

LE MODÈLE LINÉAIRE

la régression linéaire simple

$$Y = \alpha + \beta X + \epsilon$$

on observe: $Y_i = \alpha + \beta X_i + \epsilon_i$

on estime par: $\hat{Y} = a + bX$

$$\hat{Y}_i = a + bX_i$$

$\sum \epsilon_i^2$ minimum

$$\sum \epsilon_i = 0$$

$$\sum \epsilon_i \cdot X_i = 0$$

$$\bar{X} = \frac{1}{n} \sum X_i$$

$$\text{Var } X = \frac{1}{n} \sum (X_i - \bar{X})^2 = E(X^2) - (E(X))^2$$

$$\text{Cov}(X, Y) = \frac{1}{n} \sum (X_i - \bar{X})(Y_i - \bar{Y})$$

$$= \frac{1}{n} \sum X_i Y_i - \bar{X} \bar{Y}$$

$$= \frac{1}{n} \sum X_i Y_i - \bar{X} \bar{Y}$$

Estimation des coefficients

$$b = \frac{\sum_{i=1}^n X_i Y_i - n \bar{X} \bar{Y}}{\sum_{i=1}^n X_i^2 - n \bar{X}^2}$$

$$\text{et } a = \bar{Y} - b \bar{X}$$

$$b = \frac{\sum_{i=1}^n (X_i - \bar{X})(Y_i - \bar{Y})}{\sum_{i=1}^n (X_i - \bar{X})^2}$$

$$b = \frac{\text{Cov}(X, Y)}{V(X)} = r \frac{\Delta Y}{\Delta X}$$

$$\text{avec } \Delta Y^2 = \frac{1}{n-1} \sum (Y_i - \bar{Y})^2$$

$$\Delta X^2 = \frac{1}{n-1} \sum (X_i - \bar{X})^2$$

Écart-types des estimateurs

$$\text{m mde } \Delta^2 = \frac{1}{n-2} \sum (Y_i - \hat{Y}_i)^2 = \frac{n-1}{n-2} (\Delta Y^2 - b^2 \Delta X^2)$$

(estimateur sans biais de σ^2)

$$= (1-r^2) \Delta Y^2$$

$$\text{écart-type de } b: SE_b = \frac{\Delta}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}} = \sqrt{\frac{\Delta^2 / \Delta X^2 - b^2}{n-2}}$$

$$\text{écart-type de } a: SE_a = \Delta \sqrt{\frac{1}{n} + \frac{\bar{X}^2}{\sum_{i=1}^n (X_i - \bar{X})^2}}$$

Intervalle de confiance de a et b

au niveau de confiance α au seuil de 5% par test

$$\text{IC de } \alpha: a \pm t_{0,025; n-2} \cdot SE_a$$

$$\text{IC de } \beta: b \pm t_{0,025; n-2} \cdot SE_b$$

loi de Student à $n-2$ ddl.

Prediction de Y pour X donné

IC de μ_0 la moyenne de Y en X_0 :

$$a + bX_0 \pm t_{0,025} \cdot \Delta \cdot \sqrt{\frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum (X_i - \bar{X})^2}}$$

$$\text{IC pour une simple observation } X_0: a + bX_0 \pm t_{0,025} \cdot \Delta \cdot \sqrt{1 + \frac{1}{n} + \frac{(X_0 - \bar{X})^2}{\sum (X_i - \bar{X})^2}}$$

Analyse de la variance

appliquée à la régression on définit les sommes des carrés des écarts

$$\text{Totale: } SCE_T = \sum_{i=1}^n (Y_i - \bar{Y})^2$$

$$\text{Expliquée (= régression): } SCE_E = \sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2$$

$$\text{Résiduelle (= ce qui n'est pas expliqué): } SCE_R = \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = (n-2) \Delta^2$$

$$\text{on démontre: } SCE_T = SCE_E + SCE_R$$

Source de variation	Somme des carrés	Degrés de liberté	Moyenne des carrés	Fisher
Expliquée	SCE_E	1	$Q_E = \frac{SCE_E}{1}$	$F = \frac{Q_E}{Q_R}$
Résiduelle	SCE_R	$n-2$	$Q_R = \frac{SCE_R}{n-2}$ $Q_R = \Delta^2$	
Totale	SCE_T	$n-1$		

Sous l'hypothèse $\beta=0$, F suit une Fisher (1; $n-2$).

Coefficient de corrélation

$R \in [-1, 1]$
R a même signe que b

$$R^2 = \frac{SCE_E}{SCE_T} = b^2 \frac{\Delta X^2}{\Delta Y^2}$$

$$r = \beta \frac{\sigma_X}{\sigma_Y}$$

on définit

$$t_c = \frac{R \sqrt{n-2}}{\sqrt{1-R^2}}$$

$$r_{theo} = r = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X) \text{Var}(Y)}}$$

on a $F = t_c^2$

Sous l'hypothèse $\beta=0$, t_c suit une loi de Student à $(n-2)$ ddl.

Analyse de la variance à 1 facteur

Un facteur A a a modalités: $i: 1 \rightarrow a$

Pour chaque modalité on a m_i observations. $\sum m_i = n$

Source de variation	Somme des carrés des écarts	ddl	Carré moyen	F
Facteur A	$SCE_A = \sum_{i=1}^a m_i (\bar{X}_i - \bar{X})^2$	$a-1$	$Q_A = \frac{SCE_A}{a-1}$	$F = \frac{Q_A}{Q_R}$
Résiduelle	$SCE_R = \sum_{i=1}^a \sum_{k=1}^{m_i} (X_{ik} - \bar{X}_i)^2$	$\sum_{i=1}^a (m_i - 1) = n - a$	$Q_R = \frac{SCE_R}{n-a}$	
Totale	$SCE_T = \sum_{i=1}^a \sum_{k=1}^{m_i} (X_{ik} - \bar{X})^2$	$n-1$		

Sous $H_0 (\mu_1 = \mu_2 = \dots = \mu_a)$, F suit une loi de Fisher à $(a-1; n-a)$ ddl

$$\text{On a } SCE_T = SCE_A + SCE_R$$