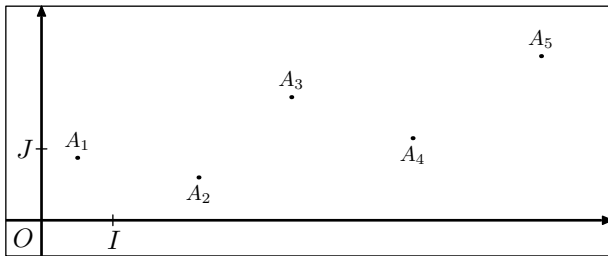


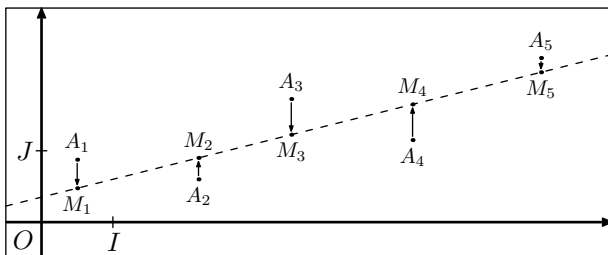
Régression linéaire

A. Introduction:

Soit $(x_i, y_i)_{0 \leq i \leq n}$ une série statistique dont le relevé graphique que les points associés à cette série statistique sont "presque alignés" autour d'une droite.



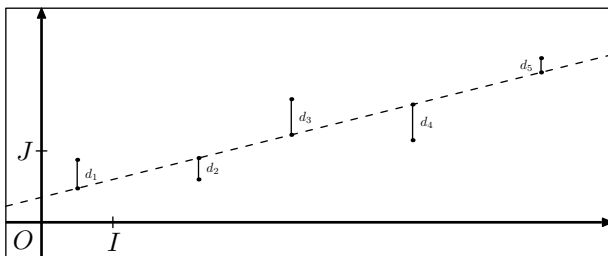
En modélisant ces données par la droite d'équation $y = a \cdot x + b$, on obtient les valeurs estimées par notre modèle qui sont les ordonnées des points M_i ci-dessous :



Pour mesurer l'erreur engendrée par cette représentation, on utilise la méthode des moindres carrés :

$$E(a, b) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \sum_{i=1}^n (y_i - a \cdot x_i - b)^2$$

Le choix de cette mesure de l'erreur entre les valeurs de la série statistique (x_i, y_i) et les valeurs prédites (x_i, \hat{y}_i) revient à mesurer la somme des carrés des variations sur l'axe des ordonnées.



On dit que la droite (d) réalise une **régression linéaire au sens des moindres carrés** lorsque les paramètres a et b minimise la valeur prise par la fonction E .

B. Quelques définitions:

Définition :

- Soit (X_i) une série statistique à une variable.
- ➔ la moyenne, notée \bar{X} , de la série (X_i) a pour valeur :

$$\bar{X} = \frac{1}{n} \cdot \sum_{i=1}^n X_i$$

- ➔ la variance, notée $\text{Var}(X)$, de la série (X_i) a pour valeur :

$$\text{Var}(X) = \frac{1}{n} \cdot \sum_{i=1}^n (X_i - \bar{X})^2$$

- Soit (X_i, Y_i) une série statistique à deux variables. La covariance de la série (X_i, Y_i) , notée $\text{Cov}(X, Y)$ a pour valeur :

$$\text{Cov}(X, Y) = \frac{1}{n} \cdot \sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})$$

C. Droite de régression linéaire:

Théorème : Soit (X_i, Y_i) une série statistique à deux variables.

Si la variance $\text{Var}(X)$ de la série statistique (X_i) est non-nulle, la droite de régression linéaire de y en x a pour équation $y = a \cdot x + b$ où :

$$a = \frac{\text{Cov}(X, Y)}{\text{Var}(X)} \quad ; \quad b = \bar{Y} - a \cdot \bar{X}$$

D. Démonstration:

Soit (x_i, y_i) une série à deux variables tels que $\text{Var}(X) \neq 0$. Considérons la fonction f définie par :

$$f : (a, b) \mapsto \sum_{i=1}^n (y_i - a \cdot x_i - b)^2$$

Proposition : la fonction f est strictement convexe

Preuve : la fonction f est une fonction quadratique (polynôme de degré 2). On en déduit qu'elle est continue et indéfiniment dérivable. On a les dérivées partielles :

- $\frac{\partial f}{\partial b}(a, b) = \sum_{i=1}^n -2x_i \cdot (y_i - a \cdot x_i - b)$

- $\frac{\partial f}{\partial b}(a, b) = \sum_{i=1}^n -2 \cdot (y_i - a \cdot x_i - b)$

- $\frac{\partial^2 f}{\partial a^2}(a, b) = \sum_{i=1}^n 2 \cdot x_i^2$

- $\frac{\partial^2 f}{\partial b^2}(a, b) = \sum_{i=1}^n 2 = 2 \cdot n$

- $\frac{\partial^2 f}{\partial a \partial b}(a, b) = \sum_{i=1}^n 2 \cdot x_i$

Nous utiliseront la propriété suivante :

f est strictement convexe

\Leftrightarrow la matrice Hessienne H de f est définie positive

Or, la fonction f admet pour matrice Hessienne :

$$H = \begin{pmatrix} \sum_{i=1}^n 2 \cdot x_i^2 & \sum_{i=1}^n 2 \cdot x_i \\ \sum_{i=1}^n 2 \cdot x_i & 2 \cdot n \end{pmatrix}$$

On a :

- $\text{Det}(H) = 2n \cdot \sum_{i=1}^n 2 \cdot x_i^2 - \left(\sum_{i=1}^n 2 \cdot x_i \right)^2$

$$= 4n^2 \cdot \left(\frac{1}{n} \cdot \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \cdot \sum_{i=1}^n x_i \right)^2 \right)$$

$$= 4n^2 \cdot \left[\frac{1}{n} \cdot \sum_{i=1}^n x_i^2 - \left(\frac{1}{n} \cdot \sum_{i=1}^n x_i \right)^2 \right]$$

$$= 4n^2 \cdot \text{Var}(X) > 0$$
- $\text{Tr}(H) = 2n + 2 \cdot \sum_{i=1}^n x_i^2 > 0$

Ces deux résultats montrent que la matrice hessienne de f est définie positive. On en déduit que la fonction f est strictement convexe.

Remarque : Voici une autre démonstration que la matrice H est définie positive.

En notant $X = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$, on a : $H = 2 \cdot X^t \cdot X$

Ainsi, pour tout $\begin{pmatrix} x \\ y \end{pmatrix}$, on a :

$$\begin{pmatrix} x \\ y \end{pmatrix}^t \cdot H \cdot \begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x \\ y \end{pmatrix}^t \cdot X^t \cdot X \cdot \begin{pmatrix} x \\ y \end{pmatrix}$$

$$= \left[\begin{pmatrix} x \\ y \end{pmatrix} \cdot X \right]^t \cdot \left[X \cdot \begin{pmatrix} x \\ y \end{pmatrix} \right] = M^t \cdot M \quad \text{où } M = X \cdot u$$

Pour déterminer les valeurs minimisant la fonction f , nous utiliserons le théorème suivant :

Théorème : condition suffisante pour un minimum global
Soit $f : \mathbb{R}^n \rightarrow \mathbb{R}$ une application différentiable et x_0 un point annulant la dérivée de f (*point critique*).

Si f est strictement convexe alors x_0 est l'unique minimum global de f .

Ainsi, le minimum global $(a; b)$ de la fonction f existe si il

vérifie : $\frac{\partial f}{\partial a}(a, b) = 0$; $\frac{\partial f}{\partial b}(a, b) = 0$

- $\frac{\partial f}{\partial b}(a, b) = 0$

$$\sum_{i=1}^n -2 \cdot (y_i - a \cdot x_i - b) = 0$$

$$-2 \cdot \sum_{i=1}^n (y_i - a \cdot x_i - b) = 0$$

$$\sum_{i=1}^n (y_i - a \cdot x_i - b) = 0$$

$$\sum_{i=1}^n y_i - a \cdot \sum_{i=1}^n x_i - n \cdot b = 0$$

$$n \cdot b = \sum_{i=1}^n y_i - a \cdot \sum_{i=1}^n x_i$$

$$b = \frac{1}{n} \cdot \sum_{i=1}^n y_i - \frac{a}{n} \cdot \sum_{i=1}^n x_i$$

$$b = -a \cdot \bar{x} + \bar{y}$$

- La fonction f peut s'exprimer par :

$$f(a, b) = \sum_{i=1}^n (y_i - a \cdot x_i - b)^2$$

$$f(a, b) = \sum_{i=1}^n [y_i - a \cdot x_i - (-a \cdot \bar{x} + \bar{y})]^2$$

$$f(a, b) = \sum_{i=1}^n (y_i - a \cdot x_i + a \cdot \bar{x} - \bar{y})^2$$

$$f(a, b) = \sum_{i=1}^n [(y_i - \bar{y}) - a \cdot (x_i - \bar{x})]^2$$

$$f(a, b) = \sum_{i=1}^n (y_i - \bar{y})^2 - 2 \cdot a \cdot (y_i - \bar{y}) \cdot (x_i - \bar{x}) + a^2 \cdot (x_i - \bar{x})^2$$

$$f(a, b) = \sum_{i=1}^n (y_i - \bar{y})^2 - 2 \cdot a \cdot \sum_{i=1}^n (y_i - \bar{y}) \cdot (x_i - \bar{x}) + a^2 \cdot \sum_{i=1}^n (x_i - \bar{x})^2$$

La dérivée partielle de la f valant 0 lorsque le minimum est réalisé vérifie :

$$\frac{\partial f}{\partial a}(a, b) = 0 - 2 \sum_{i=1}^n (y_i - \bar{y}) \cdot (x_i - \bar{x}) + 2 \cdot a \cdot \sum_{i=1}^n (x_i - \bar{x})^2$$

$$0 = 2 \cdot \left[a \cdot \sum_{i=1}^n (x_i - \bar{x})^2 - \sum_{i=1}^n (y_i - \bar{y}) \cdot (x_i - \bar{x}) \right]$$

$$a \cdot \sum_{i=1}^n (x_i - \bar{x})^2 = \sum_{i=1}^n (y_i - \bar{y}) \cdot (x_i - \bar{x})$$

$$a = \frac{\frac{1}{n} \cdot \sum_{i=1}^n (y_i - \bar{y}) \cdot (x_i - \bar{x})}{\frac{1}{n} \cdot \sum_{i=1}^n (x_i - \bar{x})^2}$$

$$a = \frac{\text{Cov}(x_i, y_i)}{\text{Var}(x_i)}$$