que sucede es que las partículas del balón han estado en contacto con las partículas del suelo y les han traspasado su energía cinética. El movimiento del balón se ha transformado en el movimiento microscópico de una cantidad inmensa de moléculas del suelo, cada una de las cuales se ha quedado con una pequeña porción de su energía cinética.

En general, las moléculas de cualquier objeto material realizan pequeños movimientos vibratorios imperceptibles, pero que acumulan una cierta energía cinética, que es lo que comúnmente se conoce como *calor*. Calentar un cuerpo equivale a aumentar la energía cinética de sus partículas constituyentes. La unidad del calor en el sistema internacional de unidades no es sino el julio. Una unidad más familiar es la *caloría*, cuya definición actual es simplemente $1 \text{ cal} = 4.1868 \text{ J}$.

No hay que confundir la cantidad de calor que posee un cuerpo con su *temperatura*, la cual es un promedio de la energía que poseen sus partículas. Por ejemplo, dos litros de agua a 50°C tienen la misma temperatura que un litro de agua a 50°C , pero la cantidad de calor que hay en dos litros de agua es el doble de la que hay en un litro. La unidad de temperatura en el sistema internacional no es el grado centígrado, sino el *Kelvin*, de modo que $T^\circ\text{C} = T + 273.15 \text{ K}$. En general, la relación entre calor y temperatura es

$$Q = mcT,$$

donde Q es la cantidad de calor de un cuerpo, T es su temperatura y c es el *calor específico*, que depende de la composición química del cuerpo, y que no es sino la cantidad de calor que alberga cada kg de la sustancia por cada K de temperatura.

Más en general, el calor y otras formas de energía potencial constituyen la llamada *energía interna* de un cuerpo. Principalmente, la energía potencial interna es de naturaleza eléctrica, pero su comportamiento depende de la estructura atómica de la materia, y algunas de sus facetas sólo pueden ser descritas en términos de la mecánica cuántica. El estudio de la energía interna de los objetos constituye la rama de la física conocida como *termodinámica*.

Conservación de la energía Las varias formas de energía que hemos descrito son tan sólo casos particulares de un concepto de energía mucho más general:

La *energía* de un objeto o sistema físico es su capacidad para realizar un trabajo. Cualquier cualidad o situación de un objeto físico que sea susceptible de ser transformada en trabajo es una forma de energía.

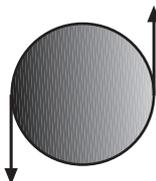
Naturalmente, esta definición no es operativa desde un punto de vista matemático, y lo que hemos hecho en esta sección es cuantificar la energía que hay que atribuir a un cuerpo en función de su velocidad, de las fuerzas que actúan en su entorno, etc. para determinar qué cantidad de trabajo está en condiciones de realizar. No obstante, hay muchas otras formas de energía.

La noción general de energía es importante porque permite enunciar una de las leyes fundamentales de la física, a saber, el *principio de conservación de la energía*, según el cual, cuantificando adecuadamente la energía de todo proceso físico, se cumple que la energía total del universo permanece siempre constante, aunque pueda transformarse de una forma a otra distinta.

Más localmente, sin apelar al universo entero, lo que afirma el principio de conservación de la energía es que la variación de la energía total de un sistema físico es igual al trabajo que realizan sobre él las fuerzas exteriores al mismo. En la práctica ahora podemos considerar un concepto más general de “trabajo”, entendido como cualquier transferencia de energía, aunque esta transferencia no se vea reflejada en la aplicación de una fuerza que haga recorrer un espacio a un objeto dado.

C.5 Dinámica de un sistema de partículas

En las secciones anteriores hemos estudiado la dinámica de las partículas puntuales, pero no toda situación dinámica puede conceptualizarse identificando cada objeto involucrado con una partícula puntual. Por ejemplo, es claro que si a una partícula puntual le aplicamos dos fuerzas opuestas de la misma intensidad la fuerza resultante será cero y su movimiento no se verá alterado, pero esto es falso si aplicamos tales fuerzas a una esfera como indica la figura:



Ciertamente, tales fuerzas no trasladarán la esfera, pero la harán girar. La clave está en que las dos fuerzas se aplican sobre puntos distintos de la esfera, cosa que no puede expresarse si la concebimos como un punto. Por ello, en esta sección vamos a estudiar el comportamiento global de sistemas formados por un número arbitrario de partículas puntuales.

Consideramos n partículas puntuales situadas en posiciones \vec{r}_i con masas m_i . Aunque nuestro objetivo a medio plazo es estudiar sólidos formados por un gran número de partículas, los resultados generales de esta sección se pueden aplicar también a sistemas formados por unas pocas partículas sin ninguna conexión entre sí más que nuestra intención de estudiarlas conjuntamente.

Para $i \neq j$, llamaremos \vec{F}_{ij} a la fuerza total que —de un modo u otro— la partícula i -ésima ejerce sobre la j -ésima. Nuestra intención es estudiar el comportamiento global del sistema sin entrar a analizar la naturaleza de estas fuerzas, a las que llamaremos *fuerzas internas*. Si las partículas representan átomos o moléculas, entonces el estudio de las fuerzas internas que se dan entre ellas pertenece al dominio de la química y de la mecánica cuántica. En cualquier caso, las fuerzas internas han de cumplir la tercera ley de Newton, que nos garantiza que $\vec{F}_{ij} = -\vec{F}_{ji}$. Por otra parte, llamaremos $\vec{F}_i = \vec{F}_{ii}$ a la resultante de todas las fuerzas exteriores ejercidas sobre la partícula i -ésima (es decir, sin contar las fuerzas \vec{F}_{ji}), y las llamaremos *fuerzas externas* del sistema. De este modo, la fuerza total que actúa sobre la partícula i -ésima es

$$\sum_{j=1}^n \vec{F}_{ji}.$$

Diremos que un sistema de partículas es un *sólido rígido* si sus fuerzas internas \vec{F}_{ij} son en cada momento las necesarias para que la distancia $\|\vec{r}_i - \vec{r}_j\|$ entre dos cualesquiera de sus partículas permanezca constante. En particular, \vec{F}_{ij} tendrá siempre la dirección de la recta que une \vec{r}_i con \vec{r}_j .

Cinemática Vamos a ver que un sistema de partículas se puede tratar hasta cierto punto como si fuera una única partícula. Para ello definimos la *masa total* del sistema como la suma de las masas de sus componentes:

$$M = \sum_{i=1}^n m_i,$$

y el *centro de masas* del sistema como la posición media de sus partículas ponderadas por sus masas, es decir:

$$\vec{R} = \frac{\sum_{i=1}^n m_i \vec{r}_i}{M}.$$

Esto significa que si tenemos dos partículas, pero una tiene el doble de masa que la otra, entonces la primera cuenta como si fueran dos partículas muy próximas, y el centro de masas del sistema estará al doble de distancia de la segunda

que de la primera. En particular, observemos que en el centro de masas de un sistema no tiene por qué haber ninguna partícula.

La velocidad del centro de masas es

$$\vec{V} = \frac{d\vec{R}}{dt} = \frac{\sum_{i=1}^n m_i \vec{v}_i}{M},$$

que puede considerarse como la “velocidad media” del sistema. Similarmente, la aceleración del centro de masas es

$$\vec{A} = \frac{d\vec{V}}{dt} = \frac{\sum_{i=1}^n m_i \vec{a}_i}{M}.$$

Momento Recordemos que el momento lineal o cantidad de movimiento de una partícula es el vector $\vec{p} = m\vec{v}$, donde m es su masa y \vec{v} su velocidad. El momento de un sistema de partículas se define como

$$\vec{P} = \sum_{i=1}^n m_i \vec{v}_i.$$

Es inmediato que tenemos la relación $\vec{P} = M\vec{V}$, que es un primer ejemplo de cómo un sistema de partículas verifica en promedio relaciones análogas a las de una única partícula. Ahora observamos que

$$\frac{d\vec{P}}{dt} = \sum_{i=1}^n m_i \vec{a}_i = \sum_{i=1}^n \sum_{j=1}^n \vec{F}_{ij} = \sum_{i=1}^n \vec{F}_i = \vec{F},$$

donde hemos cancelado los pares de fuerzas $\vec{F}_{ij} + \vec{F}_{ji} = 0$. En resumen, hemos probado la versión global de la segunda ley de Newton:

La derivada del momento de un sistema de partículas es igual a la suma de todas las fuerzas exteriores que se ejercen sobre el sistema.

Las fuerzas interiores no afectan en nada a la cantidad de movimiento del sistema. En particular, si sobre un sistema no actúan fuerzas externas, su cantidad de movimiento permanece constante.

Ejemplo Consideremos un sistema formado por dos partículas de masas M y m sometidas únicamente a la atracción gravitatoria dada por

$$\vec{F}_{12} = -\frac{GMm}{\|\vec{r}_2 - \vec{r}_1\|^3}(\vec{r}_2 - \vec{r}_1), \quad \vec{F}_{21} = -\vec{F}_{12}.$$

Como no hay fuerzas exteriores (respecto de un sistema inercial), el momento del sistema o, equivalentemente, la velocidad del centro de masas será constante. Esto implica que podemos tomar un sistema de referencia inercial con origen

en el centro de masas. Entonces, si las partículas tienen posiciones \vec{r}_1 y \vec{r}_2 , se cumple que $M\vec{r}_1 + m\vec{r}_2 = \vec{0}$, luego

$$\vec{r}_2 - \vec{r}_1 = \left(1 + \frac{m}{M}\right) \vec{r}_2.$$

Así pues, la fuerza que actúa sobre el objeto de masa m es

$$\vec{F} = -\frac{GM^*m}{r_2^3} \vec{r}_2, \quad \text{donde} \quad M^* = \frac{M}{(1 + m/M)^2}.$$

Ahora bien, el movimiento de un cuerpo afectado por una fuerza de este tipo está estudiado en la subsección 7.3.2. Allí se supone que el Sol (de masa M) permanece inmóvil, y produce esta fuerza sobre un planeta, y la conclusión es que el planeta se mueve en una órbita elíptica con el Sol en uno de sus focos (aunque en general la órbita puede ser parabólica o hiperbólica, como les sucede a algunos cometas). Ahora podemos precisar que la hipótesis de que el Sol permanece inmóvil no es exacta (o, mejor dicho, que si tomamos el Sol como origen de un sistema de referencia, éste no es inercial, porque el Sol está sometido a la atracción gravitatoria del planeta), pero sí que podemos tomar un sistema de referencia inercial con origen en el centro de masas del sistema Sol-planeta, y la conclusión es que la fuerza que actúa sobre el planeta tiene la misma forma salvo que la masa del Sol M debe reemplazarse por M^* .

Así pues, en un sistema binario de dos cuerpos sometidos únicamente a su atracción gravitatoria mutua, cada cuerpo sigue una trayectoria cónica con un foco en el centro de masas del sistema. En el caso en que uno de los dos objetos tiene masa mucho mayor que el otro, como ocurre con el Sol y cualquiera de los planetas del sistema solar, el centro de masas está dentro del Sol (pero no en su centro geométrico) y $M^* \approx M$, por lo que la diferencia es irrelevante, pero en el caso de dos objetos de masas similares, por ejemplo un sistema binario de estrellas, la corrección es relevante. ■

Momento angular El *momento angular* de un sistema de partículas se define como la suma de los momentos angulares de sus componentes:

$$\vec{L} = \sum_{i=1}^n \vec{r}_i \times \vec{p}_i.$$

La versión angular global de la segunda ley de Newton requiere una hipótesis adicional sobre el sistema. En efecto, si calculamos la derivada de \vec{L} obtenemos

$$\frac{d\vec{L}}{dt} = \sum_{i=1}^n \frac{d\vec{L}_i}{dt} = \sum_{i,j=1}^n \vec{r}_i \times \vec{F}_{ji} = \sum_{i=1}^n \vec{r}_i \times \vec{F}_i + \sum_{i<j} (\vec{r}_i - \vec{r}_j) \times \vec{F}_{ij},$$

donde hemos usado que

$$\vec{r}_i \times \vec{F}_{ji} + \vec{r}_j \times \vec{F}_{ij} = \vec{r}_i \times \vec{F}_{ji} - \vec{r}_j \vec{F}_{ji} = (\vec{r}_i - \vec{r}_j) \times \vec{F}_{ij}.$$

Vemos así que, en principio, la variación de \vec{L} depende de las fuerzas internas del sistema. No obstante, podemos prescindir de ellas si suponemos que son radiales, en el sentido siguiente:

Diremos que un sistema de partículas tiene *fuerzas internas radiales* si la fuerza \vec{F}_{ij} que la partícula i -ésima realiza sobre la j -ésima tiene la dirección del vector $\vec{r}_i - \vec{r}_j$.

Esto sucede, por ejemplo, en el caso de un sólido rígido. Bajo la hipótesis de que las fuerzas internas sean radiales llegamos a la relación

$$\frac{d\vec{L}}{dt} = \sum_{i=1}^n \vec{r}_i \times \vec{F}_i = \sum_{i=1}^n \vec{M}_i = \vec{M},$$

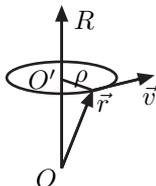
donde \vec{M} es la suma de los momentos que las fuerzas exteriores al sistema realizan sobre el sistema. Ésta es la versión angular de la segunda ley de Newton:

La variación del momento angular de un sistema de partículas con fuerzas internas radiales es igual a la suma de los momentos ejercidos por las fuerzas exteriores al sistema.

Momento de inercia de una partícula Consideremos una partícula puntual de masa m que describe una trayectoria circular de radio ρ alrededor de una recta R con velocidad escalar v constante. Que la velocidad escalar v sea constante equivale a que lo sea la velocidad angular $\omega = v/\rho$. Hemos visto que esta velocidad angular se corresponde con el concepto general de velocidad angular (escalar) respecto de cualquier sistema de referencia con origen en el centro O' de la circunferencia que describe la partícula. En tal caso, la velocidad $\vec{\omega}$ es perpendicular a la trayectoria y, como

$$\vec{L} = \vec{r} \times \vec{p} = m\rho^2\vec{\omega},$$

concluimos que el momento angular permanece constante. Desde un punto de vista dinámico, esto se debe a que la fuerza \vec{F} que mantiene a la partícula en su trayectoria apunta hacia O' , luego tiene la dirección del vector de posición, luego su momento es $\vec{M} = \vec{r} \times \vec{F} = \vec{0}$, y de aquí la invarianza del momento angular.



Ahora bien, esto ya no sería cierto si tomáramos como origen cualquier otro punto O de la recta R . En tal caso, el momento de \vec{F} ya no sería nulo y, por otra parte, es claro que \vec{L} sería en la figura un vector oblicuo, perpendicular al plano que forman \vec{r} y \vec{v} , que claramente variaría con la posición de la partícula.

Así pues, la elección de O' como origen del sistema de referencia es particularmente adecuada para estudiar el movimiento circular. Si sobre el móvil actúa una fuerza que, no obstante, lo mantiene igualmente en su trayectoria circular (de modo que ρ es constante), entonces la aceleración angular $\vec{\alpha}$ (la derivada de $\vec{\omega}$) viene determinada por la relación $\vec{M} = m\rho^2\vec{\alpha}$.

Vemos así que, para un movimiento circular, la cantidad $I = m\rho^2$ es el análogo angular a la masa en la definición del momento $\vec{P} = m\vec{v}$ y en la segunda ley de Newton $\vec{F} = m\vec{a}$. Además, la velocidad escalar del móvil es $v = \rho\omega$, luego su energía cinética es

$$E_c = \frac{1}{2}m\rho^2\omega^2 = \frac{1}{2}I\omega^2,$$

que es el análogo angular de la expresión $E = \frac{1}{2}mv^2$. Por ello, la cantidad $I = m\rho^2$ se llama *momento de inercia* del móvil respecto de O' .

Momento de inercia de un sólido rígido Cuando un sólido rígido gira, lo hace necesariamente respecto de un eje, de modo que cada uno de sus puntos describe una trayectoria circular respecto de un punto del eje que no es el mismo para todos ellos. Si queremos definir un momento de inercia asociado a la rotación del sólido, nos encontramos con que conviene definirlo respecto de un punto distinto para cada una de sus partículas componentes. Por ello, definimos el *momento de inercia* de un móvil puntual respecto de una recta R como la cantidad $I = m\rho^2$, donde ρ es la distancia del móvil a la recta. El *momento de inercia* de un sistema de partículas respecto de una recta R se define como la suma de los momentos de inercia de sus partículas:

$$I = \sum_i m_i \rho_i^2.$$

Fijemos un sistema de referencia cuyo origen O esté sobre la recta R y sea $\vec{\omega}$ un vector director de R . Podemos suponer que $\vec{\omega}$ tiene módulo 1. Entonces, si \vec{r} es el vector de posición de una partícula de masa m , el producto $\vec{r} \cdot \vec{\omega}$ es la distancia de O al punto O' de R más cercano a \vec{r} (véase la figura anterior), luego la distancia ρ de \vec{r} a R cumple que

$$\rho^2 = \vec{r}^2 - (\vec{r} \cdot \vec{\omega})^2.$$

Vamos a expresar esta distancia en términos de las coordenadas de los vectores, teniendo en cuenta que $\omega_1^2 + \omega_2^2 + \omega_3^2 = 1$:

$$\begin{aligned} \rho^2 &= x_1^2 + x_2^2 + x_3^2 - x_1^2\omega_1^2 - x_2^2\omega_2^2 - x_3^2\omega_3^2 - 2x_1x_2\omega_1\omega_2 - 2x_1x_3\omega_1\omega_3 - 2x_2x_3\omega_2\omega_3 \\ &= x_1^2\omega_1^2 + x_1^2\omega_2^2 + x_1^2\omega_3^2 + x_2^2\omega_1^2 + x_2^2\omega_2^2 + x_2^2\omega_3^2 + x_3^2\omega_1^2 + x_3^2\omega_2^2 + x_3^2\omega_3^2 \\ &\quad - x_1^2\omega_1^2 - x_2^2\omega_2^2 - x_3^2\omega_3^2 - 2x_1x_2\omega_1\omega_2 - 2x_1x_3\omega_1\omega_3 - 2x_2x_3\omega_2\omega_3 = \sum_{k,l} \omega_k A_{kl} \omega_l, \end{aligned}$$

donde

$$\begin{aligned} A_{11} &= x_2^2 + x_3^2, & A_{12} &= A_{21} = -x_1x_2, \\ A_{22} &= x_1^2 + x_3^2, & A_{13} &= A_{31} = -x_1x_3, \\ A_{33} &= x_1^2 + x_2^2, & A_{23} &= A_{32} = -x_2x_3. \end{aligned}$$

Una expresión conjunta para estas nueve cantidades es:

$$A_{kl} = \sum_i x_i^2 \delta_{kl} - x_k x_l,$$

donde δ_{kl} es la delta de Kronecker.

Por consiguiente, el momento de inercia de un sistema de partículas respecto de una recta que pase por el origen de coordenadas con vector director unitario $\vec{\omega}$ es

$$\vec{\omega} \cdot I \cdot \vec{\omega}^t = \sum_{k,l} \omega_k I_{kl} \omega_l,$$

donde

$$I_{kl} = \sum_i m_i \left(\sum_{j=1}^3 x_{ij}^2 \delta_{kl} - x_{ik} x_{il} \right).$$

La matriz simétrica I se llama *tensor de inercia* del sistema de partículas respecto del punto O . (Notemos que I depende de O , pero no del eje R porque no depende de $\vec{\omega}$.) De la definición del momento de inercia (en términos de distancias) se sigue que es independiente del sistema de referencia, por lo que lo mismo vale para el tensor de inercia, aunque esto hay que entenderlo bien: si consideramos otro sistema de referencia O' (con el mismo origen), existirá una matriz de cambio de base A tal que el vector que respecto a O tiene coordenadas $\vec{\omega}$ tiene coordenadas $\vec{\omega}'A$ respecto de O' . Entonces $\vec{\omega} I \vec{\omega}^t = \vec{\omega}' A I A^t \vec{\omega}'^t = \vec{\omega}' I' \vec{\omega}'^t$ para todo vector $\vec{\omega}$, luego $I' = A I A^t$. Cuando decimos que I sólo depende del origen O del sistema de referencia queremos decir que, al cambiar de sistema de referencia (con el mismo origen O), I se transforma según la relación $I' = A I A^t$, es decir, que su variación no es más que el mero reflejo de la variación del criterio para asignar coordenadas a los puntos.⁶

Energía cinética Vamos a obtener una expresión útil para la energía cinética de un sólido rígido. Para ello fijamos un observador inercial O y otro observador O' con origen en el centro de masas del sólido y cuyos ejes coordenados se muevan con él, de modo que todos los puntos del sólido estén en reposo respecto de O' . La fórmula (C.1) nos da que la velocidad (para O) de cualquiera de las partículas del sólido es

$$\vec{v}_i = \vec{V} + \vec{\omega} \times \vec{r}'_i, \quad (\text{C.4})$$

donde \vec{V} es la velocidad de su centro de masas, $\vec{\omega}$ es la velocidad angular de O' respecto de O y $\vec{r}'_i = \vec{r}_i - \vec{R}$, donde \vec{R} es la posición del centro de masas.

Entonces

$$\frac{1}{2} m_i v_i^2 = \frac{1}{2} m_i V^2 + m_i \vec{V} (\vec{\omega} \times \vec{r}'_i) - m_i \vec{V} (\vec{\omega} \times \vec{R}) + \frac{1}{2} m_i (\vec{\omega} \times \vec{r}'_i)^2,$$

⁶Equivalentemente, podemos pensar que el tensor de inercia es una aplicación que a cada punto O le asigna una forma bilineal simétrica $\mathcal{J}_O : \mathbb{R}^3 \times \mathbb{R}^3 \rightarrow \mathbb{R}$, de modo que la matriz I_O a la que hemos llamado tensor de inercia es, en estos términos, la matriz de \mathcal{J}_O respecto de una base ortonormal prefijada. En estos términos \mathcal{J}_O es un tensor en el sentido de la geometría diferencial, igual que lo es, por ejemplo, el tensor métrico.

y sumando obtenemos la energía cinética del sólido:

$$\begin{aligned} E_c &= \frac{1}{2}MV^2 + M\vec{V}(\vec{\omega} \times \vec{R}) - M\vec{V}(\vec{\omega} \times \vec{R}) + \frac{1}{2}\sum_i m_i(\vec{\omega} \times \vec{r}'_i)^2 \\ &= \frac{1}{2}MV^2 + \frac{1}{2}\omega^2 \sum_i m_i \rho_i^2, \end{aligned}$$

donde ρ_i es la distancia de \vec{r}'_i a la recta que pasa por O' con la dirección de $\vec{\omega}$. Si I es el tensor de inercia del sólido respecto al sistema de referencia O' , el sumatorio es igual a

$$\frac{\vec{\omega}}{\omega} \cdot I \cdot \left(\frac{\vec{\omega}}{\omega}\right)^t,$$

(porque I actúa sobre vectores unitarios) luego, en total,

$$E_c = \frac{1}{2}MV^2 + \frac{1}{2}\vec{\omega} \cdot I \cdot \vec{\omega}^t.$$

El primer sumando representa la energía de traslación del centro de masas, mientras que el segundo es la suma de las energías de rotación de las partículas del sólido.

Momento angular Calculamos ahora el momento angular de un sólido rígido. Partimos nuevamente de (C.4):

$$\begin{aligned} \vec{L} &= \sum_i m_i \vec{r}_i \times \vec{v}_{O'} + \sum_i m_i (\vec{R} + \vec{r}'_i) \times (\vec{\omega} \times \vec{r}'_i) \\ &= M\vec{R} \times \vec{V} + \sum_i \vec{R} \times (\vec{\omega} \times m_i \vec{r}'_i) + \sum_i m_i \vec{r}'_i \times (\vec{\omega} \times \vec{r}'_i) \\ &= \vec{R} \times \vec{P} + \sum_i m_i \vec{r}'_i \times (\vec{\omega} \times \vec{r}'_i), \end{aligned}$$

donde hemos usado que $\sum_i m_i \vec{r}'_i = \vec{0}$, pues es M veces la posición del centro de masas respecto de O' . Vamos a llamar $\vec{L}(\vec{\omega})$ al último término de la última expresión, que claramente es una aplicación lineal. Vamos a probar que su matriz I respecto a la base $\vec{u}_1, \vec{u}_2, \vec{u}_3$ de O' es el tensor de inercia respecto de O' . Para ello aplicamos la fórmula de Lagrange [G 5.13]. Así:

$$\vec{L}(\vec{\omega}) = \vec{\omega} \sum_i m_i r_i'^2 - \sum_i m_i \vec{r}'_i (\vec{r}'_i \cdot \vec{\omega}).$$

Ahora, I_{kl} es la coordenada l -ésima de

$$\vec{L}(\vec{u}_k) = \vec{u}_k \sum_i m_i r_i'^2 - \sum_i m_i x'_{ik} \vec{r}'_i,$$

luego

$$I_{kl} = \sum_i m_i (\vec{r}'_i{}^2 \delta_{kl} - x'_{ik} x'_{il}).$$

Esto prueba que I es el tensor de inercia del sólido. En definitiva, la expresión para el momento angular de un sólido rígido es

$$\vec{L} = \vec{R} \times \vec{P} + \vec{\omega} \cdot I.$$

El segundo término es el análogo angular de $m\vec{v}$ y contiene el momento angular asociado a la rotación del sólido, mientras que el primer término es el momento angular del centro de masas.

C.6 Distribuciones continuas de materia

Hasta ahora hemos considerado los sólidos rígidos como formados por conjuntos finitos de partículas puntuales. Esto significa admitir que una cantidad de masa puede ocupar un punto, una porción del espacio sin volumen. Esta abstracción es una forma práctica de reflejar que en muchos problemas es irrelevante el volumen concreto que ocupe una determinada cantidad de masa, pero en ocasiones es más conveniente adoptar el convenio contrario, es decir, considerar que cualquier cantidad de masa, por pequeña que sea, ha de ocupar un cierto volumen no nulo en el espacio. Esto nos lleva al concepto de densidad:

En principio, podemos definir la *densidad* de una región $V \subset \mathbb{R}^3$ del espacio de volumen⁷ finito no nulo $v(V)$ como el cociente

$$d(V) = \frac{m(V)}{v(V)}, \quad (\text{C.5})$$

donde $m(V)$ es la masa contenida en V . La unidad de densidad en el *Sistema Internacional* es el kg/m^3 . Así, por ejemplo, al dividir la masa estimada de la Tierra entre su volumen se obtiene como resultado $d = 5\,515.3 \text{ kg}/\text{m}^3$. Ahora bien, está claro que esto es un promedio, y que no podemos afirmar que cada metro cúbico del planeta contiene 5 515.3 kg de masa. No es lo mismo un metro cúbico tomado en una mina de plomo que otro de aire tomado en una caverna subterránea.

Para tener en consideración el hecho de que la densidad dentro de una misma región V del espacio puede ser distinta en distintas zonas de la misma, consideraremos una función de densidad ρ , que depende de la posición (y del tiempo, teniendo en cuenta que la materia se mueve), de modo que la masa contenida en una región del espacio V vendrá dada por

$$m(V) = \int_V \rho \, dx \, dy \, dz. \quad (\text{C.6})$$

En términos de formas diferenciales, si llamamos $dv = dx \, dy \, dz$ al *elemento de volumen* en \mathbb{R}^3 , es decir, la forma diferencial que, al integrarla sobre una región V da el volumen de V , entonces la forma diferencial $dm = \rho \, dv$ es el

⁷Técnicamente hemos de suponer que V es un conjunto medible Lebesgue, lo cual no es ninguna restricción desde un punto de vista físico.

elemento de masa, la forma diferencial que al integrarla sobre una región V da la masa que contiene.

En el lenguaje del siglo XVIII, que, al no ser riguroso, no sirve para demostrar nada, pero es útil para interpretar resultados, podemos decir que la integral (C.6) es la “suma” para los infinitos puntos (x, y, z) de V de la masa infinitesimal dm que contiene una región de volumen infinitesimal dv que contenga a (x, y, z) .

Es ésta interpretación infinitesimal la que induce a definir análogos continuos de los conceptos que hemos definido para sistemas de partículas puntuales. Por ejemplo, el *centro de masas* de un sólido determinado por una función de densidad ρ sobre una región V se define como el punto

$$\vec{R} = \frac{1}{m(V)} \int_V \vec{x} dm. \quad (\text{C.7})$$

Similarmente, su *tensor de inercia* se define como la matriz I dada por

$$I_{ij} = \int_V \left(\sum_{k=1}^3 x_k^2 \delta_{ij} - x_i x_j \right) dm. \quad (\text{C.8})$$

En principio debemos considerar estas definiciones como independientes de las correspondientes para sistemas de partículas, pero conviene observar que se puede pasar de aquéllas a éstas mediante un proceso de paso al límite.⁸

Las fórmulas para la energía cinética y el momento angular de un sistema de partículas pueden deducirse análogamente en el caso continuo.

Ejemplo Vamos a calcular la energía cinética de una rueda o una esfera homogénea de radio R que avanza sin deslizarse sobre un plano horizontal. Tomamos el sistema de referencia de modo que el centro de masas avanza en la dirección del eje X y el giro se realiza alrededor del eje Y . De este modo, la velocidad del centro de masas es de la forma $\vec{V} = (0, V, 0)$ y la velocidad angular $\vec{\omega}$ del sistema de referencia O' es $\vec{\omega} = (0, \omega, 0)$. Que el movimiento se produzca sin deslizamiento significa que el espacio recorrido por el punto de contacto sobre el suelo coincida con el recorrido sobre la circunferencia de contacto, es decir:

$$\int_{t_0}^{t_1} V dt = \int_{t_0}^{t_1} R\omega dt,$$

⁸Concretamente, suponiendo que la función ρ es continua, dado $\delta > 0$, existe una partición del espacio en cubos de forma que la integral (C.7) dista menos de δ de la suma $\sum_i \vec{r}_i \rho(\vec{r}_i) \Delta v_i$, donde Δv_i es el volumen del cubo i -ésimo y \vec{r}_i es un punto contenido en dicho cubo. Refinando la partición, podemos suponer además que la suma $\sum_i \rho(\vec{r}_i) \Delta v_i$ dista menos de δ de la integral $\int_V \rho dv = m(V)$. Por la continuidad de la función x/y , dado $\epsilon > 0$ podemos encontrar una partición en cubos tal que (C.7) dista menos de ϵ del cociente $\sum_i m_i \vec{r}_i / \sum_i m_i$, donde hemos llamado $m_i = \rho(\vec{r}_i) \Delta v_i$. Así pues, \vec{R} dista menos de ϵ del centro de masas de un cierto sistema de partículas puntuales resultante de “concentrar” en un número finito de puntos de V las masas determinadas por la función de densidad ρ .

para cualquier par de instantes t_0 y t_1 . Esto equivale a que $\omega = V/R$. Teniendo en cuenta que el momento de inercia I_0 respecto al eje de giro se calcula aplicando el tensor de inercia a un vector unitario en la dirección del eje (es decir, que $I_0 = (\vec{\omega}/\omega)I(\vec{\omega}/\omega)^t$), vemos que la energía cinética del móvil es

$$E_c = \frac{1}{2}MV^2 + \frac{1}{2}I_0\omega^2 = \frac{1}{2}MV^2 + \frac{I_0V^2}{2R^2}.$$

Para el caso de una esfera $I_0 = 2MR^2/5$ (véase la pág. 394), luego

$$E_c = \frac{7}{10}MV^2.$$

Para una rueda cilíndrica es fácil ver que $I_0 = MR^2/2$, luego

$$E_c = \frac{3}{4}MV^2.$$

■

Bibliografía

- [1] Borden, R.S. *A course in advanced calculus*. North Holand, Amsterdam, 1983.
- [2] Chernoff, P.R. *Pointwise convergence of Fourier Series*. The American Mathematical Monthly, Vol. 87, No. 5 (1980) pp. 399–400
- [3] Corwin, L.J. Szczarba, R.H. *Multivariable calculus*. Marcel Dekker, New York, 1982.
- [4] Do Carmo, M.P. *Geometría diferencial de curvas y superficies*. Alianza, Madrid, 1990.
- [5] Elsgoltz, L. *Ecuaciones diferenciales y cálculo variacional*. Ed. Mir, Moscú, 1977.
- [6] H. Flanders, *Differentiation Under the Integral Sign*, The American Mathematical Monthly, Vol. 80, No. 6 (1973) pp. 615–627
- [7] Greub, W., Halperin, S., Vanstone, R. *Connections, curvature, and cohomology* Vol. 1. Academic Press, New York, 1972.
- [8] Hu, S.T. *Differentiable manifolds*. Holt, Rinehart and Winston, New York, 1969.
- [9] John, F. *Partial differential equations*. Springer, New York, 1982.
- [10] Lang, S. *Differential Manifolds*. Addison Wesley, Reading, Mass. 1972
- [11] Perron, O. *Die Lehre von den Kettenbrüchen*. Teubner, Stuttgart, 1954
- [12] Ramo, S. Whinnery, J.R., van Duzer, T. *Campos y ondas*. Pirámide, Madrid, 1965.
- [13] Rudin, W. *Análisis real y complejo*. Mc. Graw Hill, Madrid, 1988.
- [14] Santaló, L. A. *Vectores y tensores*. Ed. Universitaria de Buenos Aires, 1968.

Índice de Materias

- abierta (aplicación), 66
- abierto, 43
 - básico, 45
- absolutamente
 - convergente (serie), 89
 - continua (medida), 375
- aceleración, 226, 439
 - angular, 293
 - geodésica, 262
- acotado, 56
- Alexandroff (compactificación de), 104
- álgebra de conjuntos, 320
- anillo (de conjuntos), 304
- arco, 109, 216
 - rectificable, 220
- argumento, 178
- argumento del seno hiperbólico, 185
- arquimediano (anillo ordenado), 4

- Banach (espacio de), 114
- Barrow (regla de), 346
- base, 45
 - de entornos, 45
 - ortonormal, 125
- Bessel (desigualdad de), 124
- Borel (álgebra, medida), 330

- cóncava (función), 150
- cantidad de movimiento, 449
- cardioide, 223
- carta, 240
- Cauchy
 - producto de, 90
 - sucesión de, 25
- celda, 302
- centrífuga (fuerza), 233
- centro de masas, 391, 455

- cerrado, 52
- Christoffel (símbolos de), 261
- cicloide, 222
- cilindro, 246
- cinemática, 437
- circunferencia osculatriz, 226
- clase monótona, 321
- clausura, 52
- compacto, 98
- compleción, 327, 407
- completo
 - conjunto ordenado, 30
 - espacio métrico, 27
- componente
 - arco-conexa, 111
 - conexa, 109
- condicionalmente convergente, 89
- conexo, 107
- conjugados (números), 367
- cono, 246
- contenido, 302
- continuidad, 59
- contractiva (aplicación), 278
- convergencia, 13
 - de funciones, 71
 - de sucesiones, 80
 - uniforme, 128
- convergente, 415
- convexa (función), 150
- convexo (conjunto), 110
- coordenadas, 240
- Coriolis (fuerza de), 233
- coseno, 175
 - hiperbólico, 183
- cubrimiento, 97, 99
- cuerpo métrico, 11
- curva parametrizada, 218

- curvatura, 225
 - geodésica, 262
 - media, 269
 - normal, 265
- D'Alembert (criterio de), 86
- densidad, 462
- denso, 18, 55
- derivada, 140
 - covariante, 260
 - direccional, 196
 - parcial, 197, 254
- difeomorfismo, 250
- diferenciabilidad, 199
- diferenciable, 250
- diferencial, 200, 251
- dinámica, 437
- Dirac (delta de), 372
- dirección principal, 268
- disco de convergencia, 164
- discreta (métrica, topología), 45
- diseminado (conjunto), 135
- distancia, 12, 41
- elemento de longitud, 257
- energía, 454
 - cinética, 451
 - potencial, 452
- entorno, 44
 - básico, 45
- esfera, 245
- espacio
 - arco-conexo, 109
 - compacto, 98
 - conexo, 107
 - de Hilbert, 118
 - discreto, 45
 - localmente conexo, 110
 - métrico, 12, 42
 - completo, 27
 - medida, 320
 - normado, 40
 - precompacto, 113
 - prehilbertiano, 118
 - σ -compacto, 360
 - tangente, 248
 - topológico, 43
 - vectorial topológico, 63
- estereográfica (proyección), 69
- Euler (constante de), 350
- Euler-Cauchy (ecuación de), 290
- exponencial, 169, 173
- extremo, 147, 212
- factorial (función), 347
- figura elemental, 303, 351
- flujo, 284
- fracción continua, 415
- Frenet
 - fórmulas de, 228
 - triedro de, 227
- frontera, 53
- fuerza, 446
 - interna/externa, 455
- Gamma (función), 347
- geodésica, 262
- geométrica (serie), 10
- gráfica, 66
- gradiente, 203
- armónica (función), 291
- Hausdorff (espacio de), 55
- Hilbert (espacio de), 118
- Hölder (desigualdad de), 368
- homeomorfismo, 65
- Hooke (ley de), 449
- indeterminación, 79
- inercia (principio de), 444
- inercial (sistema), 445
- inmersión
 - isométrica, 32
- inmersión isométrica, 32
- integrable Lebesgue (función), 342
- integral de Lebesgue, 342
- interior, 52
- intersección finita (propiedad), 98
- intervalo, 3
- isometría, 32, 259
- isomorfismo topológico, 68
- jacobiana, 202

- Kepler
 - primera ley, 292
 - segunda ley, 390
 - tercera ley, 390
- kilogramo, 446
- límite, 71
 - de una sucesión, 13
 - superior, 163
- laplaciano, 291
- Leibniz (criterio), 88
- Lipschitz, 62
- localmente compacto (espacio), 103
- logaritmo, 171, 173
- longitud de un arco, 220
- mínimo, 147, 212
- máximo, 147, 212
- marea (fuerza de), 236
- media aritmética/geométrica, 174
- medible (aplicación), 334
- medida, 318
 - completa, 321
 - de Jordan, 309
 - de Lebesgue, 329, 395, 400
 - exterior, 322
 - de Lebesgue, 329
 - finita, 320
 - finitamente aditiva, 306
 - producto, 354
 - regular, 360
 - σ -finita, 320
 - signada, 372
 - unitaria, 320
- metro, 435
- Minkowski (desigualdad de), 368
- Möbius (cinta de), 248
- momento, 294
 - angular, 294, 450, 457
 - de inercia, 459
 - de una fuerza, 450
 - lineal, 449
- mutuamente singulares (medidas), 375
- número
 - real, 411
- Newton, 447
 - ley de gravitación de, 292
 - primera ley, 444
 - primera ley de, 217
 - segunda ley, 446
 - tercera ley, 448
- norma, 40
- observador, 435
- ortogonalidad, 119
- Parseval (identidad de), 126
- parte entera, 5
- partición
 - de Hahn, 380
 - de la unidad, 106
- péndulo de Foucault, 297
- peso, 447
- precompacto, 113
- prehilbertiano, 118
- primer axioma de numerabilidad, 81
- primera categoría, 136
- primitiva, 186
- producto
 - de variedades, 248
 - escalar, 117
 - mixto, 228
- pseudoesfera, 272
- punto
 - adherente, 82
 - aislado, 55
 - de acumulación, 54
 - de Lebesgue, 384
- radio
 - de convergencia, 164
 - de curvatura, 226
- real (número), 411
- rectángulo medible, 351
- rectificable (arco), 220
- regla de la cadena, 145, 206
- regular (medida), 360
- resto de Taylor, 161
- segmento, 110
- segunda categoría, 136
- segundo, 434
- seno, 175

- hiperbólico, 183
- serie, 10
 - de potencias, 163
 - de Taylor, 161
- sólido rígido, 455
- soporte, 105
- subbase, 48
- subcubrimiento, 97
- sucesión
 - acotada, 15
 - de Cauchy, 25
 - monótona, 9
- superficie de revolución, 245
- tangente, 140, 217
 - hiperbólica, 183
- Taylor (polinomio), 159
- tensor
 - de inercia, 391, 460, 463
 - métrico, 256
- Teorema
 - de Baire, 135, 136
 - de cambio de variable, 387
 - de Cauchy, 148
 - de Fubini, 358
 - de Hahn, 379
 - de inyectividad local, 215
 - de L'Hôpital, 152–154
 - de la convergencia dominada, 343
 - de la convergencia monótona, 339
 - de la función compuesta, 145, 206
 - de la función implícita, 243
 - de la función inversa, 144, 213
 - de Lebesgue-Radon-Nikodým, 377
 - de los valores intermedios, 112
 - de Lusin, 366
 - de Meusnier, 266
 - de Riesz, 361, 381
 - de Rolle, 148
 - de Schwarz, 211
 - de Stone-Weierstrass, 131
 - de Taylor, 161
 - de Tychonoff, 101
 - del valor medio, 148
 - egregium de Gauss, 272
 - topología, 43
 - euclídea, 44, 68
 - producto, 49
 - relativa, 50
 - usual, 44
 - toro, 246
 - torsión, 227
 - trabajo, 450
 - tractriz, 224
 - transporte paralelo, 285, 286
 - trayectoria, 438
 - umbilical (punto), 268
 - Urysohn (lema), 105
 - valor absoluto, 11
 - variación total, 373
 - variedad, 240
 - tangente, 248
 - vector
 - binormal, 227
 - normal, 225
 - tangente, 218
 - velocidad, 217, 438
 - angular, 293, 443
 - escalar, 438