

Régression linéaire multiple

Définition

Soient p données statistiques $\mathbf{x}_k = (x_{1,k}, \dots, x_{n,k})$ et $\mathbf{y} = (y_1, \dots, y_n)$ des données statistiques.

Les variables x_k sont appelées **variables exogènes** ou *expliquées*.

La variable y est appelée **variable endogène** ou à *expliquer*.

Une **modélisation linéaire multiple** consiste à considérer les variables aléatoires

$$Y_i = \alpha_0 + \alpha_1 x_{i,1} + \alpha_2 x_{i,2} + \dots + \alpha_p x_{i,p} + \varepsilon_i$$

où es ε_i sont des variables aléatoires i.i.d. appelés **termes d'erreurs** et suivent une loi normale $\mathcal{N}(0, \sigma)$.

Proposition

Avec les notations de la définition précédente, un modèle linéaire multiple

$$Y_i = \alpha_0 + \alpha_1 x_{i,1} + \alpha_2 x_{i,2} + \dots + \alpha_p x_{i,p} + \varepsilon_i$$

est équivalent à

$$Y = \mathbf{x}\mathbf{m} + \mathbf{e}$$

Où \mathbf{x} est une matrice à n lignes et $p + 1$ colonnes, \mathbf{m} le vecteur à estimer de dimension $p + 1$ et \mathbf{e} est un vecteur gaussien de dimension n .

Théorème

Avec les notations précédentes, $\mathbf{e} = Y - \mathbf{x}\mathbf{m}$ et le minimum de ${}^t\mathbf{e}\mathbf{e}$ est atteint lorsque

$$\mathbf{m} = \hat{\mathbf{m}} \stackrel{\text{def}}{=} ({}^t\mathbf{x}\mathbf{x})^{-1}{}^t\mathbf{x}\mathbf{y}$$

Définition

Avec les notations précédentes, on appelle **résidu du modèle** le vecteur $\hat{\epsilon} = \mathbf{y} - \hat{\mathbf{y}}$ où $\hat{\mathbf{y}} = \mathbf{X}\hat{\mathbf{m}}$

Définition

Avec les notations précédentes on définit le **coefficient de détermination du modèle**, noté R^2 , par :

$$R^2_{x,y} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Proposition

Avec les notations précédentes

$$R^2_{x,y} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

En particulier $R^2_{x,y} \in [0, 1]$.

Proposition

Avec les notations précédentes, $\mathbb{V}(\mathbf{M}_n) = \sigma^2(\mathbf{t}_{\mathbf{X}\mathbf{X}})^{-1}$

Corollaire

Si la matrice $(\mathbf{t}_{\mathbf{X}\mathbf{X}})^{-1}$ tend vers la matrice nulle alors l'estimateur \mathbf{M}_n est convergent.

Théorème

Soit $\mathbf{e} = Y - \hat{y}$, alors

$$\mathbb{E}(\mathbf{e}) = (n - (p + 1))\sigma^2$$

Corollaire

La variable aléatoire

$$S_n = \frac{{}^t \mathbf{e} \mathbf{e}}{n - (p + 1)}$$

est un estimateur convergent et sans biais de σ^2 . De plus :

1. $S_n^M = S_n({}^t \mathbf{x} \mathbf{x})^{-1}$ est un estimateur convergent et sans biais de $\mathbb{V}(M_n)$.
2. En particulier, les coefficient diagonaux de S_n^M sont des estimateurs convergent et sans biais de la variance des coordonnées de M_n .
3. $(n - (p + 1)) \frac{S_n}{\sigma^2} \sim \chi^2(n - (p + 1))$

Corollaire

Soit $(M_n)_i$ la i -ième coordonnée de M_n et S_i^M le coefficient diagonale de la i -ième ligne de S_n^M . La variable aléatoire $\frac{M_{n,i} - m_i}{\sqrt{S_i^M}}$ suit une loi de Student à $n - (p + 1)$ degrés de libertés.

Corollaire

Soient $0 < \beta < \alpha < 1$, $t_1 = Q_{\mathcal{T}(n-(p+1))}(\beta)$ et $t_2 = Q_{\mathcal{T}(n-(p+1))}(1 - \alpha + \beta)$ alors

$$\left[(M_n)_i - t_2 \sqrt{S_i^M}; (M_n)_i - t_1 \sqrt{S_i^M} \right]$$

est un intervalle de confiance $1 - \alpha$ de m_i .