

Régression linéaire simple (continuation)

(1)

$$Y = \beta_0 + \beta_1 X + \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, \sigma^2)$$

Echantillon aléatoire $\{(x_i, y_i), i=1, \dots, n\}$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

$$\hat{\beta}_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = \hat{\rho}(x, y) \cdot \frac{s_y}{s_x}$$

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum (y_i - \hat{\beta}_0 - \hat{\beta}_1 x_i)^2 = \frac{1}{n-2} \sum_{i=1}^n \varepsilon_i^2$$

Proposition: $\hat{\beta}_0$, $\hat{\beta}_1$ et $\hat{\sigma}^2$ sont des estimateurs sans biais de β_0 , β_1 et σ^2 . de plus

$$V(\hat{\beta}_0) = \frac{\sigma^2 \sum x_i^2}{n \sum_{i=1}^n (x_i - \bar{x})^2}; \quad V(\hat{\beta}_1) = \sigma^2 \cdot \frac{1}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

$$\text{Cor}(\hat{\beta}_0, \hat{\beta}_1) = \frac{\sigma^2 \bar{x}}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Théorème [Gauss-Markov]: $\hat{\beta}_0$, $\hat{\beta}_1$, $\hat{\sigma}^2$ sont des estimateurs à variance minimale parmi tous les estimateurs linéaires en $\{y_i\}$.

Les résidus :

$$\hat{\varepsilon}_i = y_i - \hat{y}_i$$

Propriété : $\sum_{i=1}^n \hat{\varepsilon}_i = 0$

Décomposition de la variance de y :
ANOVA de la régression

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\text{variance totale}} = \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\text{résiduelle}} + \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\substack{\text{expliquée par} \\ \text{le modèle}}}$$

De plus

$$\frac{\text{variance expliquée par le modèle}}{\text{variance totale}} = R^2$$

R^2 = coefficient de détermination
(parfois noté avec R^2).

• Inference dans le modèle linéaire simple

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{V(\hat{\beta}_1)}} \sim T_{n-2}, \quad \frac{\hat{\beta}_0 - \beta_0}{\sqrt{V(\hat{\beta}_0)}} \sim T_{n-2}$$

$$\boxed{H_0: \beta_1 = 0} \Leftrightarrow \boxed{H_0: \exists \text{ lien linéaire entre } X \text{ et } Y?}$$

$$\Leftrightarrow \boxed{H_0: \rho(x, y) = 0}$$

Pour tester la signification du modèle linéaire on teste

$$\underline{H_0: \rho^2(x, y) = 0}$$

à l'aide de la statistique

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

TEST DE FISHER de la R.L.

Diagnostic de la régression linéaire simple (4)

1. Analyse des résidus : $\hat{\varepsilon}_i = y_i - \hat{y}_i$

- on fait deux hypothèses :

- résidus normaux ($\varepsilon \sim N(0, \sigma^2)$)
- indépendants et identiquement distribués (homoscédasticité)

On peut montrer que, même si l'hypothèse de homoscédasticité est vérifiée, on a :

$$\boxed{V(\hat{\varepsilon}_i) = \sigma^2(1 - h_{ii})} \quad \text{où}$$

h_{ij} sont les éléments de la matrice H telle que

$$\hat{y} = H \cdot y \quad (\text{vectoriel})$$

$$h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{\sum_i (x_i - \bar{x})^2}$$

Résidus standardisés (afin de les rendre comparables)

$$res_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma} \cdot \sqrt{1 - h_{ii}}}$$

Résidus studentisés

$$t_i = \frac{\hat{\varepsilon}_i}{s_{(i)} \sqrt{1 - h_{ii}}}$$

où $s_{(i)}^2$ = estimation de
la variance
résiduelle par
validation croisée

Remarque

Validation croisée et PRESS

Notons avec $\hat{y}_{i(i)}$ la prédiction de y_i sans
l'observation (x_i, y_i) dans l'échantillon.

Alors on peut montrer que

$$\hat{\varepsilon}_{i(i)} = y_i - \hat{y}_{i(i)} = \frac{\hat{\varepsilon}_i}{1 - h_{ii}}$$

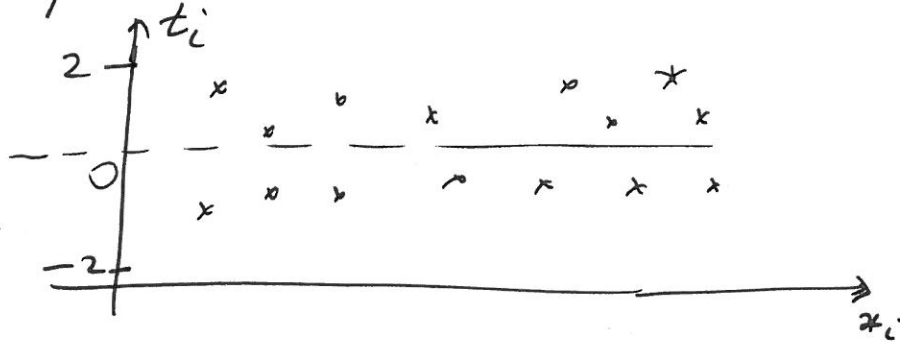
$$\text{PRESS} = \frac{1}{n} \sum \hat{\varepsilon}_{i(i)}^2$$

PRESS = predicted residual sum of squares
= leave-one-out cross-validation (validation croisée)
= mesure le coût prédictif d'un modèle
de prédiction.

1.1 Normalité des résidus : [graphique]

(6)

Si les résidus sont $N(0, \sigma^2)$ alors les résidus t_i (studentisés) sont distribués selon une loi de Student à $n-3$ ddl.
En pratique on vérifie si les $t_i \in [-2; 2]$.



On peut aussi tester la normalité des résidus à l'aide d'un test de

Kolmogorov-Smirnov : $\|F_n - \Phi\|_\infty$

Shapiro-Wilk
$$W = \frac{\left[\sum_{i=1}^{\lfloor n/2 \rfloor} a_i (x_{(n-i+1)} - x_{(i)}) \right]^2}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

Droite d'Henry ou Q-Q plot

(quantiles théoriques versus quantiles estimés)

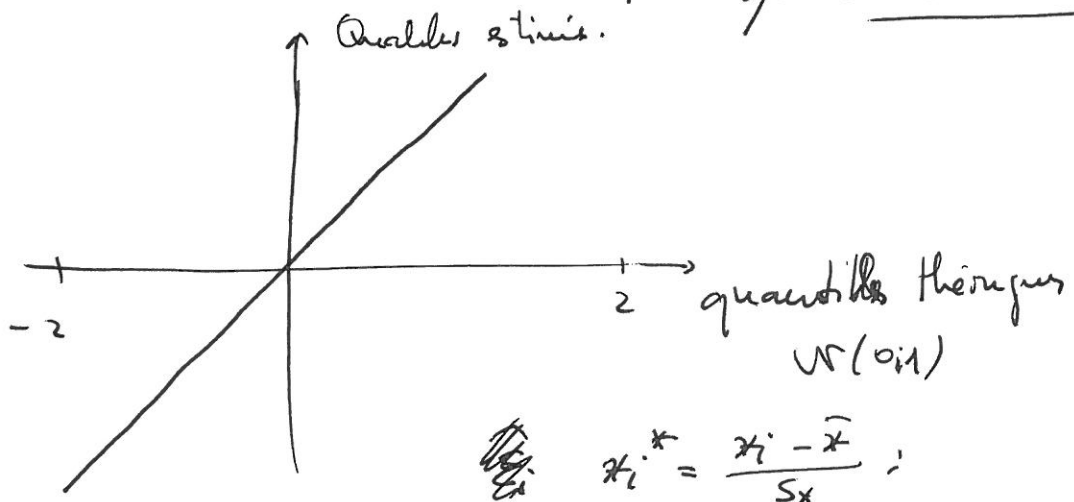
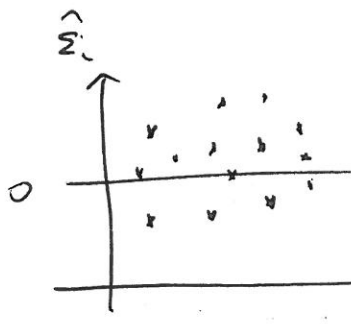
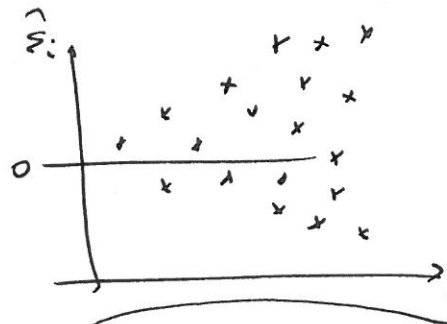


illustration en R : ks.test, shapiro.test, qqnorm + qqline

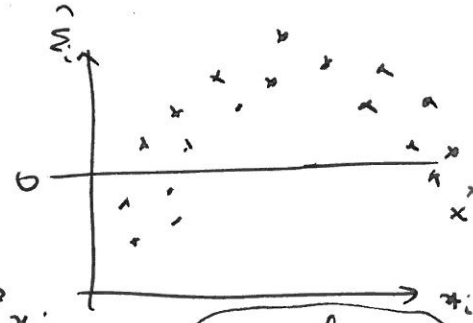
1.2 Homoscédasticité des résidus [analyse graphique] ⁽⁷⁾ (variance constante).



(OK)



variance
non-constante

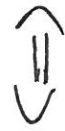


non-linéar

Solution :

transformer la variable y

$y \rightarrow \sqrt{y}$
 $y \rightarrow \log(y)$



obtenir des modèles
non-linéaires.

Pas de tests spécifiques pour l'homoscédasticité à
part le test de Breusch-Pagan : bptest dans
le package "lmtest".

SAS : PROC MODEL

1.3 Indépendance des résidus

8

Tests d'indépendance des résidus :

Autocorrélation des résidus :

$$\underline{\varepsilon_i = \rho \varepsilon_{i-1} + \eta_i}$$

H_0 : résidus non-correlés $\Leftrightarrow \rho = 0$

Test de Durbin - Watson

$$DW = \frac{\sum_{i=2}^n (\hat{\varepsilon}_i - \hat{\varepsilon}_{i-1})^2}{\sum_{i=1}^n \hat{\varepsilon}_i^2} \sim \chi^2$$

en R : `dwtest` ("lmtest" package)

en SAS : `proc Autoreg`

↑ possible tests autocorrélations
d'ordre plus grands
option : dw

... Diagnostic de la régression linéaire simple (9)
(continuation)

② Influence des observations

$$\hat{y} = H \cdot y$$

$$\begin{pmatrix} \hat{y}_1 \\ \hat{y}_i \\ \vdots \\ \hat{y}_n \end{pmatrix} = \begin{bmatrix} h_{i1} & \dots & h_{ii} & \dots & h_{in} \end{bmatrix} \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix}$$

$$\text{avec } h_{ij} = \frac{1}{n} + \frac{(x_i - \bar{x})(x_j - \bar{x})}{\sum_{i=1}^n (x_i - \bar{x})^2}$$

h_{ii} mesure l'impact de l'observation y_i dans l'estimation de \hat{y}_i

si h_{ii} est grande alors y_i influence "beaucoup" sur \hat{y}_i

Un point est un point de levier si $h_{ii} > \frac{2}{n}$ (Hoaglin)

$$> \frac{3}{n} \text{ (Weber)}$$

$$> 0.5 \text{ (Huber)}$$

Distance de Cook :

$$D_i = \frac{\sum_{j=1}^n (\hat{y}_{(i)j} - \hat{y}_j)^2}{2\Delta^2}$$

$$\Delta^2 = \frac{1}{n-2} \sum (y_i - \hat{y}_i)^2$$

$D_i > 1$ - point influent.

Régression linéaire multiple ($p \geq 2$)

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p + \varepsilon$$

Hypotheses : $\varepsilon \sim \mathcal{N}(0, \sigma^2)$.

$$\varepsilon \perp \{X_1, X_2, \dots, X_p\}$$

Données

ind	X_1	X_2		X_p	Y
1	x_{11}	x_{12}		x_{1p}	y_1
\vdots					
i	x_{i1}	x_{i2}	\dots	x_{ip}	y_i
\vdots					
n	x_{n1}	x_{n2}		x_{np}	y_n

$$X = \begin{pmatrix} 1 & x_{11} & x_{12} & \dots & x_{1p} \\ 1 & x_{21} & x_{22} & & x_{2p} \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & x_{n1} & x_{n2} & & x_{np} \end{pmatrix} \in \mathcal{M}_{n \times (p+1)}$$

$$\boxed{Y = X\beta + \varepsilon}$$

$$\text{avec } \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix}, \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

La solution qui minimise les moindres carrés (MCO)

$$\sum_{i=1}^n \left(Y_i - (\beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}) \right)^2 = \sum_{i=1}^n \varepsilon_i^2$$

vaut

$$\hat{\beta} = \underbrace{(X^T X)^{-1} X^T}_{\text{et}}$$

$$\hat{Y} = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_i \\ \vdots \\ \hat{y}_n \end{pmatrix} = X \hat{\beta} = \underbrace{X (X^T X)^{-1} X^T}_H Y$$

$$\hat{Y} = H Y \quad \text{avec}$$

$$H = X (X^T X)^{-1} X^T$$

("hat" matrix).

Remarque : Existence de $(X^T X)^{-1}$:

$$\text{cas : } \left\{ \begin{array}{l} - n < p \\ - \sum_{i=1}^p \alpha_i X_i = \alpha_0 \quad (\text{multicolinéarité}) \end{array} \right.$$

Proposition

1) $\hat{\beta}$ est un estimateur sans biais pour β

$$2) \hat{\sigma}^2 = \frac{1}{n-p-1} \sum_{i=1}^n \hat{\varepsilon}_i^2, \quad \hat{\varepsilon}_i = Y_i - \hat{Y}_i$$

est un estimateur sans biais pour σ^2

$$3) V(\hat{\beta}) = \sigma^2 (X^T X)^{-1}$$

Théorème [Gauss-Markov]

$\hat{\beta}$ est optimal parmi tous les estimateurs linéaires en Y sans biais.

$$\widehat{V(\hat{\beta})} = \hat{\sigma}^2 (X^T X)^{-1}$$

$$\widehat{V(\hat{\beta}_j)} = \hat{\sigma}^2 (X^T X)^{-1}_{[j,j]}$$

$$T_j = \frac{\hat{\beta}_j - \beta_j}{\sqrt{\widehat{V(\hat{\beta}_j)}}} \sim \text{Student}(n-p-1)$$

Tests de $H_0: \beta_j = 0$ vs $\beta_j \neq 0$

Qualité du modèle : R^2

ANOVA de la régression :

$$\underbrace{\sum_{i=1}^n (Y_i - \bar{Y})^2}_{\text{var totale}} = \underbrace{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}_{\text{var résiduelle}} + \underbrace{\sum_{i=1}^n (\hat{Y}_i - \bar{Y})^2}_{\text{var modèle}}$$

$$R^2 = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} = \rho^2(Y, \hat{Y})$$

Est-ce que le modèle est significatif : $\beta \neq 0$?

$$\begin{cases} H_0: \beta_1 = \beta_2 = \dots = \beta_p = 0 \\ H_1: \exists \beta_j \neq 0, j=1 \dots p \end{cases}$$

La statistique de Fisher :

$$F = \frac{R^2}{1 - R^2} (n - p - 1) \underset{H_0}{\sim} F(p, n - p - 1)$$

Diagnostiques de la régression linéaire multiple

14

① Analyse des résidus

- homoscedasticité
- normalité
- indépendance

Parait que dans le cas de la régression linéaire simple

② Influence des observations : Distance de Cook

③ Influence des variables

- Est-ce que la variable X_j contribue à la qualité du modèle ?

$$H_0: \beta_j = 0$$

- Test de Student

- Test de Fisher $F = \frac{\|\hat{Y}_j - \hat{Y}\|}{\|Y - \hat{Y}\|^2 / (n-p-1)} \sim F(1, n-p-1)$

Choix des variables \equiv Choix du modèle

Intèrès ?

Objectifs

- description
- estimation
- prévision