

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. Tufféry. *Data Mining et statistique décisionnelle*, Technip, 2012.
-  G. Saporta. *Probabilités et analyse de données*, 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

$$\begin{pmatrix} y_1 & x_{11} & x_{21} & \dots & x_{p1} \\ y_2 & x_{12} & x_{22} & \dots & x_{p2} \\ \vdots & & & & \vdots \\ y_n & x_{1n} & x_{2n} & \dots & x_{pn} \end{pmatrix}$$

Régression (2)

Régression

une régression de y sur x est un modèle de la forme

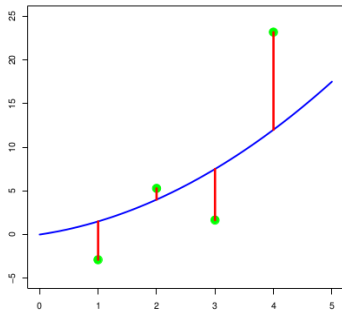
$$y = f(x) + \varepsilon$$

- f caractérise l'effet de x sur y ; la forme de f est une **hypothèse**
- ε représente l'effet des autres facteurs (non pris en compte)

Erreurs de r gression

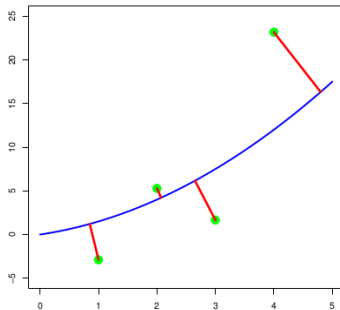
erreur «ordinaire»

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

Agr gation des erreurs de r gression

- enjeu : minimiser les erreurs, donc la norme de ε
- norme L_1 : $\|\varepsilon\|_1 = \sum_{i=1}^n |\varepsilon_i|$
 - poids des erreurs proportionnel   la valeur
 - non d rivable en z ro, donc peu adapt e   l'optimisation
- norme L_2 : $\|\varepsilon\|_2 = \sqrt{\sum_{i=1}^n \varepsilon_i^2}$
 - privil gie plusieurs petites erreurs   une grosse
 - tr s bien adapt e   l'optimisation
 - « moindres carr s »
- norme L_∞ : $\|\varepsilon\|_\infty = \max_{i=1}^n |\varepsilon_i|$
 - erreur «  gale » pour tous (id alement), contr le du pire cas
 - peu adapt e   l'optimisation
- nous nous restreindrons aux moindres carr s ordinaires

R gression lin aire simple (1)

- on se limite   une unique variable explicative x
- on se limite   une fonction lin aire

R gression lin aire simple

Le mod le de r gression lin aire simple est d fini par

$$y_i = \beta_1 x_i + \beta_0 + \varepsilon_i, i = 1..n$$

- β_1 et β_0 sont les param tres inconnus   estimer
- ε_i est l'erreur ou **r sidu** pour l'observation i

Droite des moindres carrés

Droite des moindres carrés

La droite des moindres carrés est $y = \beta_1 x + \beta_0$ avec

$$\beta_1 = \frac{C_{xy}}{s_x^2}$$

$$\beta_0 = \bar{y} - \beta_1 \bar{x}$$

- $\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$ (moyenne de x)
- $\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ (moyenne de y)
- $s_x^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2$ (variance de x)
- $C_{xy} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y}) = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}$ (covar. entre x et y)
- β_1 et β_0 sont estimateurs sans biais des variables aléatoires correspondantes
- $y = \bar{y} - \frac{C_{xy}}{s_x^2} (x - \bar{x})$
- la droite passe par le barycentre $(\bar{x}; \bar{y})$

1. R gression lin aire simple

D finitions et calculs

Travaux pratiques

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 Définissez deux fonctions qui permettent respectivement de calculer une régression linéaire et de l'évaluer sur des nouvelles données.
- 3 Calculez la régression linéaire de y en fonction de x
- 4 Ajoutez la droite de régression sur le graphique (abline).
- 5 *Astuce* : savez-vous modifier la couleur, l'épaisseur ou le style d'un tracé (col, lwd, lty) ?

TP : Calculs des erreurs

- 1 Calculez le vecteur des erreurs.
- 2 Calculez l'erreur quadratique moyenne. Calculez l'erreur moyenne, l'erreur absolue moyenne, l'erreur relative moyenne, l'erreur relative absolue moyenne.
- 3 D terminez pour combien et pour quelles valeurs de x l'erreur absolue est plus grande que 100.
- 4 Tracez la courbe $N(e)$ o  $N(e)$ repr sente le nombre de valeurs de valeurs de x pour lesquelles l'erreur est sup rieure   e .
- 5 *Astuce* : connaissez-vous `sapply` ?

2. Qualité d'une régression

Évaluation de la qualité d'une régression
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 d termination (1)

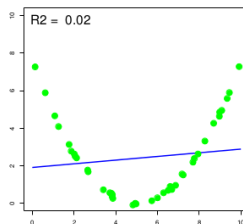
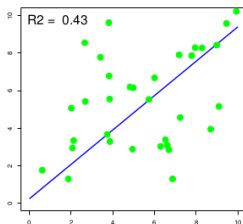
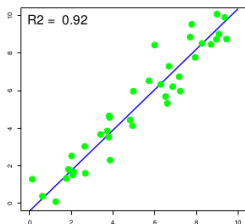
Coefficient de d termination

Le coefficient de d termination de la r gression lin aire est

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

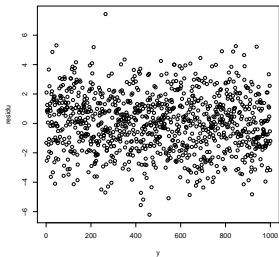
- R^2 est la proportion de la variance de y expliqu e par x
 - R^2 proche de 1 : x et y *semblent* en relation lin aire
 - R^2 proche de 0 : x et y ne sont pas li s *lin airement*
- propri t  : $R^2 = r_{xy}^2$

Coefficient de d termination (2)

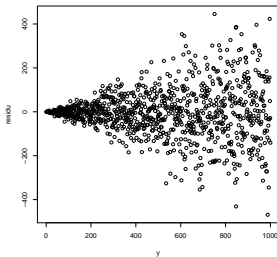


Analyse des r siduals

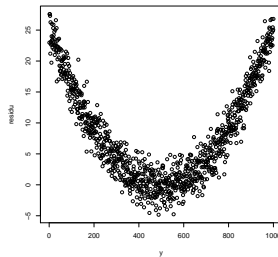
-   r siduals en fonction de x (ou y ou \hat{y})



normal



h t rosc dasticit 



composante
non lin aire

2. Qualit  d'une r gression

 valuation de la qualit  d'une r gression
Travaux pratiques

TP : R s dus

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

3. R gression lin aire multiple

G n ralisation de la r gression lin aire
Travaux pratiques

R gression lin aire multiple (1)

- on  tend la r gression lin aire   p variables explicatives
- on se limite   des fonctions lin aires (en chaque variable)

R gression lin aire multiple

Le mod le de r gression lin aire simple est d fini par

$$y_i = \beta_p x_{p,i} + \dots + \beta_1 x_{1,i} + \beta_0 + \varepsilon_i, i = 1..n$$

- les β_i sont les param tres inconnus   estimer
- ε_i est le r sidu pour l'observation i

Conclusion partielle

- nous savons calculer des modèles de régression linéaire
- nous savon évaluer ces modèles, dans une certaine mesure
- encore beaucoup à dire sur la construction et l'interprétation de modèles linéaires
 - autres critères de qualité (que le R_{adj}^2)
 - tests statistiques sur la significativité d'une variable
 - méthodes de construction et de sélection des variables
 - analyse des erreurs, intervalles de confiance des prédictions
 - ...

3. Régression linéaire multiple

Généralisation de la régression linéaire
Travaux pratiques

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 v hicules (2)

TP not  !

-     faire en bin me.
-     rendre sur Chamilo le dimanche 6 mars 2016.
-   rendu : script R **comment **.
 -   Les commandes permettant de r pondre   une question doivent  tre fournies,   l'exclusion de toutes autres.
 -   Les graphiques doivent  tre soign s, clairs et complets.
 -   Une r ponse explicite, pour chaque question, est requise. Cette r ponse sera r dig e convenablement (hormis fonctions ou graphiques demand s).