

DN9 : La régression linéaire

Leçons : 230, 231, 437 ; Référence : NES C

Notation : Si X est une variable aléatoire qui admet une espérance on note $\bar{X} = X - \mathbb{E}(X)$ la variable centrée.

Théorème 1 : On suppose que $\Omega = \{\omega_i | i \in \llbracket 1, n \rrbracket\}$ est fini de cardinal $n \geq 2$ et muni de la probabilité uniforme \mathbb{p} sur $\mathcal{P}(\Omega)$, soit $M = (X, Y)$ un couple de variables aléatoires non constantes définies sur Ω , on note $\forall i \in \llbracket 1, n \rrbracket, x_i = X(\omega_i), y_i = Y(\omega_i)$ et $M_i = M(\omega_i)$, on se place dans un repère orthonormé (O, \vec{i}, \vec{j}) et on note pour toute droite \mathcal{D} du repère $\phi_{\mathcal{D}} = \frac{1}{n} \sum_{i=1}^n (M_i H_i)^2$ ou H_i est le projeté de M_i sur \mathcal{D} parallèlement au vecteur \vec{j} , il existe une unique droite Δ dans le repère tel que $\phi_{\Delta} = \text{Min}\{\phi_{\mathcal{D}} | \mathcal{D} \text{ droite du plan}\}$ et alors $\phi_{\Delta} = \text{v}(Y)(1 - \text{Cor}(X, Y)^2)$ avec $\text{Cor}(X, Y) = \text{Cov}(X, Y) / (\sigma_X \sigma_Y)$. L'équation de la droite Δ est $t \rightarrow \text{Cov}(X, Y) / \text{v}(X) \cdot (t - \mathbb{E}(X)) + \mathbb{E}(Y)$.

• Si on note $\langle . | . \rangle$ l'application $(u, v) \rightarrow \frac{1}{n} \sum_{i=1}^n u_i v_i$ de \mathbb{R}^n vers \mathbb{R} alors on peut vérifier facilement qu'il s'agit d'un produit scalaire sur \mathbb{R}^n et on note $\| . \|$ la norme associée. Si on note $x = (x_i)_{1 \leq i \leq n}$ et $y = (y_i)_{1 \leq i \leq n}$ ainsi que $\mathcal{P} = \text{Vect}\{\hat{1}, x\}$ ou $\hat{1}$ désigne le vecteur de \mathbb{R}^n dont toutes les composantes sont égales à 1 alors on a $\{\phi_{\mathcal{D}} | \mathcal{D} \text{ droite du plan}\} = \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - ax_i - b)^2 \mid (a, b) \in \mathbb{R}^2 \right\} = \{ \|y - ax - b \cdot \hat{1}\|^2 \mid (a, b) \in \mathbb{R}^2 \} = \{ \|y - z\|^2 \mid z \in \mathcal{P} \}$ qui est non vide et minoré par 0.

• On a $\text{Inf}\{\phi_{\mathcal{D}} | \mathcal{D} \text{ droite du plan}\} = \text{Inf}\{\|y - z\|^2 \mid z \in \mathcal{P}\} = \text{Inf}\{\|y - z\| \mid z \in \mathcal{P}\}^2 = d(y, \mathcal{P})^2$ et d'après le théorème de projection sur un convexe complet on sait que $d(y, \mathcal{P}) = \|y - y'\|$ ou y' est le projeté orthogonal de y sur \mathcal{P} . Le vecteur y' est déterminé par les équations suivantes :

$$\begin{cases} y' \in \mathcal{P} \\ y - y' \perp \hat{1} \\ y - y' \perp x \end{cases} \Leftrightarrow \begin{cases} \exists (a, b) \in \mathbb{R}^2, y' = ax + b \cdot \hat{1} \\ \langle y - y' | \hat{1} \rangle = 0 \\ \langle y - y' | x \rangle = 0 \end{cases} \Leftrightarrow \begin{cases} \exists (a, b) \in \mathbb{R}^2, y' = ax + b \cdot \hat{1} \\ \langle y | \hat{1} \rangle = \langle y' | \hat{1} \rangle \\ \langle y | x \rangle = \langle y' | x \rangle \end{cases}$$

• Avec la deuxième équation du système on obtient $\mathbb{E}(Y) = \frac{1}{n} \sum_{i=1}^n y_i = \langle y | \hat{1} \rangle = \langle ax + b \cdot \hat{1} | \hat{1} \rangle = a \cdot \langle x | \hat{1} \rangle + \langle b \cdot \hat{1} | \hat{1} \rangle = a \cdot \mathbb{E}(X) + b$ donc $b = \mathbb{E}(Y) - a \cdot \mathbb{E}(X)$. La dernière équation du système donne $\mathbb{E}(YX) = \langle y | x \rangle = \langle ax + b \cdot \hat{1} | x \rangle = a \cdot \langle x | x \rangle + b \cdot \langle x | \hat{1} \rangle = a \cdot \mathbb{E}(X^2) + b \cdot \mathbb{E}(X) = a \cdot \mathbb{E}(X^2) + \mathbb{E}(Y) \cdot \mathbb{E}(X) - a \cdot \mathbb{E}(X)^2 = a \cdot \text{v}(X) + \mathbb{E}(Y) \cdot \mathbb{E}(X)$ donc $a = \text{Cov}(X, Y) / \text{v}(X)$. La droite Δ est déterminée et son équation est $t \rightarrow a \cdot t + b$.

• Puisque $y - y' \perp y'$ on a $\phi_{\Delta} = \text{Min}\{\phi_{\mathcal{D}} | \mathcal{D} \text{ droite du plan}\} = \|y - y'\|^2 = \langle y - y' | y - y' \rangle = \langle y - y' | y \rangle = \mathbb{E}(Y^2) - \langle y' | y \rangle = \mathbb{E}(Y^2) - \langle ax + b \cdot \hat{1} | y \rangle = \mathbb{E}(Y^2) - a \cdot \mathbb{E}(XY) - b \cdot \mathbb{E}(Y) = \mathbb{E}(Y^2) - a \cdot \mathbb{E}(XY) - (\mathbb{E}(Y) - a \cdot \mathbb{E}(X)) \cdot \mathbb{E}(Y) = \text{v}(Y) - a \cdot \text{Cov}(X, Y) = \text{v}(Y) \left(1 - \frac{\text{Cov}(X, Y)^2}{\text{v}(X) \text{v}(Y)} \right)$ de plus on remarque en notant $\bar{x} = (x_i - \mathbb{E}(X))_{1 \leq i \leq n}$ et $\bar{y} = (y_i - \mathbb{E}(Y))_{1 \leq i \leq n}$ que :

$$\text{Cor}(X, Y) = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y} = \frac{\mathbb{E}(\bar{X} \bar{Y})}{\sqrt{\mathbb{E}(\bar{X}^2) \mathbb{E}(\bar{Y}^2)}} = \frac{\langle \bar{x} | \bar{y} \rangle}{\|\bar{x}\| \|\bar{y}\|} = \text{Cos}(\bar{x}, \bar{y})$$

A partir d'un angle de 30° la corrélation linéaire est jugée bonne.

DEVELOPPEMENT SOUS LA FORME D'UN EXERCICE :

- 1) Déterminer un produit scalaire $\langle . | . \rangle$ sur \mathbb{R}^n tel que $\{\phi_{\mathcal{D}} | \mathcal{D} \text{ droite du plan}\} = \{ \|y - z\|^2 \mid z \in \mathcal{P} \}$ ou y est un vecteur de \mathbb{R}^n associé à Y et \mathcal{P} un plan vectoriel de \mathbb{R}^n .
- 2) Montrer que $\text{Min}\{\phi_{\mathcal{D}} | \mathcal{D} \text{ droite du plan}\} = \|y - y'\|^2$ ou y' est le projeté orthogonal de y sur \mathcal{P} .
- 3) Déterminer l'équation de la droite Δ et $\phi_{\Delta} = \text{Min}\{\phi_{\mathcal{D}} | \mathcal{D} \text{ droite du plan}\}$.
- 4) Montrer que $\text{Cor}(X, Y)$ est le cosinus d'un certain angle que l'on déterminera.