Gauss, la estadística y el planeta enano Ceres

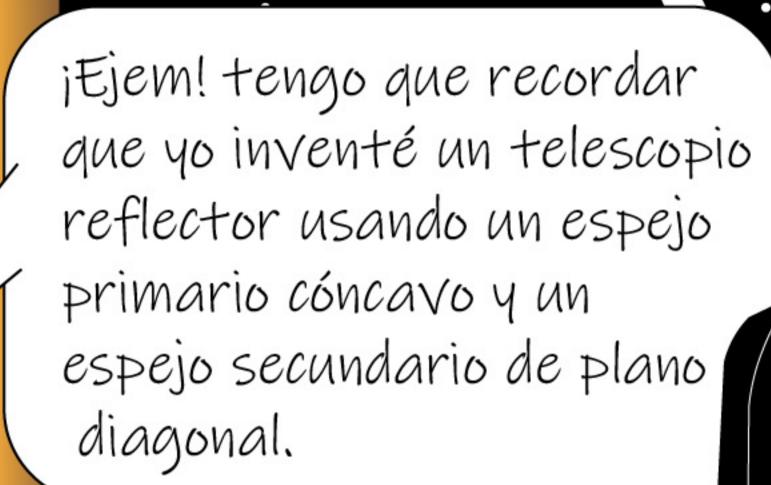
Cuando pensamos en agrias disputas matemáticas siempre nos viene a la memoria el enfrentamiento que hubo entre Sir **Isaac Newton** y **Gottfried Wilhelm Leibniz** en el siglo XVII por la primacía en el descubrimiento del cálculo infinitesimal. En este cómic os contaremos la historia de otra rivalidad entre otros insignes matemáticos y en otro campo de las matemáticas: La **ESTADÍSTICA**.



El 1 de enero de 1801, recién estrenado el siglo XIX, el astrónomo italiano Giuseppe Piazzi anunció el descubrimiento de un nuevo planeta que giraba alrededor del Sol en una órbita entre las de Marte y Jupiter. El planeta, al que bautizaron CERES como la diosa romana de la agricultura y la fecundidad, era mucho más pequeño que los siete planetas conocidos hasta entonces.



Durante 40 días, hasta la noche del 11 de febrero, Piazzi siguió al objeto en su viaje por el espacio.



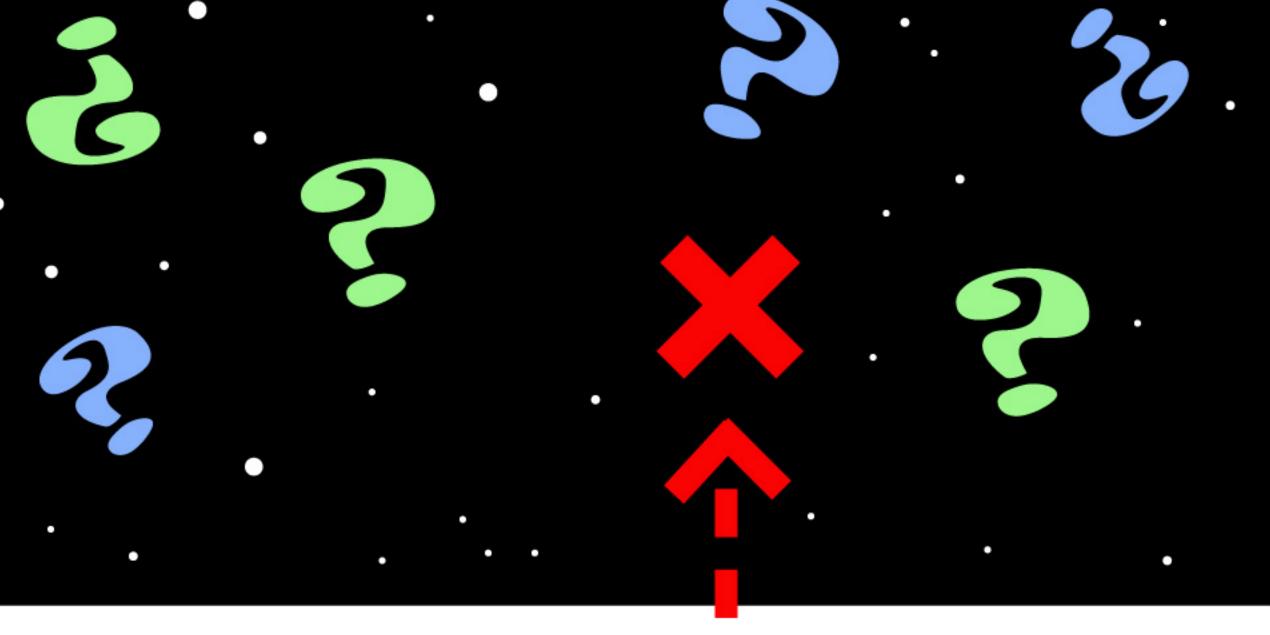
¡Si Newton!, pero ya se habían fabricado otros telescopios. Algunos por el mismísmo **Galileo Galilei** a quien se le ocurrió no utilizar

estos artilugios para fines topográficos o militares ¡sino para observar el cielo!.

Archimedes' Tub



Pero tras regresar del otro lado del Sol Ceres desapareció del cielo nocturno perdido en la inmensidad del firmamento. A pesar de disponer de Leyes que describían las órbitas y movimiento de los planetas los astrónomos del siglo XIX no disponían



matemáticas

a partir de los

de suficientes herramientas para calcular su trayectoria escasos datos recogidos por Piazzi.

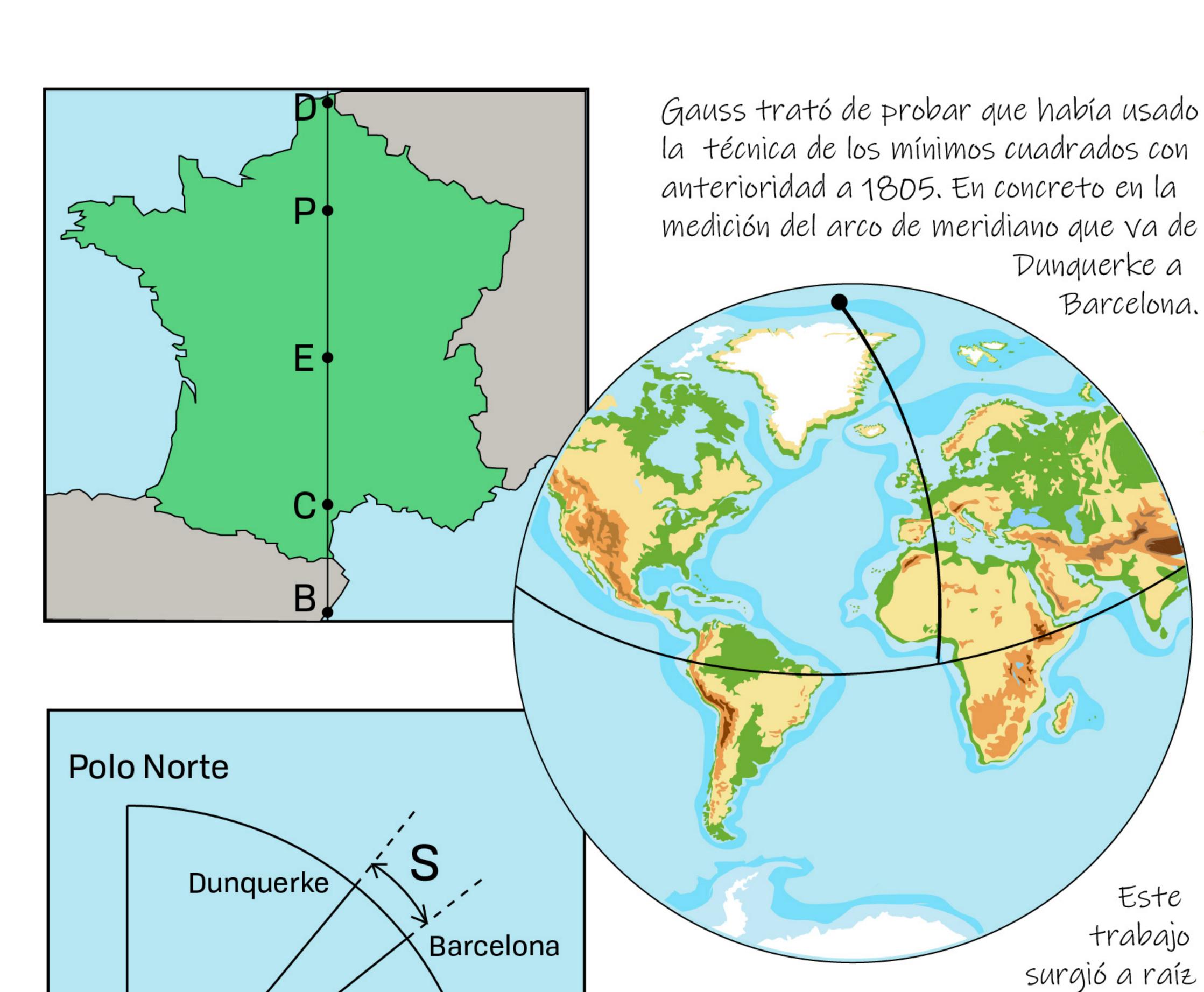
¡Por cierto! La Ley de Gravitación Universal de la que se deducen las Leyes de Kepler a partir del cálculo infinitesimal también es obra mía...

iiBasta Newton!! Hoy no hemos venido a hablar de ti, sino del único científico que con tan solo 24 años señaló el lugar exacto en el que debían apuntar sus telescopios los astrónomos para encontrar de nuevo a Ceres.

oy no plar de ntífico raños to en rar sus ónomos nuevo

Archimedes' Tub





de que la Academia de Ciencias Francesa decidiera en 1793 basar el sistema métrico en una unidad, el metro, que se definiría como la 10.000.000-ésima parte del cuadrante de meridiano.

Este

Mi sola

La mala fortuna quiso que los cálculos de Gauss se perdieran y en 1831 se le sugirió la necesidad de que repitiera dichos cálculos con el fin de probar que utilizó para ellos el método de los mínimos cuadrados.

Ecuador

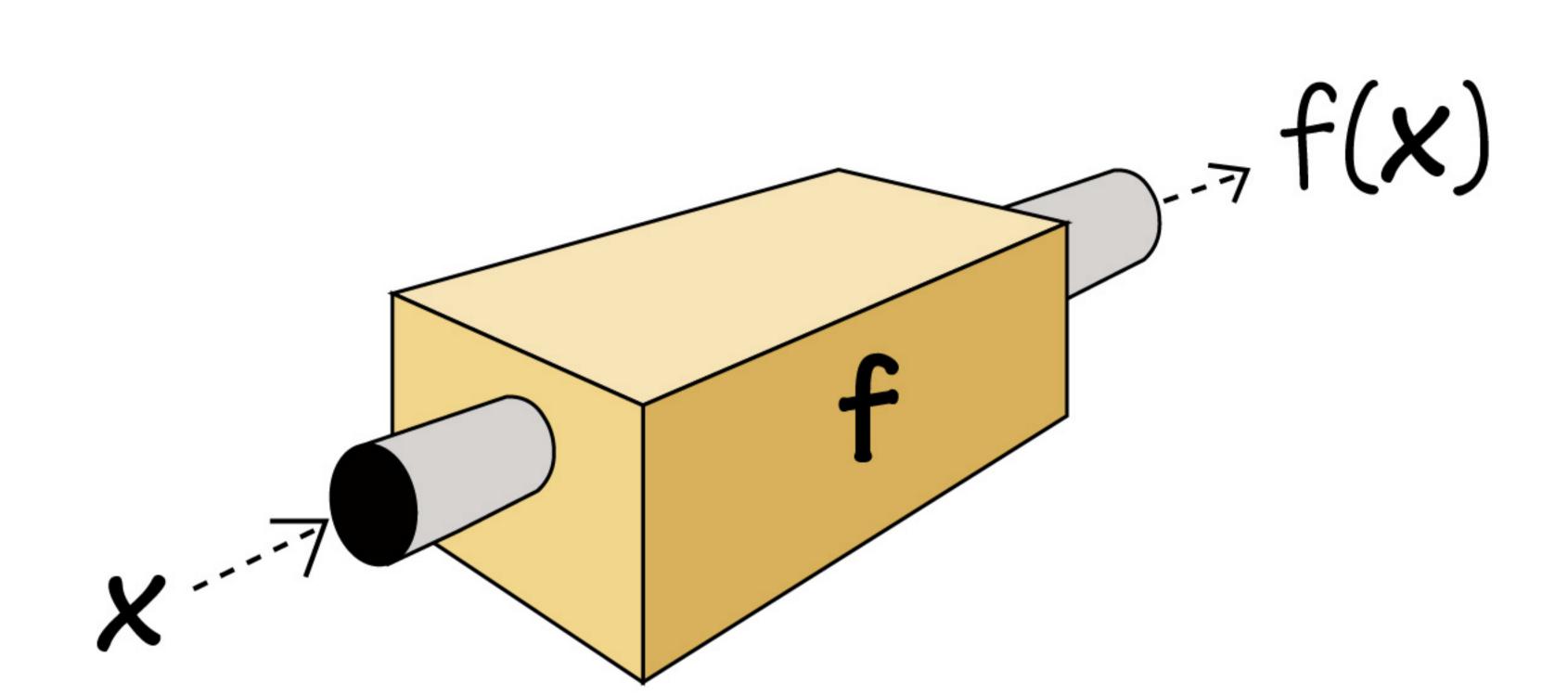
Gauss se negó categóricamente:





¿En qué consiste exáctamente el método de los MÍNIMOS CUADRADOS?

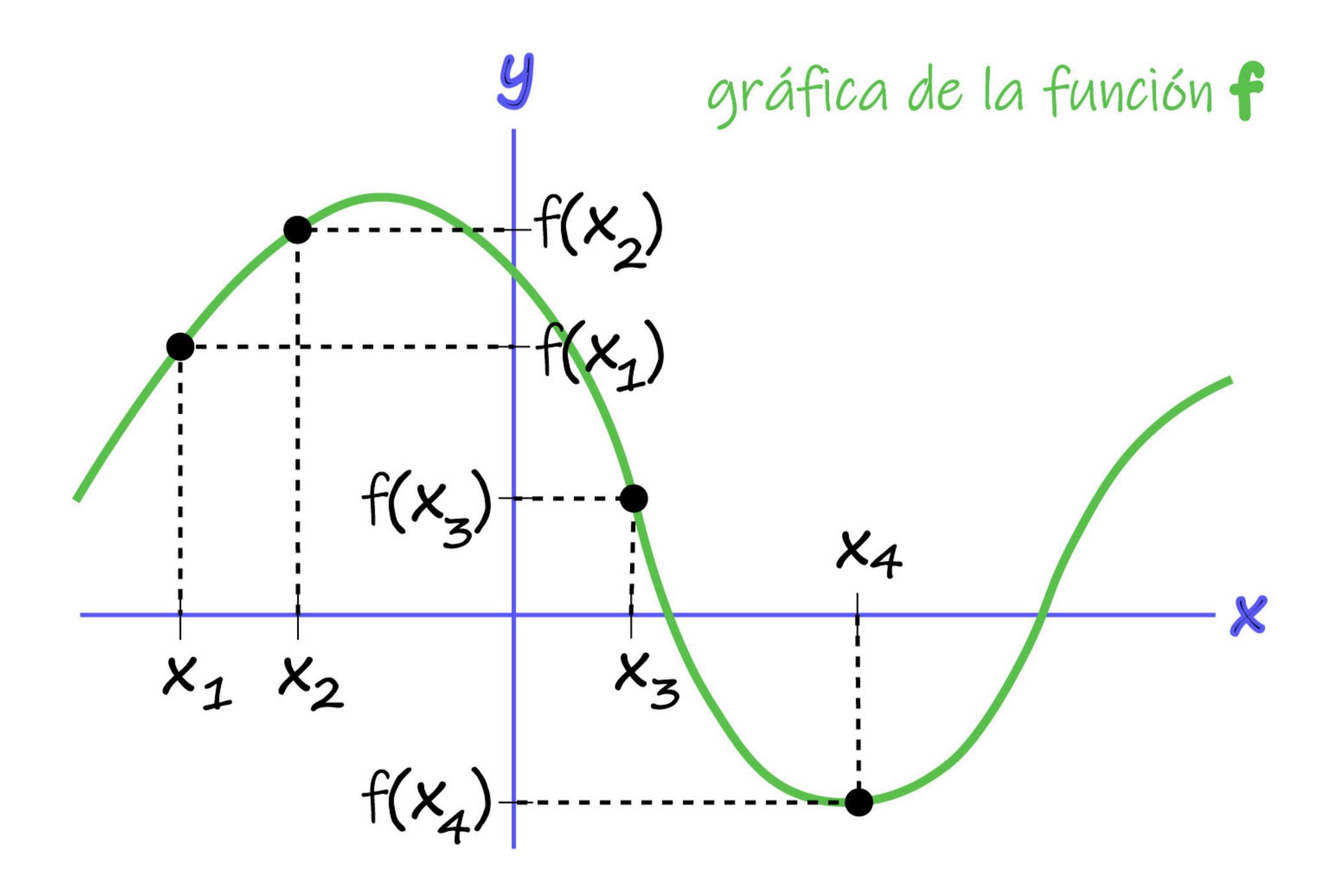
Todos estamos familiarizados con el concepto de **función**, que no es otra cosa que una regla, que denotaremos por f, que permite asociar a un cierto valor numérico x un único valor f(x).



Podemos Visualizar una funcion \mathbf{f} como si de una "máquina" se tratase. En ella introducimos un número \mathbf{x} , que tras ciertas manipulaciones se convierte en otro número $\mathbf{f}(\mathbf{x})$.

Otra forma de visualizar una función es a través de su gráfica.

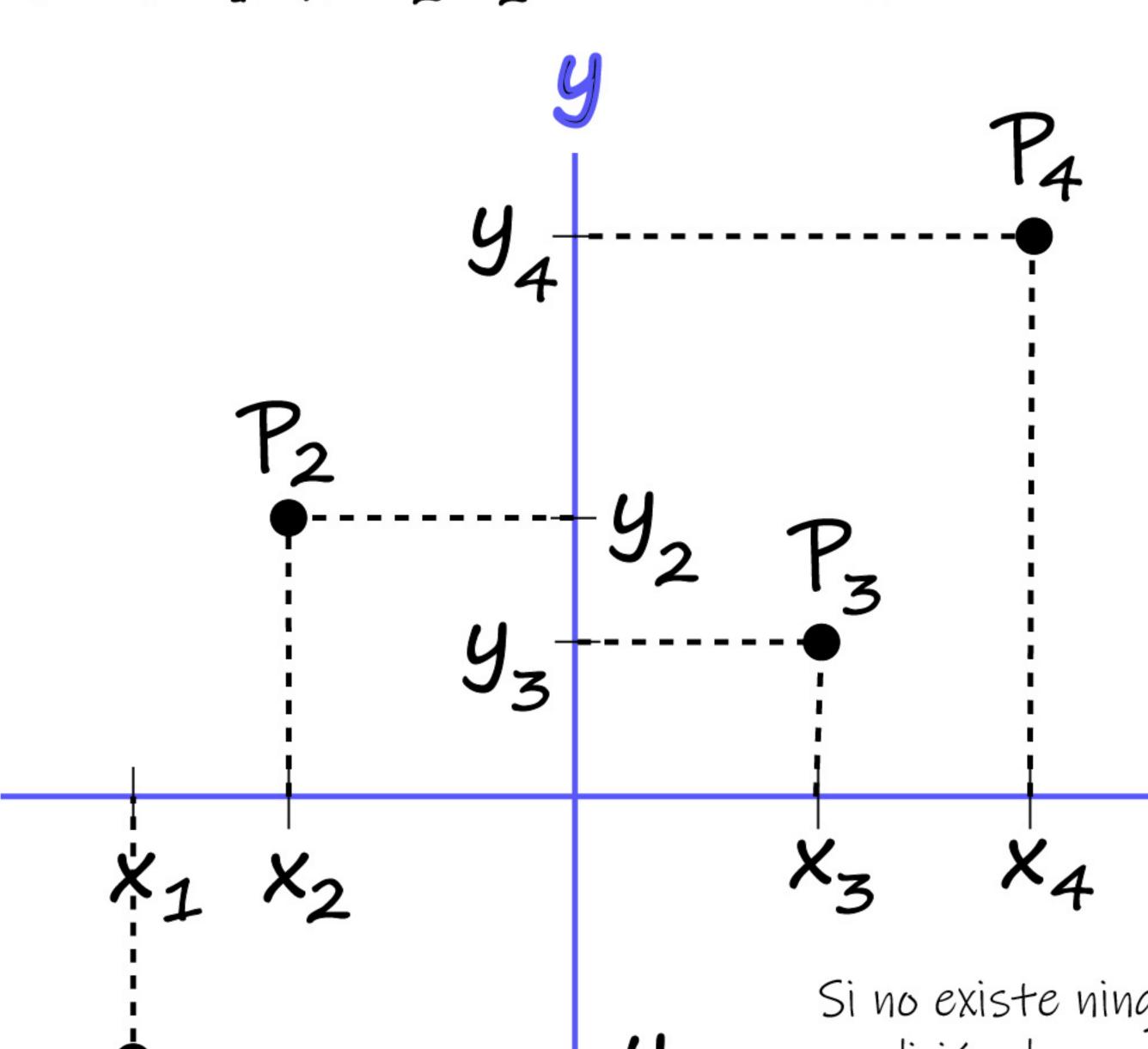
Si representamos en un eje horizontal. los valores x que se introducen en la función f, y en un eje vertical los valores f(x) que se van obteniendo, podemos dibujar cada par (x, f(x)) como un punto en el plano. Todos estos puntos determinan la **gráfica** de la función f.





Supongamos que tenemos una colección de puntos del plano

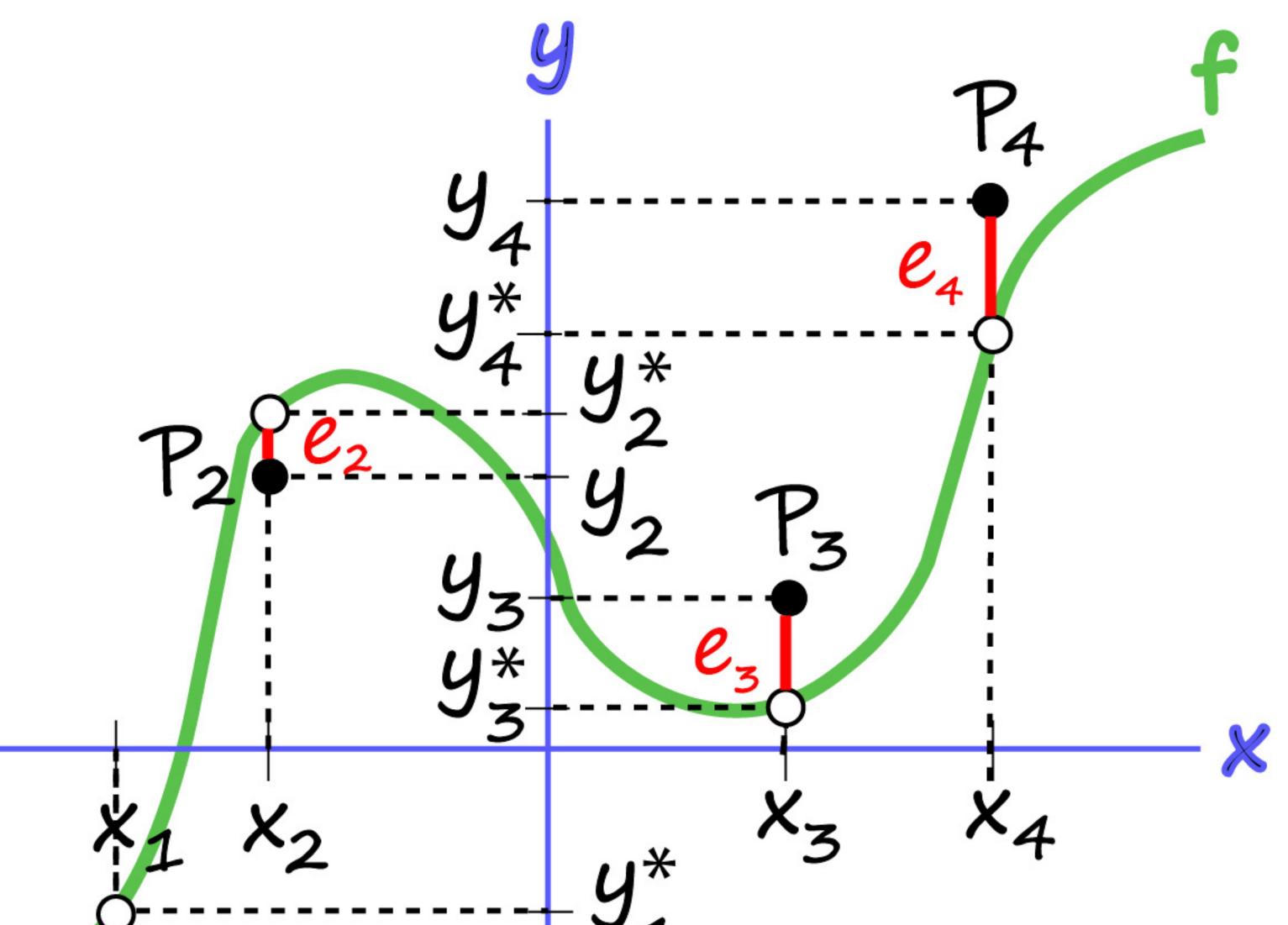
 $P_1 = (x_1, y_1), P_2 = (x_2, y_2), \dots, P_n = (x_n, y_n), y$ una determinada familia de funciones.



Nos gustaría saber si existe una función f de esa familia que pase por todos los puntos, esto es, tal que $f(x_1) = y_1$, $f(x_2) = y_2$, ..., $f(x_n) = y_n$.

Por ejemplo, los puntos pueden ser las anotaciones tomadas por Piazzi sobre la posición de Ceres, y la función **f** podría ser la órbita que queremos averiguar.

Si no existe ninguna función en la familia que cumpla la condición de pasar por todos los puntos, nos gustaría encontrar una que esté lo más "cerca" posible de cumplirla. ¿Cómo se consigue esto? Si tenemos que $f(x_1) = y_1^*$, $f(x_2) = y_2^*$, ..., $f(x_n) = y_n^*$. Podemos considerar el error como la diferencia $e_i = y_i^* - y_i$.



De este modo buscamos la función **f**, de la familia dada que minimize los errores **e**.

Para ello nótese que no sería muy inteligente minimizar la suma

$$\sum_{i=1}^{n} e_{i}$$

pues los errores pueden ser de diferente signo y darse el caso de que el sumatorio es nulo a pesar de que la función no pase por ningún punto **P**,.

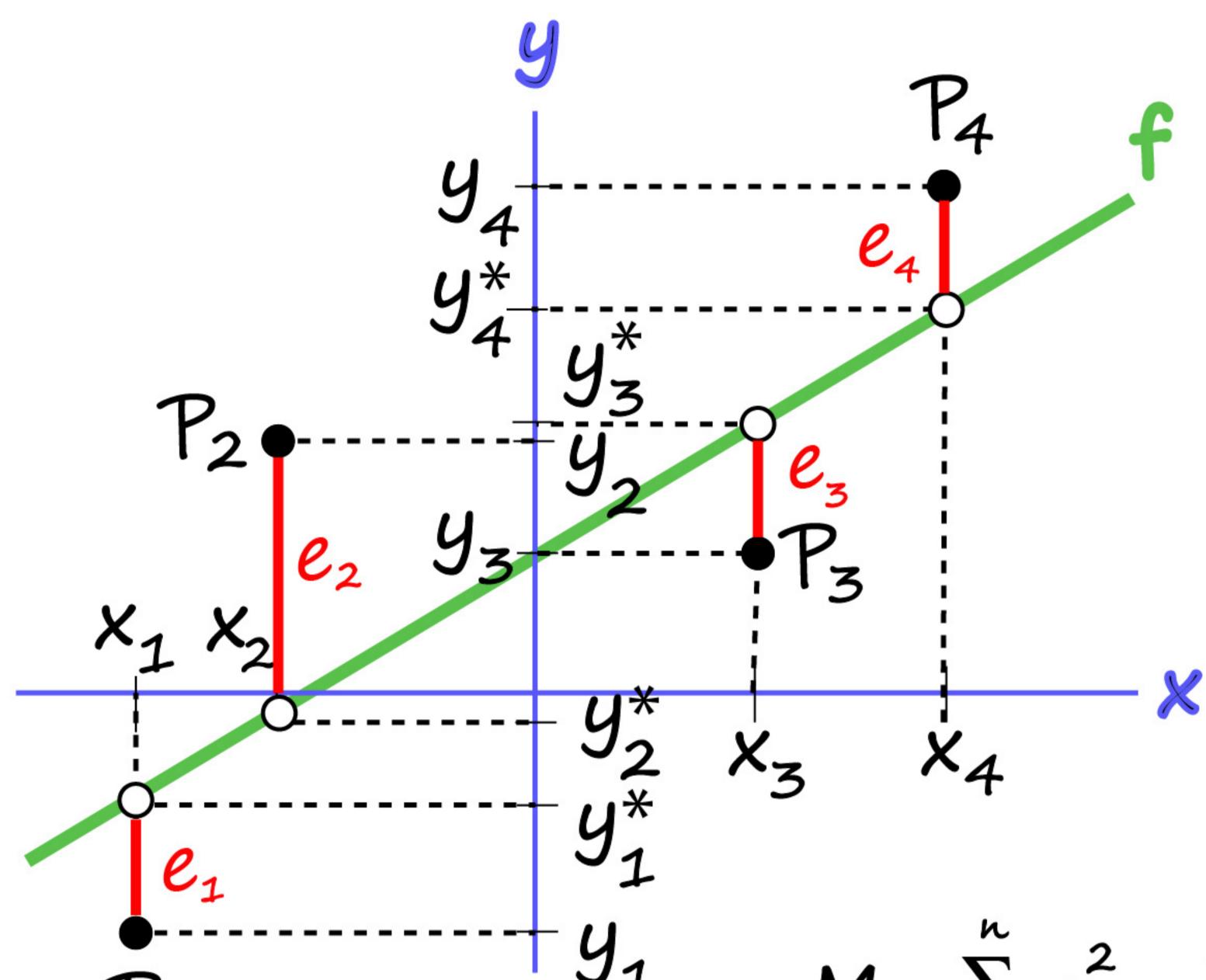
Para evitar este problema se minimizará la suma de los cuadrados de los errores, esto es:



$$M = \sum_{i=1}^{n} e_i^2$$

REGRESIÓN LINEAL

El caso particular en el que aproximamos nuestros datos utilizando una familia de funciones lineales, esto es, rectas, recibe el nombre de **Regresión Lineal**, y es de especial relevancia en el campo de la **Estadística**.



Tendríamos que las funciones de la familia son de la forma

$$f(x_i) = y_i^* = a + b x_i,$$

y por tanto los errores que nos aparecen son

$$e_i = a + b x_i - y_i$$

Siguiendo el método de los mínimos cuadrados, la expresión que hay que minimizar es

$$M = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (a + b x_i - y_i)^2$$

Esta expresión depende de dos parametros a y b, por tanto, para minimizarla derivaremos M respecto de a y b, e igualaremos a O. El sistema resultante tendrá como incógnitas a y b.

Por simplicidad en la notación los sumatorios los escribiremos simplemente como Σ .

$$\frac{\partial M}{\partial a} = \sum 2(a + b x_i - y_i) = 0$$

$$\frac{\partial M}{\partial b} = \sum 2(a + b x_i - y_i) x_i = 0$$

Sacando del sumatorio los factores que no tengan índices, operando y teniendo en cuenta que

$$\sum_{i=1}^{n} a = a n,$$

obtenemos el siguiente sistema de ecuaciones lineales:

$$an + b\sum x_i = y_i$$

$$a\sum x_i + b\sum x_i^2 = \sum x_i y_i$$

Este sistema recibe el nombre de Ecuaciones Normales.



El sistema de Ecuaciones Normales podemos resolverlo fácilmente por la regla de **Cramer**. En efecto, el determinante de la matriz de coeficientes es

$$\Delta = \begin{vmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{vmatrix} = n \sum x_i^2 - (\sum x_i)^2$$

Sustituyendo la primera y segunda columna de este determinante por la columna de términos independientes y cocientando por el determinante D obtenemos los valores de las incógnitas a y b respectivamente:

$$\mathbf{A} = \frac{1}{\Delta} \begin{vmatrix} \mathbf{\Sigma} \mathbf{y}_{i} & \mathbf{\Sigma} \mathbf{x}_{i} \\ \mathbf{\Sigma} \mathbf{x}_{i} \mathbf{y}_{i} & \mathbf{\Sigma} \mathbf{x}_{i} \end{vmatrix} = \frac{\mathbf{\Sigma} \mathbf{x}_{i}^{2} \mathbf{\Sigma} \mathbf{y}_{i} - \mathbf{\Sigma} \mathbf{x}_{i} \mathbf{\Sigma} \mathbf{x}_{i} \mathbf{y}_{i}}{\mathbf{n} \mathbf{\Sigma} \mathbf{x}_{i}^{2} - (\mathbf{\Sigma} \mathbf{x}_{i})^{2}}$$

$$b = \frac{1}{\Delta} \begin{vmatrix} n & \sum y_i \\ \sum x_i & \sum x_i y_i \end{vmatrix} = \frac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - (\sum x_i)^2}$$

Podemos simplificar considerablemente estas expresiones para a y b utilizando las nociones de media, varianza y covarianza.

Recordemos que si $m{x}_i$ son los datos correspondientes a una variable aleatoria $m{x}$, una medida de concentración de la variable está dada por la media

$$\overline{X} = \frac{\sum_{i=1}^{N} X_{i}}{N}$$
 donde N es el número de datos recogidos.

Dividiendo por N^2 en el numerador y denominador de la expresión obtenida para b nos queda:

$$b = \frac{\frac{1}{N} \sum_{i} x_{i} y_{i} - \overline{x} \overline{y}}{\frac{1}{N} \sum_{i} x_{i}^{2} - \overline{x}^{2}}$$



Si además tenemos en cuenta cómo estaban definidas las **medidas de dispersión**, todavía podemos escribir la expresión anterior de una forma más compacta.

La **varianza** de la variable aleatoria \mathbf{x} , mide cómo se desvían los datos de la variable respecto de la media $\overline{\mathbf{x}}$, y se define por la fórmula:

$$\sigma_{x}^{2} = \frac{\sum_{i=1}^{N} (x_{i} - \overline{x})^{2}}{N} = \frac{\sum_{i=1}^{N} x_{i}^{2}}{N} - \overline{x}^{2},$$

Una definición similar a la de la varianza pero mezclando la desviación de dos variables \mathbf{x} e \mathbf{y} es la **covarianza**, definida por la fórmula:

$$\sigma_{xy}^{2} = \frac{\sum_{i=1}^{n} (x_{i} - \overline{x})(y_{i} - \overline{y})}{n} = \frac{\sum_{i=1}^{n} x_{i} y_{i}}{n} - \overline{x} \overline{y}$$

vemos que la expresión para b puede escribirse en término de media, varianza y covarianza simplemente como:

$$\frac{\sigma_{xy}^{2}}{\sigma_{x}^{2}}$$

Finalmente, si dividimos entre **n** la primera ecuación del sistema de Ecuaciones Normales obtenemos:

$$a + b \overline{x} = \overline{y}$$

de donde obtenemos despejando, la fórmula para calcular a en términos de b y las medias $\overline{\mathbf{x}}$ e $\overline{\mathbf{y}}$:

$$a = \overline{y} - b \overline{x}$$

