

MOOC Statistique pour ingénieur

Thème 4 : Régression linéaire

Vidéo 1 : Mettre en œuvre la régression linéaire simple

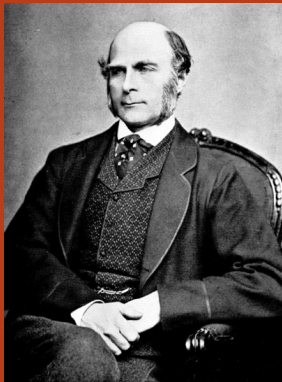
Anca Badea Lomig Hamon François Seyte Audrey Villot

Institut Mines-Télécom
Mines Saint-Étienne, Mines Nantes, Mines Alès, Mines Nantes

Sommaire

- 1 Introduction
- 2 Notations et vocabulaire
- 3 Estimation des paramètres
- 4 Analyse de la variance

Origine

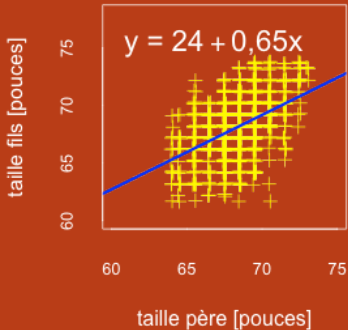


Francis Galton
(1822-1911)

Regression towards mediocrity in hereditary stature

Journal of the Anthropological Institute 15 :
246-63 (1886)

Origine



la taille des enfants nés des parents très grands (ou petits) se rapproche de la taille moyenne de la population → elle **régresse**

Objet

analyser la **relation** entre

- une **variable expliquée**
 - variable dépendante
 - réponse
 - variable endogène

et

- une ou plusieurs **variables explicatives**
 - variables indépendantes
 - prédicteurs
 - variables exogènes

dans un but

- explicatif
- de prévision
- ...

sous l'hypothèse que la
relation est
linéaire en ses paramètres

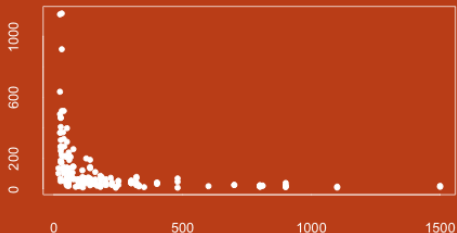
et
à partir de données

Exemples

- industrie

performance relative à un benchmark

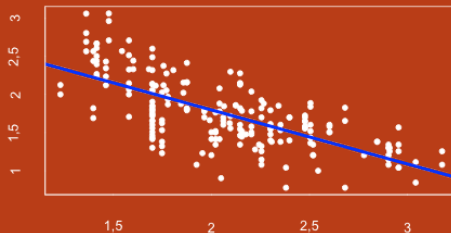
Performances des ordinateurs (P. Ein-Dor, J. Feldmesser (1987))



temps du cycle [nanosecondes]

$\log_{10}(\text{performance})$

Performances des ordinateurs (P. Ein-Dor, J. Feldmesser (1987))



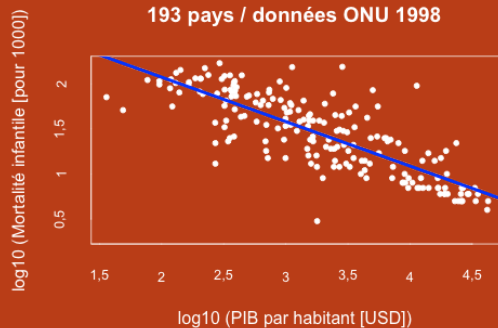
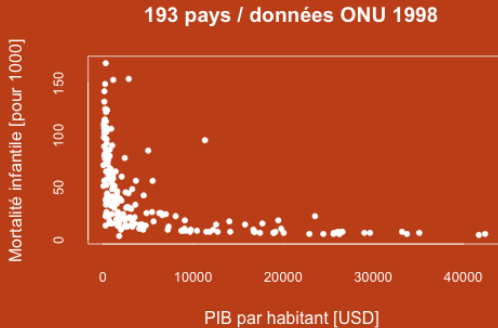
$\log_{10}(\text{temps du cycle [nanosecondes]})$

Exemples

- économie

Exemples

- économie



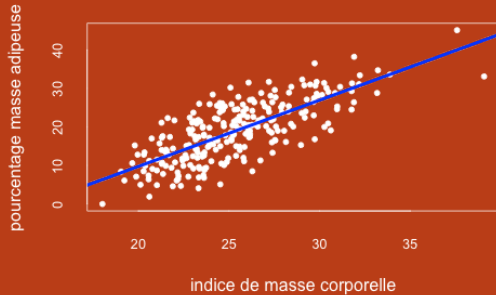
Exemples

- médecine

Exemples

- médecine

Mesures corporelles (Roger W. Johnson (1996))



- contrôle qualité
- sociologie
- métrologie
- marketing
- ...

Sommaire

- 1 Introduction
- 2 Notations et vocabulaire
- 3 Estimation des paramètres
- 4 Analyse de la variance

Vocabulaire

- régression linéaire simple :
 - une seule variable expliquée
 - une seule variable explicative
- régression linéaire multiple :
 - une seule variable expliquée
 - plusieurs variables explicatives

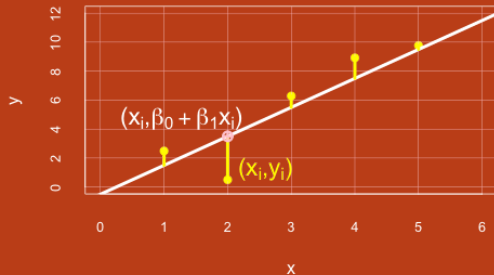
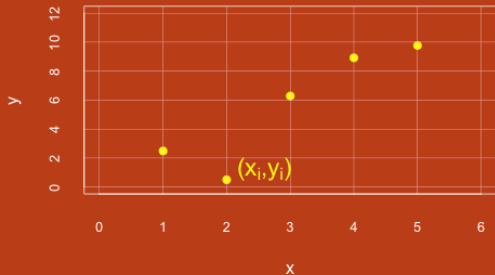
Notations

- y : la variable expliquée
- x : la variable explicative
- $y \approx \beta_0 + \beta_1 x$: la relation approximative entre y et x
- β_0, β_1 : les paramètres (coefficients) du modèle
- $\hat{\beta}_0, \hat{\beta}_1$: les estimations des paramètres
- $y = \beta_0 + \beta_1 x + \varepsilon$ avec ε : l'erreur
- $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$: la prédiction de la variable expliquée

Sommaire

- 1 Introduction
- 2 Notations et vocabulaire
- 3 Estimation des paramètres
- 4 Analyse de la variance

Les paramètres du modèle sont à **estimer** à partir des données $(x_i, y_i)_{i=1, \dots, n}$



interprétation : β_0 : l'ordonnée à l'origine, β_1 : la pente

Méthode des moindres carrés ordinaires

$$(\hat{\beta}_0, \hat{\beta}_1) = \operatorname{argmin} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$$

$$\begin{cases} \frac{\partial}{\partial \beta_0} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = 0 \\ \frac{\partial}{\partial \beta_1} \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2 = 0 \end{cases}$$

$$\begin{cases} \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 = \frac{\operatorname{Cov}(x, y)}{s_x^2} \end{cases}$$

les estimations des paramètres

Méthode des moindres carrés ordinaires

notations usuelles

- moyennes : $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$
- variances : $s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$ $s_y^2 = \frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2$
- covariance : $\text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$

Exemple

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

$$(1; 3, 36), (2; 0, 71), (3; 5, 27), (4; 7, 55), (5; 9, 01)$$

$$\hat{\beta}_0 = -0,27 ; \hat{\beta}_1 = 1,82$$

Sommaire

- 1 Introduction
- 2 Notations et vocabulaire
- 3 Estimation des paramètres
- 4 Analyse de la variance

Somme des carrés

$$y_i = \hat{y}_i + \varepsilon_i$$

avec

\hat{y}_i la part expliquée par le modèle et

ε_i la part inexpliquée : l'erreur

Et les sommes des carrés correspondantes ?

Somme des carrés

- somme des carrés totaux : $SCT = \sum_{i=1}^n (y_i - \bar{y})^2$
- somme des carrés expliqués : $SCE = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$
- somme des carrés résiduels : $SCR = \sum_{i=1}^n (y_i - \hat{y}_i)^2$

Equation d'analyse de la variance

$$SCT = SCE + SCR$$

Indicateur de qualité d'une régression

coefficient de détermination : $R^2 = \frac{SCE}{SCT} = 1 - \frac{SCR}{SCT}$

R^2 doit être proche de 1

rôle du statisticien dans l'interprétation

