

La régression linéaire

Compréhension et pratique

Sommaire

I. Interpretation Géométrique.....	3
II. Cas d'erreurs gaussienne.....	4
III. L'utilisation du log	5
Références	6

I. Interpretation Géométrique

Dans cette partie, nous allons aborder la régression linéaire simple du point de vue géométrique. Pour cela nous allons commencer par transformer nos échantillons en vecteur.

Vu que nous avons n échantillons de la forme (x_i, y_i) , soient les vecteur suivant :

$$Y = [y_1, \dots, y_n] \text{ et } X = [x_1, \dots, x_n]$$

Or, pour tout i , nous avons supposé qu'il existe a, b deux réels vérifiant :

$$y_i = ax_i + b + \epsilon_i$$

Cette relation devient donc :

$$Y = aX + b\mathbb{I} + E \text{ avec } \begin{cases} \mathbb{I} = [1, \dots, 1] \\ E = [\epsilon_1, \dots, \epsilon_n] \end{cases}$$

Or, $E \cdot \mathbb{I} = \sum_{i=1}^n \epsilon_i = 0$ et $X \cdot \mathbb{I} = 0$ car, par hypothèse, X et \mathbb{I} sont indépendants.

De plus les vecteur \mathbb{I} et X ne sont pas colinéaires, donc il forment une base de \mathbb{R}^n de dimension 2. Soit $\tilde{Y} = aX + b\mathbb{I}$. \tilde{Y} est donc le projeté orthogonal de Y sur le plan formé par la base (X, \mathbb{I}) (vu que le vecteur E est orthogonal aux vecteurs \mathbb{I} et X).

D'où, par définition du projeté orthogonal, \tilde{Y} minimise la distance entre Y et le plan (X, \mathbb{I}) .

Ce qui veut dire :

$$d(Y, \tilde{Y}) = \min_{Z \in P} d(Y, Z) \text{ avec } P \text{ le plan formé par } (X, \mathbb{I})$$

Donc, on rejoint le résultat obtenu par la méthode des moindres carrés :

$$\|Y - \tilde{Y}\|^2 = \sum_{i=1}^n (y_i - (ax_i + b))^2 = \min_{c, d \in \mathbb{R}^2} \sum_{i=1}^n (y_i - (cx_i + d))^2$$

II. Cas d'erreurs gaussienne

Dans le cours précédent, des hypothèses, dites faibles, ont été faites sur les erreurs ϵ .

$$(\mathcal{H}) \begin{cases} E(\epsilon_i) = 0 \forall i \in \{1, \dots, n\} \\ Var(\epsilon_i) = \sigma^2 \forall i \in \{1, \dots, n\} \\ Cov(\epsilon_i, \epsilon_j) = 0 \text{ si } i \neq j \end{cases}$$

Afin d'obtenir de meilleurs résultats de prédiction, nous allons effectuer des hypothèses plus lourdes sur ces erreurs (qu'il faudrait vérifier ultérieurement bien sur) :

$$(\mathcal{H}') \begin{cases} \epsilon_i \sim N(0, \sigma^2) \forall i \in \{1, \dots, n\} \\ Cov(\epsilon_i, \epsilon_j) = 0 \text{ si } i \neq j \end{cases}$$

Cette hypothèse faite, la loi des y_i devient donc :

$$y_i \sim N(ax_i + b, \sigma^2)$$

Et comme les ϵ_i sont mutuellement indépendants donc les y_i le sont aussi.

Connaissant donc la loi régissant les ϵ_i et les y_i nous pouvons donc prédire plus sûrement en imposant des intervalles de confiance et en calculant de manière plus précise nos estimateurs.

(Pour plus de détails sur la suite de la procédure, voir le cours (1) mis en référence à la fin de ce document)

III. L'utilisation du log

Malheureusement, les cas exposés depuis le début de ce cours sont des cas peut fréquents lors de l'analyse de données réelles et s'approchent d'une situation idéale plus qu'une situation réelle (les relations linéaires entre la variables à prédire et la ou les variable(s) explicatives sont assez rares).

C'est pour cela que les transformations semi-logarithmiques ou bi-logarithmiques sont utilisés, afin de ramener des modèles potentiellement non-linéaires, ou bien qui ne vérifient pas les hypothèses faites sur les résidus, au modèle linéaire que nous venons d'étudier.

Les deux méthodes consistent à mener l'études linéaire en remplaçant les couple (X,Y) par :

$$(X, \ln(Y)) \text{ (transformation semi – logarithmique)}$$

$$(\ln(X), \ln(Y)) \text{ (transformation bi – logarithmique)}$$

Ces deux transformations permettent de linéariser (grâce aux propriété du log) des relations quadratiques, exponentielles, ... Et donc d'obtenir des meilleurs résultats de prédiction après application de la régression linéaire pour ces types de relations sans pour autant erroner les résultats si la relation entre X et Y est bel et bien linéaire.

Références

- (1) <http://www.lsta.upmc.fr/guyader/files/teaching/Regression.pdf>
- (2) http://grasland.script.univ-paris-diderot.fr/STAT98/stat98_7/stat98_7.htm