

Statistique Cours 9

Modèles linéaires

$$y = a_1 x_1 + a_2 x_2 + \dots + a_m x_m + b$$

← terme pondéré
 $a_i \in \mathbb{R}$
 $b_i \in \mathbb{R}$

Inconnues $a_i, i \in 1, \dots, m$

↓
Estimation $\hat{a}, \hat{b} \Rightarrow$ Régression multi-linéaire

$$\begin{bmatrix} y_1 \\ \vdots \\ y_m \end{bmatrix} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1m} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & \dots & \dots & x_{mm} \end{bmatrix} \begin{bmatrix} a_1 \\ a_2 \\ \vdots \\ a_m \\ b \end{bmatrix} \quad y = Xa$$

hyp: $m \gg n$
2° $x_1, \dots, x_m \forall a_i$

variables explicatives

3° y dépend linéairement des x_1, \dots, x_m
↳ variable à expliquer

$$y = Xa$$
$$\hat{a} = X^+ y$$

où $X^+ = (X^T X)^{-1} X^T$

$$\hat{y} = X \hat{a}$$

Analyse des résidus

$$\hat{y} = \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_m \end{bmatrix}, \quad e_i = y_i - \hat{y}_i$$

erreur de l'estimation

$$4^\circ E(e) = 0$$

$$V = E[(X - E(X))(X - E(X))^T]$$
$$= \begin{bmatrix} V_{xx} & V_{xy}^T \\ V_{xy} & V_{yy} \end{bmatrix} \quad \begin{array}{l} \text{matrice des variance} \\ \text{- covariances.} \end{array}$$

$V_{xx}(i, i) \Rightarrow$ variance de x_i avec x_j

$V_{yy} \Rightarrow$ ————— de y

$V_{xy}(i, j) \Rightarrow$ covariance entre x_i et x_j .

$$\text{Corrélation: } R^2 = \frac{V_{xy}^T V_{xx}^{-1} V_{xy}}{V_{yy}}$$

Variance de régression $S^2 = \frac{m}{m-m-1} V_{yy} (1-R^2)$

Tests de signification

Si $\frac{m-m-1}{m} \cdot \frac{R^2}{1-R^2} \geq F_{\alpha, m, m-m-1} \Rightarrow$ il y a au moins un coef parmi les a_i qui est significatif ($a_i \neq 0$)

$$\text{Id}(\hat{a}_i) = [\hat{a}_i - t_{\alpha, m-m-1} \sqrt{V(\hat{a}_i)}, \hat{a}_i + t_{\alpha, m-m-1} \sqrt{V(\hat{a}_i)}]$$

$$\text{avec } V(\hat{a}_i) = S = (X^T X)^{-1}$$

$x_1 \dots x_q$ X_q sous matrice de X créée à partir de x_1, x_2, \dots, x_q .

Linéarisation des données

Fonct
 $y = \alpha x^\beta$

Transf.
 $y' = \log y$
 $x' = \log x$

Fonct
 $y' = \log \alpha + \beta x'$

$y = \alpha e^{\beta x}$

$y' = \ln y$

$y' = \ln \alpha + \beta x'$

$y = \alpha + \beta \log(x)$

$x' = \log(x)$

$y' = \alpha + \beta x'$

$y = \frac{\alpha}{\alpha x - \beta}$

$y' = \frac{1}{y}$

$y' = \alpha - \beta x'$

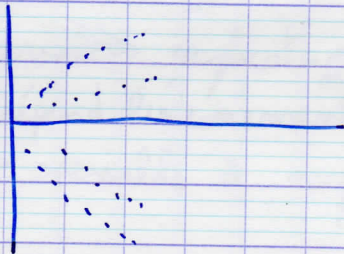
$x' = \frac{1}{x}$

$y = \frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}}$

$y' = \ln\left(\frac{y}{1-y}\right)$

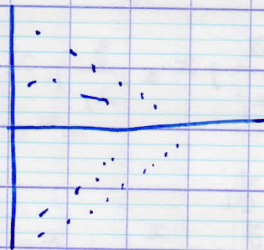
$x' = \alpha + \beta x$

Détection de non-linéarité par examen des résidus



pb grandes
valeurs de y

variance non
constante



pb
petites valeurs
de y

Transformation sur les valeurs de y.

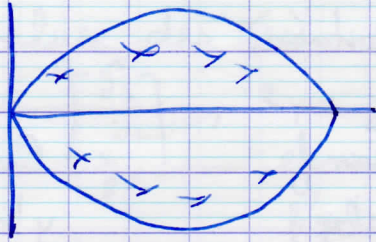
$y' = \sqrt{y}$

$y' = \log(y)$

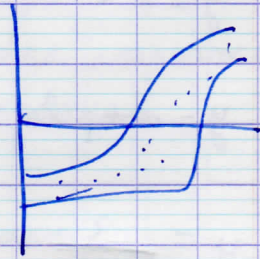
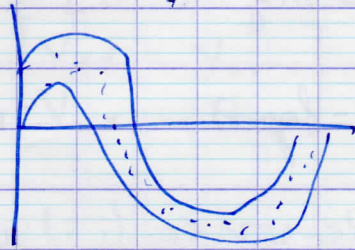
$y' = \log(y+1)$

$y' = \frac{1}{y}$

$y' = \frac{1}{y+1}$



$$y' = \log(y) \quad y' = \log(y+1)$$



variance non constante
pas de réponse.

Algo de la régression multi-linéaire

$$1^{\circ} \begin{matrix} y & | & x_1 & \dots & x_m \\ \frac{y}{x_1} & | & \frac{y}{x_2} & \dots & \frac{y}{x_m} \end{matrix}$$

x_i test de signification $F_{1, 1, m-2}$ est le meilleur

$$2^{\circ} \frac{y}{x_i}, x_1, \frac{y}{x_i}, x_2 \dots \frac{y}{x_i}, x_m$$

$$\textcircled{1} \begin{matrix} y_i \\ \frac{y}{x_i} \end{matrix}, x_j, x_2 \dots$$