

Fiche de synthèse : Régression linéaire

Cette feuille regroupe l'essentiel de ce que nous avons vu sur les sujets MMI 2017 et 2019 dans la feuille 24 d'informatique.

Cette notion apparait dans plusieurs chapitres du programme de mathématiques :

- Cours de statistiques : Coefficients de la droite de régression linéaire d'une série double de données.
- Projection orthogonale : Calcul de la distance d'un vecteur u à un plan $\text{Vect}(v_1, v_2)$.
- Fonction de \mathbb{R}^2 dans \mathbb{R} : Calcul d'un point critique de $(a, b) \mapsto F(a, b)$.

Objectif :

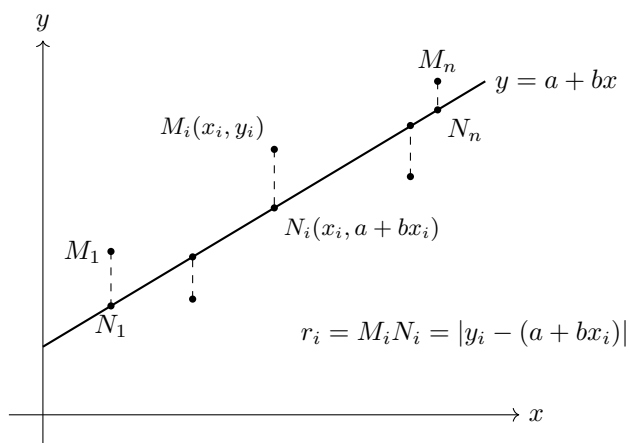
Déterminer le meilleur ajustement affine au sens des moindres carrés du nuage de points :

$$((x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)) \quad (\text{on suppose que les } x_i \text{ ne sont pas tous égaux})$$

Autrement dit :

on veut trouver la droite d'équation $y = a + bx$ qui minimise la somme $\sum_{i=1}^n r_i^2$ où les r_i sont définis sur

la figure ci-dessous :

**Notations usuelles :**

$$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i \quad \bar{y} = \frac{1}{n} \sum_{i=1}^n y_i \quad \overline{x^2} = \frac{1}{n} \sum_{i=1}^n x_i^2 \quad \overline{xy} = \frac{1}{n} \sum_{i=1}^n x_i y_i$$

$$\sigma_x^2 = \overline{x^2} - (\bar{x})^2 \quad (\text{variance empirique de } (x_i)) \quad \text{et} \quad \text{cov}(x, y) = \overline{xy} - \bar{x} \bar{y} \quad (\text{covariance empirique de } (x_i, y_i))$$

Remarque : les x_i ne sont pas tous égaux donc $\sigma_x^2 \neq 0$.

Résultat fondamental

Les coefficients de cette droite appelée droite de régression linéaire sont : $b = \frac{\text{cov}(x, y)}{\sigma_x^2}$ et $a = \bar{y} - b \bar{x}$

• Première approche.

On note $u = (y_1, \dots, y_n)$, $v_1 = (1, \dots, 1)$ et $v_2 = (x_1, \dots, x_n)$

(la famille (v_1, v_2) est une famille libre car les x_i ne sont pas tous égaux)

on note H le plan $\text{vect}(v_1, v_2)$.

Remarque : J'appelle ici ce plan H et non F pour éviter les confusions avec la fonction F .

On sait que la distance de u à H est le minimum de $\{\|u - av_1 - bv_2\| \mid (a, b) \in \mathbb{R}^2\}$ et ce minimum est atteint lorsque $av_1 + bv_2$ est le projeté orthogonal de u sur H .

$$\text{Or, pour tout } (a, b) \in \mathbb{R}^2, \|u - av_1 - bv_2\|^2 = \sum_{i=1}^n (y_i - (a + bx_i))^2 = \sum_{i=1}^n r_i^2$$

Ainsi, minimiser $\sum_{i=1}^n r_i^2$ revient à projeter u sur H . Autrement dit trouver (a, b) tel que : $p_H(u) = av_1 + bv_2$

$$\begin{aligned}
 av_1 + bv_2 = p_H(u) &\iff u - (av_1 + bv_2) \in H^\perp \\
 &\iff \begin{cases} u - (av_1 + bv_2) \perp v_1 \\ u - (av_1 + bv_2) \perp v_2 \end{cases} \\
 &\iff \begin{cases} \langle u - (av_1 + bv_2) | v_1 \rangle = 0 \\ \langle u - (av_1 + bv_2) | v_2 \rangle = 0 \end{cases} \\
 &\iff \begin{cases} \langle u | v_1 \rangle = a \langle v_1 | v_1 \rangle + b \langle v_2 | v_1 \rangle \\ \langle u | v_2 \rangle = a \langle v_1 | v_2 \rangle + b \langle v_2 | v_2 \rangle \end{cases} \\
 &\iff \begin{pmatrix} \langle v_1 | v_1 \rangle & \langle v_1 | v_2 \rangle \\ \langle v_2 | v_1 \rangle & \langle v_2 | v_2 \rangle \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \langle u | v_1 \rangle \\ \langle u | v_2 \rangle \end{pmatrix}
 \end{aligned}$$

on remarque que

$$\langle v_1 | v_1 \rangle = n \quad \langle v_1 | v_2 \rangle = n\bar{x} \quad \langle v_2 | v_2 \rangle = n\bar{x}^2 \quad \langle u | v_1 \rangle = n\bar{y} \quad \langle u | v_2 \rangle = n\bar{x}\bar{y}$$

on obtient le système :

$$\begin{aligned}
 \begin{pmatrix} n & n\bar{x} \\ n\bar{x} & n\bar{x}^2 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} n\bar{y} \\ n\bar{x}\bar{y} \end{pmatrix} &\iff \begin{pmatrix} 1 & \bar{x} \\ \bar{x} & \bar{x}^2 \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} \bar{y} \\ \bar{x}\bar{y} \end{pmatrix} \\
 &\iff \begin{cases} a + \bar{x}b = \bar{y} \\ \bar{x}a + \bar{x}^2b = \bar{x}\bar{y} \end{cases} \\
 &\iff \begin{cases} a + \bar{x}b = \bar{y} \\ (\bar{x}^2 - (\bar{x})^2)b = \bar{x}\bar{y} - \bar{x}\bar{y} \end{cases} \quad L_2 \leftarrow L_2 - \bar{x}L_1
 \end{aligned}$$

On obtient alors (car $\sigma_x^2 \neq 0$) :

$$av_1 + bv_2 = p_H(u) \iff \begin{cases} a = \bar{y} - b\bar{x} \\ b = \frac{\text{cov}(x, y)}{\sigma_x^2} \end{cases}$$

• Deuxième approche.

On définit $F : (a, b) \mapsto \sum_{i=1}^n (y_i - a - bx_i)^2$,

On remarque que $F(a, b) = \sum_{i=1}^n r_i^2$, on cherche donc a et b tels que $F(a, b)$ est minimal.

La première approche justifie que F admette un minimum donc on cherche le(s) point(s) critique(s) de F .

$F : (a, b) \mapsto \sum_{i=1}^n (y_i - a - bx_i)^2$ est polynomiale en a et en b donc elle admet des dérivées partielles.

$$\begin{aligned}
 \frac{\partial F}{\partial a}(a, b) &= \sum_{i=1}^n -2(y_i - a - bx_i) & \frac{\partial F}{\partial b}(a, b) &= \sum_{i=1}^n -2x_i(y_i - a - bx_i) \\
 &= -2 \left(\sum_{i=1}^n y_i - a \sum_{i=1}^n 1 - b \sum_{i=1}^n x_i \right) & &= -2 \left(\sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i - b \sum_{i=1}^n x_i^2 \right)
 \end{aligned}$$

donc

$$\begin{aligned}
 \begin{cases} \frac{\partial F}{\partial a}(a, b) = 0 \\ \frac{\partial F}{\partial b}(a, b) = 0 \end{cases} &\iff \begin{cases} n\bar{y} - na - bn\bar{x} = 0 \\ \sum_{i=1}^n x_i y_i - an\bar{x} - b \sum_{i=1}^n x_i^2 = 0 \end{cases} \\
 &\iff \begin{cases} a + \bar{x}b = \bar{y} \\ \bar{x}a + \bar{x}^2b = \bar{x}\bar{y} \end{cases} \\
 &\iff \begin{cases} a = \bar{y} - b\bar{x} \\ b = \frac{\text{cov}(x, y)}{\sigma_x^2} \end{cases}
 \end{aligned}$$