

Analyse de données pour le génie industriel




Régression linéaire simple

Iragaël Joly Pierre Lemaire

Grenoble INP – Génie Industriel, 2A

2015–2016

Bibliographie

-  P.-A. Cornillon et al. *Statistiques avec R*, PU Rennes, 2012.
-  S. Tuff ry. *Data Mining et statistique d cisionnelle*, Technip, 2012.
-  G. Saporta. *Probabilit s et analyse de donn es*, Technip, 2011.

1. R gression lin aire simple

D finitions et calculs

Travaux pratiques

R gression (1)

- les donn es sont une matrice $n \times (p + 1)$
 - un  chantillon de n **observations**
 - chaque observation est constitu e de $p + 1$ variables
 - une **r ponse** (ou variable d pendante ou   expliquer) : y
 - des **variables explicatives** (ou ind pendantes) : x_1, x_2, \dots, x_p

$$\left(\begin{array}{c|cccc} y_1 & x_{11} & x_{21} & \dots & x_{p1} \\ y_2 & x_{12} & x_{22} & \dots & x_{p2} \\ \vdots & & & & \vdots \\ y_n & x_{1n} & x_{2n} & \dots & x_{pn} \end{array} \right)$$

R gression

une r gression de y sur x est un mod le de la forme

$$y = f(x) + \varepsilon$$

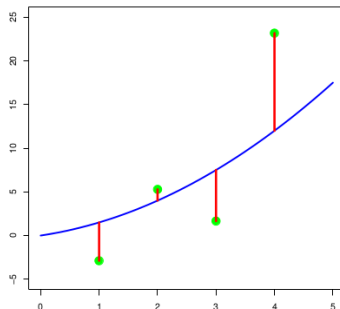
R gression (2)

- f caract rise l'effet de x sur y ; la forme de f est une **hypoth se**
- ε repr sente l'effet des autres facteurs (non pris en compte)

Erreurs de r gression

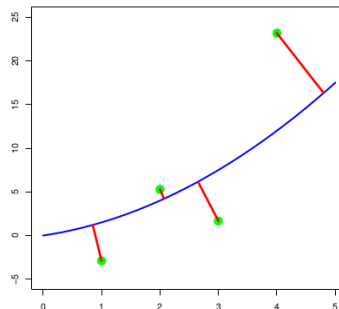
erreur «ordinaire»

x est parfaitement connu, toute
l'erreur est sur y



erreur orthogonale

l'erreur est partag e entre x
et y



- mod le $y = f(x) + \varepsilon$ implique une erreur ordinaire
- l'erreur ordinaire est une hypoth se   valider

TP :  chauffement

- 1 G n rez $x = (1; 2; \dots; 100)$ et y de la forme $y = f(x) + \varepsilon$
on prendra $y = 0.05x^2 - 3x + 7 + \mathcal{N}(\mu = 0, \sigma = 50)$
- 2 Tracez y en fonction de x ; donnez un nom aux axes et au graphique.
- 3 Calculez la variance de x et la covariance de x et y (faire vos propres fonctions). Comparez aux r sultats des fonctions `var` et `cov`.
- 4 *Astuce* : savez-vous afficher la valeur de la variance de x avec exactement 3 d cimales (`sprintf`) ?

TP : Calculs de base

- 1 Calculez β_1 et β_0 (faire des fonctions)
- 2 D finissez deux fonctions qui permettent respectivement de calculer une r gression lin aire et de l  valuer sur des nouvelles donn es.
- 3 Calculez la r gression lin aire de y en fonction de x
- 4 Ajoutez la droite de r gression sur le graphique (abline).
- 5 *Astuce* : savez-vous modifier la couleur, l' paisseur ou le style d'un trac  (`col`, `lwd`, `lty`) ?

Travaux pratiques

Coefficient de corrélation linéaire

Coefficient de corrélation

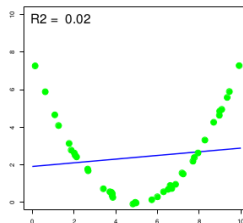
Le coefficient de corrélation linéaire (empirique) entre les x_i et les y_i est

$$r_{xy} = \frac{C_{xy}}{s_x s_y}$$

où C_{xy} est la covariance entre x et y , s_x et s_y sont les écarts-types de x et y

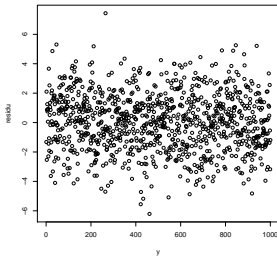
- $\delta^2 = s_y^2(1 - r_{xy}^2)$: minimiser δ est équivalent à maximiser r_{xy}^2
- $r_{xy} = \pm 1$ ssi les points sont parfaitement alignés
- $r_{xy} \in [-1; 1]$

Coefficient de détermination (2)

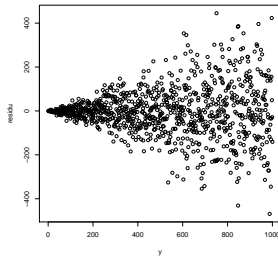


Analyse des résidus

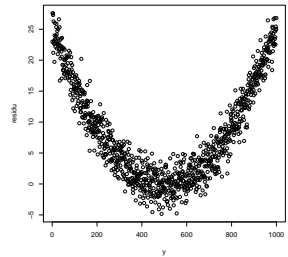
- résidus en fonction de x (ou y ou \hat{y})



normal

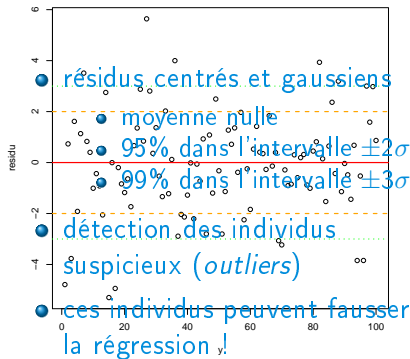


hétéroscédasticité



composante
non linéaire

Observations suspectes



Évaluation de la qualité d'une régression

Travaux pratiques

TP : Qualit  de la r gression

- 1 Calculez le coefficient de corr lation lin aire (faire une fonction).
- 2 Calculez la variance totale, la variance expliqu e, la variance r siduelle (faire des fonctions). V rifiez la d composition de la variance.
- 3 Calculez le coefficient de d termination (faire une fonction).
- 4 *Astuce* : savez-vous ajouter ces grandeurs sur le graphique (text) ?

TP : Résidus

- ① Tracez les résidus en fonction de x .
- ② Ces résidus vous semblent-ils «normaux» ? Y a-t-il des points suspects ?
- ③ Calculez la régression $y = ax^2 + b$.
- ④ Est-elle meilleure ?
- ⑤ Parmi les régression de la forme $y = ax^k + b$, laquelle vous semble la meilleure ?

26/37

TP : lm (et predict)

-   R cup rez le fichier 003_regression_lin aire_lm.R
-   Ex cutez-le pas   pas, en prenant le temps de comprendre chaque instruction propos e.
-   V rifiez que vous savez
 -   importer ou g n rer des donn es, tracer des graphiques pour les repr senter ;
 -   calculer une r gression lin aire avec une ou plusieurs variables, transformer les variables ;
 -   d terminer les caract ristiques essentielles d'une r gression lin aire (coefficients, R^2 , r sidus, ...) ;
 -   calculer des pr dictions pour les r gressions calcul es.

Régression linéaire multiple (2)

- hypothèse : les résidus sont centrés
 - si les résidus sont indépendants, de même loi et même variance σ^2 , on parle de **moindres carrées ordinaires**
 - sinon on parle de **moindres carrées généralisés**
- il n'y a pas d'hypothèse sur les régresseurs (en particulier, ils n'ont pas à être indépendants)

- écriture matricielle : $Y = X\beta + \varepsilon$

$$Y = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_{11} & \dots & x_{p1} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1n} & \dots & x_{pn} \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_p \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

Coefficient de détermination

Coefficient de détermination

Le coefficient de détermination de la régression linéaire est

$$R^2 = \frac{SSR}{SST} = 1 - \frac{SSE}{SST}$$

- identique à la régression linéaire simple
- « *multiple R-squared* » sous R
- propriété : $R^2 = r^2$
corollaire : R^2 augmente avec le nombre de variables...

35/37

TP : Consommation de véhicules (1)

Jeu de données mpg.dat (UCI Machine Learning Repository) :

- mpg (*miles per gallon*) : la distance parcourue
- cylinders : le nombre de cylindres.
- displacement (pouces-cube) : la cylindrée
- horsepower : puissance du moteur
- weight (livres) : poids du véhicule
- acceleration : secondes pour atteindre 100 m/h
- model year : année du modèle.
- origin : lieu (1 - USA ; 2 - EU ; 3 - Japon).
- name : nom du modèle (unique).

TP : Consommation de véhicules (2)

TP noté !

- à faire en binôme.
- à rendre sur Chamilo le dimanche 6 mars 2016.
- rendu : script R **commenté**.
 - Les commandes permettant de répondre à une question doivent être fournies, à l'exclusion de toutes autres.
 - Les graphiques doivent être soignés, clairs et complets.
 - Une réponse explicite, pour chaque question, est requise. Cette réponse sera rédigée convenablement (hormis fonctions ou graphiques demandés).