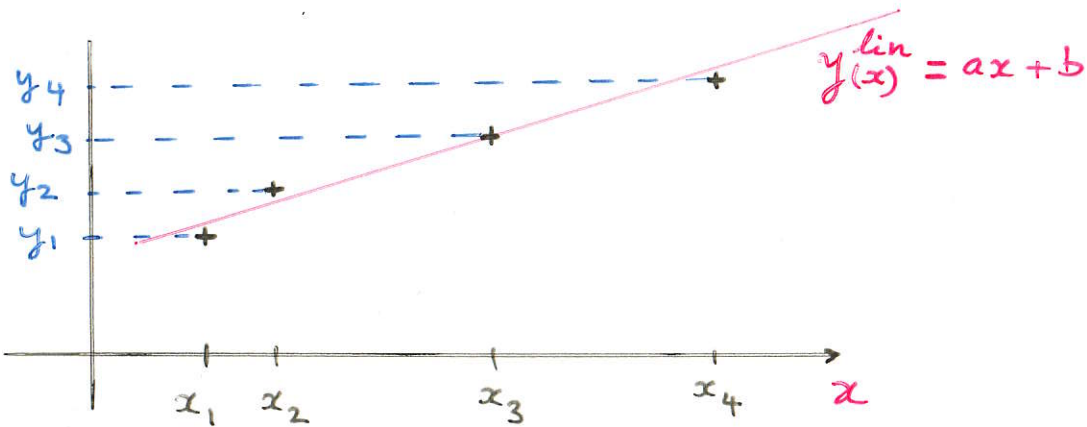


Régression linéaire



• Soient $\{x_i\}_{i=1, N}$ des nombres connus (par exemple les

concentrations de solutions étalons)

• Soient $\{y_i\}_{i=1, N}$ les valeurs mesurées pour chaque x_i avec une incertitude δ_i

• Comment trouver la "meilleure" régression linéaire, c'est-à-dire les coefficients a et b tels que la fonction $y^{\text{lin}}(x) = ax + b$ est le plus proche possible de y_i lorsque $x = x_i$?

Autrement dit, comment minimiser $\Delta y_i = y_i - y^{\text{lin}}(x_i)$ pour toutes les valeurs de i ?

• On définit la fonction d'erreur comme suit

$$E(a, b) = \sum_{i=1}^N \left(\frac{\Delta y_i}{\delta_i} \right)^2 = \sum_{i=1}^N \frac{(y_i - ax_i - b)^2}{\delta_i^2}$$

• Recherche du minimum: pour une fonction $f(x)$ d'une seule variable qui a un minimum en $x = x_0$, le développement limité autour de x_0 s'écrit

$$f(x) = f(x_0) + \left. \frac{df}{dx} \right|_{x=x_0} (x-x_0) + \frac{1}{2} \left. \frac{d^2f}{dx^2} \right|_{x=x_0} (x-x_0)^2 + \dots$$

si $\left. \frac{df}{dx} \right|_{x=x_0} > 0$, $f(x) < f(x_0)$ pour de très faibles valeurs négatives de $(x-x_0)$

si $\left. \frac{df}{dx} \right|_{x=x_0} < 0$, $f(x) < f(x_0)$ pour de très faibles valeurs positives de $(x-x_0)$

absurde!

donc $\left. \frac{df}{dx} \right|_{x=x_0} = 0$

De plus $\left. \frac{d^2f}{dx^2} \right|_{x=x_0} > 0$ (pour avoir un minimum)

• Les valeurs a_0 et b_0 qui minimisent $\xi(a, b)$ sont donc telles

que

$$\left. \frac{\partial \xi(a, b)}{\partial a} \right|_{a_0, b_0} = 0 = \left. \frac{\partial \xi(a, b)}{\partial b} \right|_{a_0, b_0}$$

$$\frac{\partial \xi(a, b)}{\partial a} = \sum_{i=1}^N \frac{1}{\delta_i^2} \times 2 \times (-x_i) (y_i - ax_i - b)$$

$$\frac{\partial \xi(a, b)}{\partial b} = \sum_{i=1}^N \frac{1}{\delta_i^2} \times 2 \times (-1) (y_i - ax_i - b)$$

En introduisant les notations

$$S = \sum_{i=1}^N \frac{1}{\delta_i^2}$$

$$S_y = \sum_{i=1}^N y_i / \delta_i^2$$

$$S_x = \sum_{i=1}^N x_i / \delta_i^2$$

$$S_{xx} = \sum_{i=1}^N x_i^2 / \delta_i^2$$

$$S_{xy} = \sum_{i=1}^N x_i y_i / \delta_i^2$$

on obtient

$$\frac{\partial \xi(a, b)}{\partial a} = 2 \left(-S_{xy} + a S_{xx} + b S_x \right)$$

$$\frac{\partial \xi(a, b)}{\partial b} = 2 \left(-S_y + a S_x + b S \right)$$

d'où

$$\begin{cases} -S_{xy} + a_0 S_{xx} + b_0 S_x = 0 & (1) \\ -S_y + a_0 S_x + b_0 S = 0 & (2) \end{cases}$$

2/RL

⇕

$$\begin{cases} -SS_{xy} + a_0 SS_{xx} + b_0 SS_x = 0 & (1) \times S \rightarrow (3) \\ -S_x S_y + a_0 S_x^2 + b_0 S_x S = 0 & (2) \times S_x \rightarrow (4) \end{cases}$$

⇓

$$\left\{ -SS_{xy} + S_x S_y + a_0 (SS_{xx} - S_x^2) = 0 \right. \quad (3) - (4)$$

Ainsi

$$a_0 = \frac{SS_{xy} - S_x S_y}{SS_{xx} - S_x^2}$$

De plus $S_x \times (1) - S_{xx} \times (2)$

↓

$$-S_x S_{xy} + S_x^2 b_0 + S_{xx} S_y - S_{xx} S b_0 = 0$$

d'où

$$b_0 = \frac{S_{xx} S_y - S_x S_{xy}}{SS_{xx} - S_x^2}$$

Incertitudes sur a et b :

soit $\delta a = a - a_0$ et $\delta b = b - b_0$.

$$\begin{aligned} \xi(a, b) &= \xi(a_0 + \delta a, b_0 + \delta b) = \sum_{i=1}^N \frac{(y_i - a_0 x_i - \delta a x_i - b_0 - \delta b)^2}{\delta_i^2} = \sum_{i=1}^N \frac{1}{\delta_i^2} \left[(y_i - a_0 x_i - b_0)^2 \right. \\ &\quad \left. + 2(y_i - a_0 x_i - b_0)(-\delta a x_i - \delta b) \right. \\ &\quad \left. + (\delta a x_i + \delta b)^2 \right] \\ &= \xi(a_0, b_0) + \underbrace{\delta a \frac{\partial \xi(a, b)}{\partial a}}_{a_0, b_0} + \underbrace{\delta b \frac{\partial \xi(a, b)}{\partial b}}_{a_0, b_0} \end{aligned}$$

$$+ \sum_{i=1}^N \frac{1}{\delta_i^2} (\delta a x_i + \delta b)^2$$

soit $\xi(a, b) = \xi(a_0, b_0) + \sum_{i=1}^N \frac{\delta a^2 x_i^2}{\delta_i^2} + 2 \frac{\delta a \delta b x_i}{\delta_i^2} + \frac{\delta b^2}{\delta_i^2}$

$$\xi(a, b) = \xi(a_0, b_0) + \delta a^2 S_{xx} + 2 \delta a \delta b S_x + \delta b^2 S$$

L'écart maximal autorisé sur a par rapport à a₀ s'obtient en écrivant $\xi(a, b) = \xi(a_0, b_0) + S \left[\delta b^2 + 2 \frac{S_x}{S} \delta a \delta b + \frac{S_{xx}}{S} \delta a^2 \right]$

$(\delta b + \frac{S_x}{S} \delta a)^2 - \frac{S_x^2}{S^2} \delta a^2$

$$\text{soit } \xi(a, b) = \xi(a_0, b_0) + S \left(\delta b + \frac{S_x}{S} \delta a \right)^2 + \left(S_{xx} - \frac{S_x^2}{S} \right) \delta a^2$$

En choisissant $\delta b = -\frac{S_x}{S} \delta a$ on réduit l'erreur mais il reste toujours le terme

$$\left(S_{xx} - \frac{S_x^2}{S} \right) \delta a^2$$

Remarque :

$$S_{xx} - \frac{S_x^2}{S} = \frac{S_{xx} S - S_x^2}{S} = \frac{-\Delta}{S}$$

où $\Delta = S_x^2 - S S_{xx}$

$$\begin{aligned} &= \left(\sum_{i=1}^N \frac{x_i}{\delta_i^2} \right) \left(\sum_{j=1}^N \frac{x_j}{\delta_j^2} \right) - \left(\sum_{i=1}^N \frac{1}{\delta_i^2} \right) \left(\sum_{j=1}^N \frac{x_j^2}{\delta_j^2} \right) \\ &= \sum_{i,j=1}^N \frac{1}{\delta_i^2 \delta_j^2} [x_i x_j - x_j^2] \\ &= \sum_{i,j=1}^N \frac{x_i x_j}{\delta_i^2 \delta_j^2} - \frac{1}{2} \sum_{i,j=1}^N \frac{x_j^2}{\delta_i^2 \delta_j^2} - \frac{1}{2} \sum_{i,j=1}^N \frac{x_i^2}{\delta_i^2 \delta_j^2} \end{aligned}$$

Soit $\Delta = \frac{1}{2} \sum_{i,j=1}^N \frac{1}{\delta_i^2 \delta_j^2} (2x_i x_j - x_i^2 - x_j^2)$

$$\Delta = -\frac{1}{2} \sum_{i,j=1}^N \frac{(x_i - x_j)^2}{\delta_i^2 \delta_j^2} < 0$$

• Ainsi l'erreur $(S_{xx} - \frac{S_x^2}{S}) \delta_a^2 = -\frac{\Delta}{S} \delta_a^2$ est bien positive et devient significative lorsque $-\frac{\Delta}{S} \delta_a^2 \gg 1$

La valeur de δ_a pour laquelle l'erreur reste acceptable sera donc définie comme

$$\delta_a = \sqrt{-\frac{S}{\Delta}} \leftarrow \text{interprète' comme l'incertitude sur la détermination de } a_0.$$

• Incertitude sur la détermination de b_0 :

$$\begin{aligned} \mathcal{E}(a,b) &= \mathcal{E}(a_0, b_0) + S_{xx} \left[\delta_a^2 + 2 \frac{\delta_a \delta_b S_x}{S_{xx}} \right] + S \delta_b^2 \\ &= \left(\delta_a + \frac{\delta_b S_x}{S_{xx}} \right)^2 + \delta_b^2 \left(S - \frac{S_x^2}{S_{xx}} \right) \end{aligned}$$

En choisissant $\delta_a = -\frac{\delta_b S_x}{S_{xx}}$ on minimise

l'erreur. Pour que cette dernière soit acceptable il faut que $\delta_b^2 \times \left(-\frac{\Delta}{S_{xx}} \right) \ll 1$

d'où $\delta_b = \sqrt{-\frac{S_{xx}}{\Delta}}$

\leftarrow interprète' comme l'incertitude sur la détermination de b_0 .

• Coefficient de régression linéaire :

Afin de rendre compte de la qualité de la régression linéaire on calcule le coefficient dit de régression linéaire

$$r = \sqrt{1 - \frac{\mathcal{E}(a_0, b_0)}{v_y}}$$

où $v_y = \sum_{i=1}^N \frac{(y_i - \bar{y})^2}{\delta_i^2}$

$$\bar{y} = \frac{1}{N} \sum_{i=1}^N y_i$$

$$\mathcal{E}(a_0, b_0) = \sum_{i=1}^N \frac{(y_i - a_0 x_i - b_0)^2}{\delta_i^2}$$

\bar{y} est la valeur moyenne des $\{y_i\}_{i=1, \dots, N}$

v_y indique l'écart des y_i à leur valeur moyenne.

Lorsque la régression est exacte, $\mathcal{E}(a_0, b_0) = 0 \Rightarrow r = 1$