

Nociones sobre cuadrados mínimos

por Dr. Horacio Bruzzone

1. Introducción

Supongamos que, a través de una serie de mediciones, se han determinado un conjunto de n pares de valores de dos magnitudes físicas, X e Y , asociadas a un cierto fenómeno. Es decir, se tiene el siguiente conjunto de datos experimentales:

$$\{x_i \pm \Delta x_i, y_i \pm \Delta y_i; i = 1 \text{ hasta } n\}$$

Supongamos además que hay motivos suficientes como para conjeturar que existe una relación funcional lineal entre X e Y (si la relación funcional es más compleja, lo que sigue puede aún aplicarse de varias maneras), es decir:

$$Y = a + bX$$

En este caso, cabe plantearse un doble problema: ¿me permiten mis datos afirmar que la relación lineal es correcta (al menos dentro del rango de valores medidos)? Y, supuesto que esto tiene respuesta afirmativa ¿cómo determino los parámetros a y b , y cuales son sus incertezas?

La técnica más usual para responder a esas preguntas, desarrollada originalmente por Gauss, se conoce como el método de ajuste por cuadrados mínimos, y se basa en lo que sigue.

Sean a y b los valores que se quieren encontrar. Para cada par de datos x_i, y_i , se define la cantidad: $d_i = y_i - (a + bx_i)$. Es intuitivo ver que, si la relación es efectivamente lineal, y los a y b son los correctos, cada d_i debe ser “pequeño” (luego precisaremos mejor esto). No podemos esperar que sean cero, porque hay incertezas de medición tanto en X como en Y , pero lo razonable es esperar que, si los valores de a y b elegidos son los que corresponden a la recta buscada, todos los d_i sean chicos. Como criterio que satisface esto, se puede proponer que la suma de los d_i sea mínima (en realidad, esto da el criterio para evaluar a y b). Sin embargo, la suma simple no es buen criterio, porque los d_i pueden ser positivos y negativos, y cancelarse mutuamente haciendo que la suma sea mínima aunque los d_i sean grandes en módulo. Para evitar este problema, se pueden sumar los módulos o, lo que es más conveniente, los d_i elevados al cuadrado. Entonces, la primera aproximación del método consiste en elegir los valores de a y b que minimizan a la suma:

$$\Sigma' = \sum_{i=1}^n d_i^2$$

2. Cuadrados mínimos ponderados

Antes de ver como se obtienen a y b usando ese proceso de minimización, vamos a dar un paso más para decidir que suma minimizar. Para ello, vamos a suponer (de momento arbitrariamente, luego se precisará esto) que los errores de la magnitud X son despreciables frente a los errores de Y . Nótese, de paso, que d_i es la distancia, en la dirección del eje y , entre y_i y el valor $y'_i (= a + bx_i)$ que es la coordenada y del punto de la recta buscada cuya abscisa es x_i . Es claro que, si los x_i “no tienen error”; es cierto que la relación es lineal y que sus parámetros son a y b , entonces las diferencias entre y_i e y'_i son únicamente atribuibles a los errores en los y_i . Pero si estos errores son distintos entre si, parece sensato realizar el procedimiento de

minimización asignando una mayor importancia a los d_i provenientes de valores de los y_i que tengan errores más chicos. Esto se logra por el procedimiento de “ponderar”, que consiste en multiplicar a cada d_i por una cantidad que sea mayor cuanto más pequeño sea el correspondiente Δy_i . La forma más sencilla de hacer esto es dividir a cada d_i por Δy_i : cuanto menor es el error, mayor es el coeficiente que tiene el d_i . Usando este criterio, lo que se debe minimizar es:

$$\Sigma = \sum_{i=1}^n \left(\frac{d_i}{\Delta y_i} \right)^2 = \sum_{i=1}^n \left(\frac{y_i - a - bx_i}{\Delta y_i} \right)^2$$

Al método de obtención de a y b usando la minimización de Σ' se lo llama “cuadrados mínimos ordinarios”, mientras que al que minimiza a Σ se lo llama “cuadrados mínimos ponderados”. Como la técnica de minimización es común, describiremos sólo la de cuadrados mínimos ponderados.

Una parte considerable del problema original que tenemos planteado consiste en responder a esta pregunta: dados n valores de x_i, y_i , y Δy_i ¿cuales son los valores de a y b que minimizan a Σ ? El análisis matemático enseña que son la solución del siguiente sistema de ecuaciones:

$$\frac{\partial \Sigma}{\partial a} = 0 \quad ; \quad \frac{\partial \Sigma}{\partial b} = 0$$

Haciendo esas derivadas se obtiene:

$$\begin{aligned} \frac{\partial \Sigma}{\partial a} &= \frac{\partial}{\partial a} \sum_i \left(\frac{y_i - a - bx_i}{\Delta y_i} \right)^2 = \sum_i \frac{\partial}{\partial a} \left(\frac{y_i - a - bx_i}{\Delta y_i} \right)^2 = -2 \sum_i \left(\frac{y_i - a - bx_i}{\Delta y_i^2} \right) = 0 \\ \frac{\partial \Sigma}{\partial b} &= \frac{\partial}{\partial b} \sum_i \left(\frac{y_i - a - bx_i}{\Delta y_i} \right)^2 = -2 \sum_i x_i \left(\frac{y_i - a - bx_i}{\Delta y_i^2} \right) = 0 \end{aligned}$$

Desarrollando las sumatorias, y definiendo (para simplificar notación):

$$\begin{aligned} S &= \sum \frac{1}{(\Delta y_i)^2} \quad ; \quad S_x = \sum \frac{x_i}{(\Delta y_i)^2} \quad ; \quad S_y = \sum \frac{y_i}{(\Delta y_i)^2} \\ S_{xx} &= \sum \frac{x_i^2}{(\Delta y_i)^2} \quad ; \quad S_{yy} = \sum \frac{y_i^2}{(\Delta y_i)^2} \quad ; \quad S_{xy} = \sum \frac{x_i y_i}{(\Delta y_i)^2} \end{aligned}$$

el sistema de ecuaciones se escribe como:

$$\begin{aligned} S_y - aS - bS_x &= 0 \\ S_{xy} - aS_x - bS_{xx} &= 0 \end{aligned}$$

cuya solución es:

$$b = \frac{SS_{xy} - S_x S_y}{SS_{xx} - S_x^2} \quad ; \quad (\Delta b)^2 = \frac{S}{SS_{xx} - S_x^2} \quad ; \quad a = \frac{S_{xx} S_y - S_x S_{xy}}{SS_{xx} - S_x^2} \quad ; \quad (\Delta a)^2 = \frac{S_{xx}}{SS_{xx} - S_x^2}$$

También se han agregado expresiones para las incertezas de a y b. No se dará aquí el detalle de su cálculo, basta comentar que se obtienen de las expresiones de a y b, propagando los errores Δy_i , que afectan a cada y_i que aparece en las sumatorias

$$((\Delta a)^2 = \sum_j \left(\frac{\partial a}{\partial y_j} \Delta y_j \right)^2 \text{ y otra equivalente para } \Delta b).$$

3. Confiabilidad

Queda todavía por contestar la pregunta de si los datos permiten confiablemente asumir que las variables X e Y guardan una relación lineal. Este tema está un poco confuso en la literatura, y es frecuente encontrar que el criterio para sostener la relación lineal es que un coeficiente, r^2 , llamado coeficiente de correlación de los datos y definido como sigue:

$$r^2 = \frac{(SS_{xy} - S_x S_y)^2}{(SS_{xx} - S_x^2)(SS_{yy} - S_y^2)}$$

sea “suficientemente cercano a 1”. Haciendo un poco de maniobras algebraicas, se puede demostrar que si $r^2=1$, entonces todos los puntos (x_i, y_i) están alineados en una recta. Obviamente, si se cumple esa condición, la relación lineal queda asegurada. Pero el problema es que no se puede decir nada cuando r^2 es menor que 1, en el sentido que no hay criterio sólido para saber que es “cercano a 1”.

Existe un criterio más elaborado, que es aplicable cuando los errores en Y son de naturaleza puramente estadística (y gaussiana). Su derivación es bastante complicada, y requiere el uso de teoría estadística, por lo que no se dará en este curso. Por otra parte, la médula de la argumentación empleada en ese criterio es extensible a casos no estadísticos, y la discutiremos a continuación.

Ya se mencionó que, gráficamente, d_i es la distancia en la dirección del eje y entre el dato i-ésimo y la recta obtenida usando el método de cuadrados mínimos (es decir, aquella determinada por los valores de a y b que se obtienen del método). Por lo tanto, $d_i/\Delta y_i$ es una “medida” de cuanto vale esa distancia comparada con el error en Y en ese punto. Supongamos que los errores en Y son de naturaleza instrumental. Es claro que, si ese cociente es menor o igual que 1, para ese dato particular, se puede decir que el punto pertenece (experimentalmente) a la recta. Si los errores fueran estadísticos gaussianos, esta condición no es suficiente, porque hay una probabilidad respetable (33%) de que el valor medio utilizado para y_i caiga fuera del intervalo definido por una dispersión, y en estos casos se debería agregar un intervalo de probabilidad en esa condición. Dejando de lado este problema, consideremos ahora que pasa con el conjunto de datos. En principio, es posible que algunos de los puntos tengan ese cociente menor que 1 mientras que otros lo tengan mayor que 1. Parece razonable usar como criterio para decidir si se acepta o no la relación lineal el que la mayoría de los puntos satisfaga esa condición. Y esto se puede hacer sencillamente calculando el valor que toma Σ , una vez calculados a y b. En efecto, Σ es la suma de los cuadrados de esos cocientes, de modo que Σ/n es el valor medio de ese cuadrado. Luego, $(\Sigma/n)^{1/2}$ es una estimación de cuanto vale el promedio de los $d_i/\Delta y_i$, y si resulta ser menor o igual que 1, la relación lineal debe ser aceptada; si es mucho mayor que 1, debe ser rechazada. Entre el valor 1 y “mucho mayor que 1” queda un margen de indecisión (entre otras cosas, en decidir cuanto es “mucho”). La ventaja de la teoría

estadística mencionada más arriba consiste en que, haciendo una comparación que es substancialmente equivalente, da resultados en términos de probabilidades para cualquier valor que se obtenga en la comparación. Con todo, el criterio expuesto aquí tiene la ventaja de ser sencillo de entender y válido para cualquier tipo de errores.

¿Que relación tiene esto con el r^2 ? Si se reemplazan los resultados de a y b dados más arriba en la expresión de Σ , y se usa la definición de r^2 , haciendo un poco de maniobras algebraicas se encuentra que:

$$\Sigma = (1 - r^2)(S_{yy} - S_y^2 / S)$$

en donde se ve que si r^2 vale 1, Σ es cero y por tanto, el criterio dado anteriormente se satisface. Sin embargo, el hecho que r^2 sea cercano a 1 (=0,95, por ejemplo) no alcanza para garantizar que Σ sea pequeño (menor o del orden de n) porque el otro factor $(S_{yy} - S_y^2/S)$ es habitualmente un número mucho mayor que 1.

Todavía no se dijo nada sobre n, el número de datos. Si los datos fueran sólo 2, es obvio que determinan un único valor de pendiente y ordenada al origen, y consecuentemente, $\Sigma=0$. Pero también es claro que, con sólo dos datos, no se puede pretender asegurar una relación lineal entre las variables X e Y. También es claro que, cuantos más datos se tengan, mayor confianza se puede tener en determinar la existencia (o inexistencia) de una relación lineal. En un tratamiento estadístico existen criterios adicionales sobre el número de datos requeridos, pero para lo que nos interesa acá basta decir que un mínimo razonable para n es $n = 10$. De todos modos, hay que agregar otra consideración: no basta con tener un valor de n razonable; es necesario además que los rangos de valores de X e Y cubiertos por los datos sean también razonables, al menos que cubran un orden de magnitud. Nótese que, en un rango de valores relativamente pequeño, cualquier función puede ser bien aproximada por una recta.

Una vez llegados a este punto, también puede verse cual es la diferencia esencial entre el método de cuadrados mínimos ponderados y el ordinario. El ponderado admite dar un criterio para aceptar o no la hipótesis lineal, el ordinario no, porque aún que se pueda definir en él un coeficiente equivalente al r^2 , hemos visto que esto no sirve para garantizar la hipótesis lineal. Es frecuente encontrar en calculadoras de mano y en paquetes de computación métodos de “correlación lineal” que, ante un conjunto de datos x_i, y_i dan valores de a, b, r^2 , e “incertezas” de a y b. Todos estos programas usan cuadrados mínimos ordinarios, y como no tienen incertezas en los datos Y, asumen que estas incertezas son constantes e iguales a $[\Sigma/(n-2)]^{1/2}$, usando este valor en las fórmulas de propagación. Pero, salvo el que dividen por n-2 en lugar de dividir por n (debido a razones estadísticas), lo que están usando como error en Y es la distancia promedio a la recta, y se pierde así toda posibilidad de dar un criterio. En efecto, este método supone que esos apartamiento de la recta son debidos a un error, en lugar de compararlos con los errores de medición. Este uso puede ser aceptable en algunas disciplinas (ciencias sociales, por ejemplo, en las que no es fácil ni habitual determinar errores), pero de ningún modo es aceptable en física, en donde cada dato debe conocerse con su error.

4. Errores en las dos variables

Acabamos de mencionar la importancia y conveniencia de incluir los errores de los datos en el método de cuadrados mínimos. Y nos queda pendiente la discusión de que criterio usar para decidir que una de las variables tiene errores despreciables,

que es lo que justifica el tratamiento anterior. El modo de aclarar este punto es replantear el problema para el caso general de errores en ambas variables, cosa que es relativamente nueva (D. York, Canadian Journal of Physics, **44** (1966) p. 1079, completado por J. Williamson, Canadian Journal of Physics, **46** (1968) p. 1845, y con trabajos hasta la fecha que aclaran y mejoran el método). A continuación se dará un resumen de las ideas en las que se basa.

Hay dos cosas que modificar en el tratamiento anterior: una es incluir en los factores de peso a los Δx_i , la otra es modificar consistentemente la “distancia” a la recta que se usa. Es claro que ya no se justifica minimizar la distancia en y (¿porque no en x ?), y la solución es minimizar la distancia perpendicular a la recta. Sea d_i' esa distancia (va a depender de la pendiente b , porque la pendiente de la recta perpendicular es $-1/b$ y también de los valores x_i e y_i), por analogía al caso anterior, se eligen los pesos como los $\Delta d_i'$ evaluados propagando los errores Δx_i y Δy_i (notar que con la definición anterior de d_i , Δy_i es su error si no hay error en x). Este cálculo es sencillo, y se obtiene la siguiente sumatoria a minimizar:

$$\Sigma_2 = \sum_i \frac{(y_i - a - bx_i)^2}{(\Delta y_i)^2 + (b\Delta x_i)^2}$$

El efecto final (engañoso) es el de minimizar la suma de los viejos d_i^2 pero divididos por coeficientes que incluyen el error en x . Nótese que ahora aparece un nuevo término que depende de b en el denominador, por lo que al derivar e igualar a cero, las ecuaciones serán más complejas. De hecho, en el caso general de errores todos distintos entre si, no hay solución analítica para a y b , sino que deben obtenerse soluciones en forma numérica. En cambio, las incertezas de a y b si tienen soluciones analíticas (pero muy arduas de calcular), supuesto que se conocen a y b .

En lo que sigue, nos limitaremos a un caso más sencillo, llamado de errores proporcionales ($\Delta x_i/\Delta y_i = \text{constante}$ para todo i). Esta limitación se adopta porque permite disponer de soluciones analíticas, y sin perder generalidad, permite entender que significa que uno de los errores sea despreciable.

Llamando d al cociente de errores podemos escribir Σ_2 en la forma:

$$\Sigma_2 = \frac{1}{1+b^2d^2} \sum \frac{(y_i - a - bx_i)^2}{(\Delta y_i)^2}$$

Siguiendo el mismo procedimiento que en el caso anterior, se forma un sistema de 2 ecuaciones para a y b , que pueden ser despejados y resulta una ecuación cuadrática para b , cuyas soluciones son:

$$db_{+,-} = \frac{1}{2} \left[-\frac{1}{db^*} + \frac{db^*}{r^2} \pm \sqrt{\frac{1}{(db^*)^2} + \frac{(db^*)^2}{r^4} + 4 - \frac{2}{r^2}} \right] ; \quad a_{+,-} = \frac{S_y}{S} - \frac{S_x}{S} b_{+,-}$$

en donde se ha llamado b^* a la solución ya conocida del caso ponderado y r^2 también es el del caso ponderado, lo mismo que las sumatorias “ S ” que aparecen en $a_{+,-}$. Para decidir cual signo tomar en la solución, hay que hacer un análisis de cual solución es un mínimo y cual un máximo. Se encuentra que esto depende del signo de b^* : si b^* es positivo, se debe tomar el signo $+$ de la raíz y viceversa.

Tomaremos en adelante la solución b_+ y, para simplificar escritura, en lo que sigue se usa esa raíz sin el subíndice. También se pueden calcular las incertezas de estas soluciones, y se obtiene:

$$\sigma_b^2 = \sigma_{b^*}^2 f(b, b^*, d, r^2) \quad ; \quad \sigma_a^2 = \sigma_{a^*}^2 (1 + b^2 d^2) f(b, b^*, d, r^2) + [1 - f(b, b^*, d, r^2)] \frac{1 + b^2 d^2}{S}$$

$$\text{donde } f(b, b^*, d, r^2) = (1 + b^2 d^2)^2 \frac{1 + b^{*2} d^2 / r^2}{[1 + d^2 (2bb^* - b^{*2} / r^2)]^2}$$

De estas expresiones se puede ver que, en el límite $d \rightarrow 0$, $b \rightarrow b^*$ (basta sacar factor común $1/(b^*d)^2$ en la raíz y desarrollar en serie); $a \rightarrow a^*$; y como $f \rightarrow 1$, las incertezas tienden a las clásicas. También se puede ver que si $r^2 \rightarrow 1$, $b \rightarrow b^*$ y $a \rightarrow a^*$ pero ahora las incertezas se modifican, resultando:

$$\sigma_{a,b}^2 = \sigma_{a^*,b^*}^2 (1 + d^2 b^{*2})$$

es decir, son mayores que las del caso clásico. Nótese que el incremento relativo de las incertezas es proporcional a $(db^*)^2 = (b^* \Delta x_i / \Delta y_i)^2$. Como las incertezas clásicas son, substancialmente, proporcionales a una especie de valor medio de $(\Delta y_i)^2$, lo que esto dice es que se debe agregar otro término que es substancialmente proporcional a b^* multiplicado por una especie de valor medio de $(\Delta x_i)^2$, lo cual es razonable. Por otro lado, esto pone en evidencia que criterio se debe usar para decir que un error es ignorable (o suficientemente pequeño) comparado con el otro: lo que cuenta no es que d sea chico, sino que db^* lo sea. Dicho de otra forma, el modo de asegurar a priori si los errores en alguna de las variables son despreciables (para errores proporcionales) consiste en estimar aproximadamente la pendiente de la recta y luego comparar los Δy con el producto de esa pendiente con los Δx . Si resulta Δy mucho mayor, los Δx son ignorables; si mucho menor, los Δy son ignorables, si son del mismo orden, entonces se deben incluir ambos.