

Sommaire

Introduction

A-Définition et hypothèses

B-Estimation

C-Lois des estimateurs

D-Tests

E-Validation

Introduction

Références

A-Définition et hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Introduction

- ▶ Expliquer ou prédire une variable aléatoire **réponse** Y en fonction d'une liste de variables **explicatives** $\mathbf{X} = (X_1, \dots, X_p)$ **observées** sur des individus

$$Y = f(\mathbf{X})$$

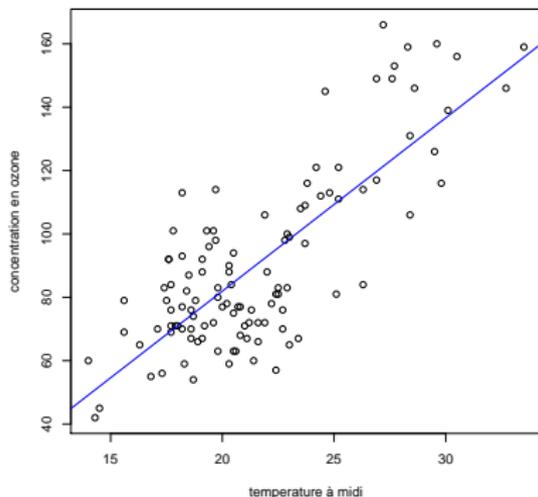
- ▶ Modéliser Y par une fonction des variables explicatives \mathbf{X} , par exemple, modèle à bruit additif

$$Y_{|\mathbf{X}=x} = m(x) + \varepsilon$$

$m(x) = \mathbb{E}(Y|\mathbf{X} = x)$ est la **fonction de régression**

- ▶ choix est guidé par la connaissance d'un phénomène physique, biologique,...
- ▶ par l'observation
- ▶ ou non...
- ▶ apprentissage **supervisé** : à partir d'un échantillon étiqueté $((\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n))$, **estimer** $m(x)$, **prédire** Y pour une nouvelle condition x ... avec **le moins d'erreur**

Exemple : régression linéaire (simple) ¹



$$Y_i = \underbrace{\theta_1 + \theta_2 X_i}_{\text{déterministe}} + \underbrace{\varepsilon_i}_{\text{bruit aléatoire}}, \quad \mathbf{E}_{X_i}(\varepsilon_i) = 0$$

Introduction

Références

A-Définition et hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des estimateurs

RC

D-Tests

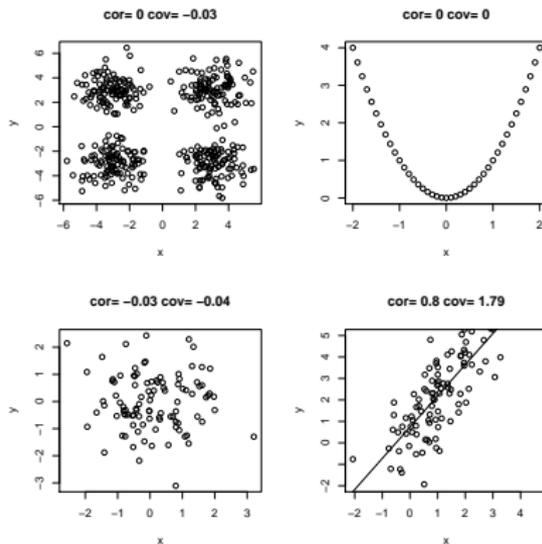
de Student
de Wald
de Fisher

E-Validation

1. <http://www.agrocampus-ouest.fr/math/livreR/>

Que cherche-t-on à détecter ?

La corrélation entre Y et $x...$ qui n'est que l'expression d'une **liaison linéaire**



Une corrélation nulle n'indique pas forcément une absence de liaison

Introduction

Références

A-Définition et hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des estimateurs

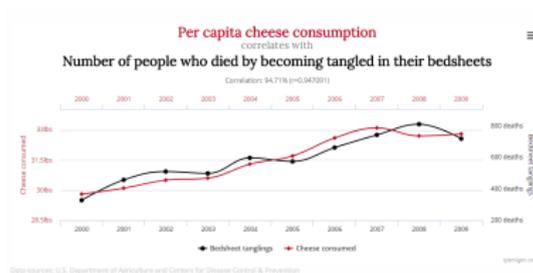
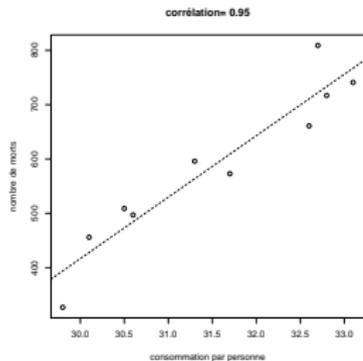
RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Mais attention : corrélation \neq causalité



Introduction

Références

A-Définition et hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des estimateurs

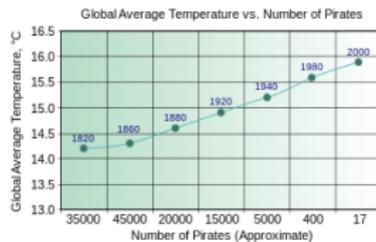
RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

source : <http://tylervigen.com/spurious-correlations>



source : Wikipédia

- ▶ **Données d'apprentissage** :
 $\mathcal{D} = \{(\mathbf{X}_1, Y_1), \dots, (\mathbf{X}_n, Y_n)\}$ (i.i.d. $\sim \mathbf{P}$)
- ▶ **Prédicteur** : $f : \mathcal{X} \rightarrow \mathcal{Y}$ mesurable
- ▶ **Fonction de perte ou de coût** : $\ell(f(\mathbf{X}), Y)$ mesure avec quelle qualité $f(\mathbf{X})$ "prédit" Y , d'où le **risque** :

$$\mathcal{R}(f) = \mathbf{E}\ell(Y, f(\mathbf{X}))$$

↪ souvent $\ell(f(\mathbf{X}), Y) = \|f(\mathbf{X}) - Y\|^2$ ou
 $\ell(f(\mathbf{X}), Y) = \mathbf{1}_{Y \neq f(\mathbf{X})}$

But : apprendre une règle pour construire un **classifieur** / **régresseur** $\hat{f} \in \mathcal{F}$ à partir des données d'apprentissage \mathcal{D}_n avec un **risque** $\mathcal{R}(\hat{f})$ **petit en moyenne**.

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Fonction de régression

Approcher Y par une fonction $h(X)$ de la variable X .

Risque quadratique :

$$R(h) = \mathbb{E}[(Y - h(X))^2] = \int (y - h(x))^2 dP(x, y).$$

Théorème

La fonction qui minimise le risque quadratique $R(\cdot)$ est $m(X) = \mathbb{E}(Y|X)$, l'espérance de Y conditionnellement à X

$$\forall x, m(x) = \mathbb{E}(Y|X = x).$$

*Elle est appelée **fonction de régression**.*

- ▶ PB : on ne sait pas la calculer dans le cas général
- ▶ Si X est déterministe, on travaille également à $X = x$ fixé... et les développements sont identiques.

↪ poser des hypothèses de modélisation.

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

J. Pagès. *Statistique générales pour utilisateurs*. Presses Universitaires de Rennes, Rennes, 2005.

J.-M. Azais et J.-M. Bardet. *Le modèle linéaire par l'exemple*. Dunod, Paris, 2005.

P.-A. Cornillon et E. Matzner-Løber. *Régression Théorie et applications*. Springer, Paris, 2007.

P.-A. Cornillon et al. *Statistiques avec R*. Presses Universitaires de Rennes, Rennes, 2008.

Sommaire

Introduction

A-Définition et hypothèses

B-Estimation

C-Lois des estimateurs

D-Tests

E-Validation

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Régression linéaire : Y est une variable quantitative

à partir de n observations $(x_1, Y_1), \dots, (x_n, Y_n)$

1. Modèle à **bruit additif** posant une relation de **linéarité de l'espérance** $\mathbf{E}(Y_i|x_i)$ en un **paramètre** θ :

$$\begin{aligned} Y_i &= x_{i1}\theta_1 + \dots + x_{ip}\theta_p + \varepsilon_i \\ &= x_i\theta + \varepsilon_i, \end{aligned}$$

2. Les **résidus sont centrés** : $\mathbf{E}(\varepsilon_i|x_i) = 0$.
3. La **variance de l'erreur est constante** : $\text{var}(\varepsilon_i|x_i) = \sigma^2$.
4. Les résidus sont **décorrélés** : $\text{cov}(\varepsilon_i, \varepsilon_j|x_i, x_j) = 0$ pour $i \neq j$.
5. Eventuellement : ε_i est gaussien

$\hookrightarrow x_i$ est la valeur **fixée** des variables explicatives pour la i -ème observation. Elle n'est pas aléatoire

$\hookrightarrow m(x_i) = x_i\theta$ est appelée **fonction de régression**

[Introduction](#)[Références](#)[A-Définition et hypothèses](#)[Exemples](#)[Identifiabilité](#)[B-Estimation](#)[EMC](#)[Propriétés](#)[EMV](#)[C-Lois des estimateurs](#)[RC](#)[D-Tests](#)[de Student](#)[de Wald](#)[de Fisher](#)[E-Validation](#)

- ▶ Matriciellement : $Y = X\theta + \varepsilon$

$$Y_{n \times 1} = X_{n \times p} \theta_{p \times 1} + \varepsilon_{n \times 1};$$

$$\mathbb{E}(\varepsilon_{n \times 1}) = 0; \quad \text{var}(\varepsilon_{n \times 1}) = \sigma^2 Id_n,$$

X est la **matrice du plan d'expérience**, concaténation des n vecteurs lignes x_i ou des p variables colonnes $X_j = (x_{1j}, \dots, x_{nj})'$.

Remarque : En général, $\forall i, x_{i1} = 1$, X_1 est l'**intercept**.

- ▶ Modèle linéaire gaussien

$$Y = m + \varepsilon; \quad \varepsilon \sim \mathcal{N}(0, \Sigma); \quad (\Sigma = \sigma^2 Id_n)$$

où $m \in V$, sous espace vectoriel de \mathbb{R}^n de dimension p .

Introduction

Références

A-Définition et
hypothèses

Exemples

Identifiabilité

B-Estimation

EMC

Propriétés

EMV

C-Lois des
estimateurs

RC

D-Tests

de Student

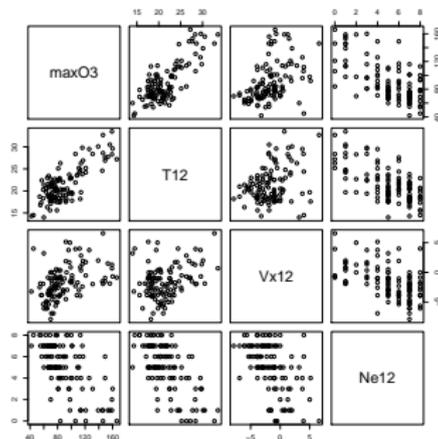
de Wald

de Fisher

E-Validation

Régression multiple

- ▶ les covariables X sont quantitatives
- ▶ si $p = 2$ dont un intercept : régression **simple**
 $m(x_i) = \theta_1 + \theta_2 x_i$; $x_i = (1 \ x_i)$; $\theta = (\theta_1, \theta_2)'$
- ▶ si $p > 2$, régression **multiple**



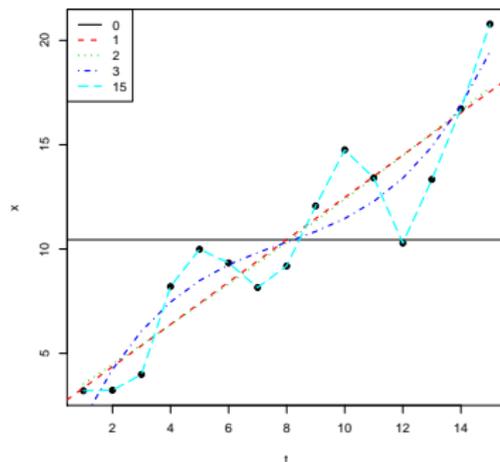
$$\text{maxO3} = \theta_1 + \text{T12} \theta_2 + \text{Vx12} \theta_3 + \text{Ne12} \theta_4 + \varepsilon$$

Régression polynomiale

Un cas particulier de régression multiple : régression polynômiale

$$Y_t = \theta_1 + t\theta_2 + t^2\theta_3 + t^3\theta_4 + \varepsilon_t$$

$$X = \begin{pmatrix} 1 & t_1 & \cdots & t_1^{p-1} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & t_n & \cdots & t_n^{p-1} \end{pmatrix}$$



Un régression linéaire (en θ) peut permettre de représenter un phénomène non-linéaire (en t)

Identifiabilité

Définition : deux paramètres différents définissent deux lois différentes

Théorème (CNS d'identifiabilité)

La paramétrisation $\mathbb{E}(Y) = X\theta$ est *identifiable* si et seulement si l'une des propriétés équivalentes suivantes est vérifiée :

- ▶ les colonnes de X sont indépendantes,
- ▶ X est de rang plein,
- ▶ X est injective,
- ▶ $\text{Ker}(X) = \{0\}$.

Définition

On appelle *dimension du modèle* la dimension de $\text{Im}(X)$.

On a : $\dim(\text{Im}(X)) \leq \dim(\theta)$

↪ On suppose dans la suite que le modèle est identifiable ou *régulier*.

Sommaire

Introduction

A-Définition et hypothèses

B-Estimation

C-Lois des estimateurs

D-Tests

E-Validation

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Estimateur des Moindres Carrés (EMC)

Recherche un estimateur de θ sous forme d'un minimiseur de la **somme des carrés résiduels** $SCR(\theta) = \|Y - X\theta\|^2$:

$$\hat{\theta} = \arg \min_{\theta \in \Theta} \|Y - X\theta\|^2.$$

On a : $\hat{\theta} = (X'X)^{-1}X'Y$

Proposition

$\hat{Y} := \widehat{\mathbf{E}(Y)} = X\hat{\theta} = H_X Y$ est la projection orthogonale de Y sur $Im(X)$:

$$H_X = X(X'X)^{-1}X'$$

- ▶ H_X : **Hat matrix**
- ▶ \hat{Y} : **valeurs ajustées**
- ▶ $\hat{\varepsilon} = Y - \hat{Y}$: **résidus** (estimés)
- ▶ la droite de régression $y = x\hat{\theta}$ passe par le point moyen $(\bar{Y}, \bar{X}^{(1)}, \dots, \bar{X}^{(p)})$

[Introduction](#)[Références](#)[A-Définition et hypothèses](#)[Exemples](#)[Identifiabilité](#)[B-Estimation](#)[EMC](#)[Propriétés](#)[EMV](#)[C-Lois des estimateurs](#)[RC](#)[D-Tests](#)[de Student](#)[de Wald](#)[de Fisher](#)[E-Validation](#)

Exemple : régression linéaire (simple)

Le modèle s'écrit $Y = \theta_1 \mathbb{1} + \theta_2 X_2 + \varepsilon$, où $\mathbb{1}$ est l'intercept, et $X_2 = (x_1, \dots, x_n)'$

$$X'X = \begin{pmatrix} n & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{pmatrix}$$

$$(X'X)^{-1} = \frac{1}{n \sum_i x_i^2 - \sum_i x_i \sum_i x_i} \begin{pmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & n \end{pmatrix}$$

soit

$$\hat{\theta} = (X'X)^{-1}X'Y = \begin{pmatrix} \hat{\theta}_1 = \bar{Y} - \bar{x}\hat{\theta}_2 \\ \hat{\theta}_2 = \frac{(\frac{1}{n}\sum_i x_i Y_i) - \bar{x}\bar{Y}}{(\frac{1}{n}\sum_i x_i^2) - \bar{x}^2} \end{pmatrix}$$

La droite de régression **estimée** $\hat{y} = \hat{\theta}_1 + \hat{\theta}_2 x$ passe par le **point moyen** (\bar{x}, \bar{Y}) .

Introduction

Références

A-Définition et
hypothèses

Exemples

Identifiabilité

B-Estimation

EMC

Propriétés

EMV

C-Lois des
estimateurs

RC

D-Tests

de Student

de Wald

de Fisher

E-Validation

Propriétés ne nécessitant pas l'hyp. gaussienne

Proposition

L'EMC $\hat{\theta}$ de θ est

- ▶ *sans biais* : $\mathbb{E}(\hat{\theta}) = \theta$
- ▶ *de variance* $\text{var}(\hat{\theta}) = \sigma^2(X'X)^{-1}$.

De plus, $\hat{\theta}$ vérifie la propriété de **Gauss-Markov** :

Théorème (Gauss-Markov)

Parmi les estimateurs *linéaires et sans biais* de θ , l'EMCO $\hat{\theta}$ est de variance minimum.

Estimateurs de la variance σ^2 :

- ▶ $\hat{s}^2 = \text{SCR}(\hat{\theta})/n$ est biaisé à distance finie, et asymptotiquement sans biais
- ▶ $\hat{\sigma}^2 = \text{SCR}(\hat{\theta})/(n - p)$ est sans biais

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Estimateur du Maximum de Vraisemblance (EMV)

Si la loi de ε est connue (gaussienne), la **vraisemblance** de l'échantillon est

$$L_n(\theta, \sigma^2; Y) = \frac{1}{(\sqrt{2\pi\sigma^2})^n} \exp\left(-\frac{1}{2\sigma^2} \underbrace{\|Y - X\theta\|^2}_{SCR(\theta)}\right).$$

L'EMV $\hat{\beta} = (\hat{\theta}, \hat{\sigma}^2)$ est

$$\hat{\beta} = \arg \max_{\beta \in \mathbb{R}^p \times \mathbb{R}^{+*}} \log L_n(\beta; Y).$$

D'où

$$\hat{\theta} = (X'X)^{-1}X'Y.$$

... identique à l'EMC. L'EMV de σ^2 est

$$\hat{\sigma}^2 = \frac{1}{n} \|Y - X\hat{\theta}\|^2 = \frac{1}{n} SCR(\hat{\theta})$$

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Sommaire

Introduction

A-Définition et hypothèses

B-Estimation

C-Lois des estimateurs

D-Tests

E-Validation

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Théorème

L'EMV $\hat{\beta} = (\hat{\theta}, \hat{s}^2)$ du modèle linéaire **gaussien**

$$Y = X\theta + \varepsilon, \quad \varepsilon \sim \mathcal{N}_n(0, \sigma^2 I_n), \quad \theta \in \mathbb{R}^p, \quad \sigma \in \mathbb{R}^{+*},$$

supposé régulier, suit la loi suivante :

$$\hat{\theta} \sim \mathcal{N}(\theta, \sigma^2(X'X)^{-1}); \quad \hat{s}^2 = \frac{\|Y - X\hat{\theta}\|^2}{n} \sim \frac{\sigma^2}{n} \chi^2(n-p).$$

De plus, $\hat{\theta}$ et \hat{s}^2 sont indépendants

Csq : Soit A une matrice $q \times p$ de rang $q \leq p$

$$A\hat{\theta} \sim \mathcal{N}_q(A\theta, \sigma^2 A(X'X)^{-1}A')$$

Introduction

Références

A-Définition et
hypothèses

Exemples

Identifiabilité

B-Estimation

EMC

Propriétés

EMV

C-Lois des
estimateurs

RC

D-Tests

de Student

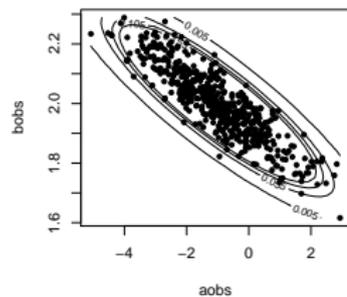
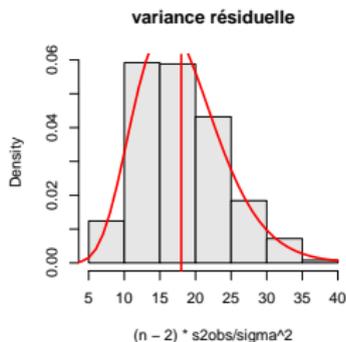
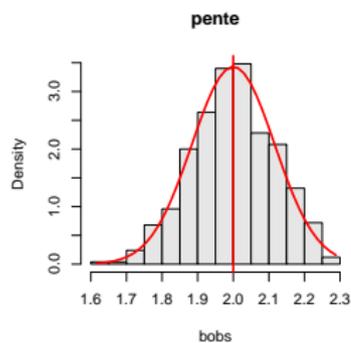
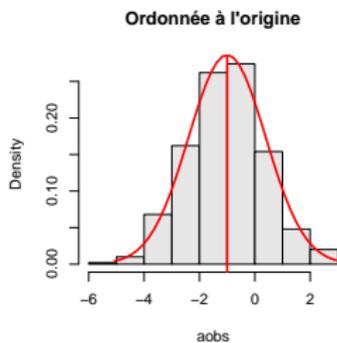
de Wald

de Fisher

E-Validation

Illustration simulée

500 échantillons (Y, X) de taille $n = 30$, avec $x = 1 : 30$ et $m(x) = -1 + 2x$



Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Cas gaussien

Quand ε est **gaussien**, on utilise plutôt l'estimateur non biaisé de la variance $\hat{\sigma}^2 = \|Y - X\hat{\theta}\|^2 / (n - p)$

Résultats

- ▶ $\hat{\theta}_{MV} = \hat{\theta}_{MC}$ et $\hat{\sigma}^2$ sont **indépendants**
- ▶ $\hat{\theta} \sim \mathcal{N}(\theta, \sigma^2(X'X)^{-1})$, $(n - p)\hat{\sigma}^2 / \sigma^2 \sim \chi^2(n - p)$
- ▶ pour toute matrice L non nulle de dimension $(1 \times p)$:
$$T_n = (L(X'X)^{-1}L')^{-1/2}(L\hat{\theta} - L\theta) / \hat{\sigma} \sim \mathcal{T}(n - p)$$
- ▶ pour toute matrice A de dimension $r \times p$ et de rang $r \leq p$:
$$F_n = (A(\hat{\theta} - \theta))' [A(X'X)^{-1}A']^{-1} A(\hat{\theta} - \theta) / (r\hat{\sigma}^2) \sim \mathcal{F}(r, n - p)$$
- ▶ si $r = 1$, $F_n = T_n^2$

Introduction

Références

A-Définition et
hypothèsesExemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMVC-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Théorème

Soit L une matrice de dimension $1 \times p$. Un **intervalle de confiance** de niveau $1 - \alpha$ d'une forme linéaire de $L\theta$ est donné par

$$\left[L\hat{\theta} \pm q_{1-\alpha/2}^{T_{n-p}} \hat{\sigma} \sqrt{L(X'X)^{-1}L'} \right],$$

où $q_{1-\alpha/2}^{T_{n-p}}$ est le quantile d'ordre $1 - \alpha/2$ de la loi de Student à $n - p$ degrés de liberté.

Applications : IC d'une composante de θ , IC de $\mathbb{E}(Y|x_0)$

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

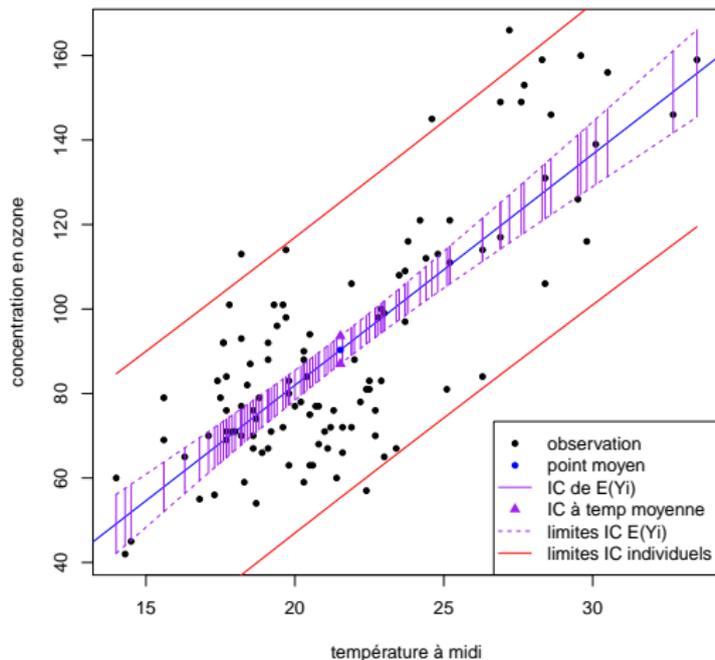
RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

IC d'une valeur moyenne ou d'une prédiction individuelle



Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Intervalle de confiance d'une prévision

A partir de n observations ($\hookrightarrow \hat{\theta}$), prédire une **nouvelle** observation **individuelle indépendante** sous la condition d'expérience $x_{n+1} = (x_{n+1,1}, \dots, x_{n+1,p})$

- ▶ modèle : $Y_{n+1} = x_{n+1}\theta + \varepsilon_{n+1}$,
- ▶ observation prédite : $\hat{Y}_{n+1}^p = x_{n+1}\hat{\theta}$
- ▶ erreur de prévision

$$\hat{\varepsilon}_{n+1}^p = Y_{n+1} - \hat{Y}_{n+1}^p,$$

$$\mathbf{E}(\hat{\varepsilon}_{n+1}^p) = 0, \quad \text{var}(\hat{\varepsilon}_{n+1}^p) = \sigma^2(x_{n+1}(X'X)^{-1}x_{n+1}' + 1).$$

- ▶ IC de prévision de niveau $1 - \alpha$:

$$\left[x_{n+1}\hat{\theta} \pm q_{1-\alpha/2}^{T_{n-p}} \hat{\sigma} \sqrt{x_{n+1}(X'X)^{-1}x_{n+1}' + 1} \right].$$

Introduction

Références

A-Définition et
hypothèsesExemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMVC-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Région de confiance simultanée (Bonferroni)

Soit θ_1 et θ_2 deux composantes de θ .

On cherche une région de confiance de la forme

$$IC(\theta_1) \times IC(\theta_2)$$

$$\mathbb{P}((\theta_1, \theta_2) \notin RC_\alpha(\theta_1, \theta_2))$$

$$\leq \mathbb{P}(\theta_1 \notin IC_{1-\alpha/2}(\theta_1)) + \mathbb{P}(\theta_2 \notin IC_{1-\alpha/2}(\theta_2)) \leq \alpha.$$

K intervalles de confiance de Bonferroni de risque α/K forment une région de confiance de risque simultané α .

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Région de confiance simultanée (autre méthode)

Théorème

Soit A une matrice de dimension $r \times p$ et de rang r . Une **région de confiance** de niveau $1 - \alpha$ de $A\theta$ est donnée par

$$RC_{\alpha}(A\theta) =$$

$$\left\{ u \in \mathbb{R}^p, \frac{1}{r\hat{\sigma}^2} (A\hat{\theta} - Au)' [A(X'X)^{-1}A']^{-1} (A\hat{\theta} - Au) \leq f_{r, n-p, 1-\alpha} \right\}$$

où $f_{r, n-p, 1-\alpha}$ est le quantile d'ordre $1 - \alpha$ de la loi de $\mathcal{F}(r, n - p)$.

Introduction

Références

A-Définition et
hypothèses

Exemples

Identifiabilité

B-Estimation

EMC

Propriétés

EMV

C-Lois des
estimateurs

RC

D-Tests

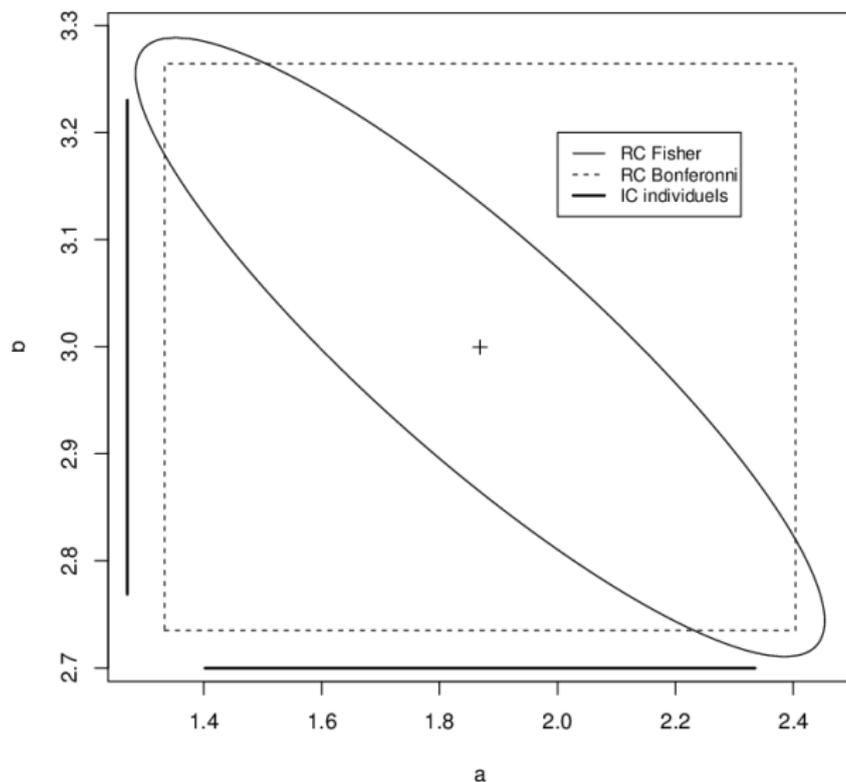
de Student

de Wald

de Fisher

E-Validation

Comparaison de régions de confiance



Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Sommaire

Introduction

A-Définition et hypothèses

B-Estimation

C-Lois des estimateurs

D-Tests

E-Validation

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Test de Student

Test d'une **relation affine** des composantes du paramètre $L\theta = c$, avec L matrice ligne de dimension p .

- ▶ test bilatéral : $H_0 : L\theta = c$ contre $H_1 : L\theta \neq c$
- ▶ statistique de test

$$T = \frac{L\hat{\theta} - c}{\hat{\sigma} \sqrt{L(X'X)^{-1}L'}} \sim_{H_0} \mathcal{T}(n-p),$$

- ▶ région de rejet du test bilatère de risque α

$$\mathcal{R}_\alpha = \{|T| > t_{n-p, 1-\alpha/2}\}.$$

- ▶ utilisé pour tester la **significativité** (non-nullité) d'une composante de θ

Remarque : test unilatéral : $H_0 : L\theta = c$ contre $H_1 : L\theta < c$
de région de rejet à gauche :

$$\mathcal{R}_\alpha = \{T < t_{n-p, \alpha}\}.$$

$$\text{Exemple : } Y = \mu + \beta_1 V1 + \beta_2 V2 + \beta_3 V3 + \varepsilon$$

```
> out=lm(Y~V1+V2+V3, data=df); summary(out)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	H0	H1
(Intercept)	-20.3129	51.8246	-0.392	0.701		
V1	-3.2241	3.6100	-0.893	0.386	~V2+V3	vs ~V1+V2+V3
V2	0.2334	0.2768	0.843	0.412	~V1+V3	vs ~V1+V2+V3
V3	20.3895	1.2662	16.103	7.1e-11 ***	~V1+V2	vs ~V1+V2+V3

Residual standard error: 4.163 on 15 degrees of freedom

F-statistic: 584.6 on 3 and 15 DF, p-value: 9.386e-16

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Test de Wald avec loi exacte

Modèle emboîté $\omega \subset \Omega$ défini par des hypothèses **affines** :
 $A\theta = c$ contre $A\theta \neq c$, A matrice $r \times p$ de rang r ,
 $\dim(\omega) = p - r = q$

$$W_{\text{exact}} = W/r = W = \frac{1}{r\hat{\sigma}^2} (A\hat{\theta} - c)' [A(X'X)^{-1}A']^{-1} (A\hat{\theta} - c) \sim_{(H_0)} \mathcal{F}(r, n-p)$$

- tester la significativité **globale** de la régression : tous les coeff. sont nuls sauf celui de l'intercept : $(H_0) : A\theta = 0$ contre $(H_0) : A\theta \neq 0$

$$A = \begin{pmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}$$

F-statistic: 584.6 on 3 and 15 DF, p-value: 9.386e-16

- tester des modèles **emboîtés**.

Introduction

Références

A-Préfixion et hypothèses

Exemples

Identifiabilité

B-Estimation

EMC

Propriétés

EMV

C-Lois des estimateurs

RC

D-Tests

de Student

de Wald

de Fisher

E-Validation

Petite interprétation géométrique

Soit ω un sev de $\Omega \in \mathbb{R}^n$, $\dim(\omega) = q < \dim(\Omega) = p$.

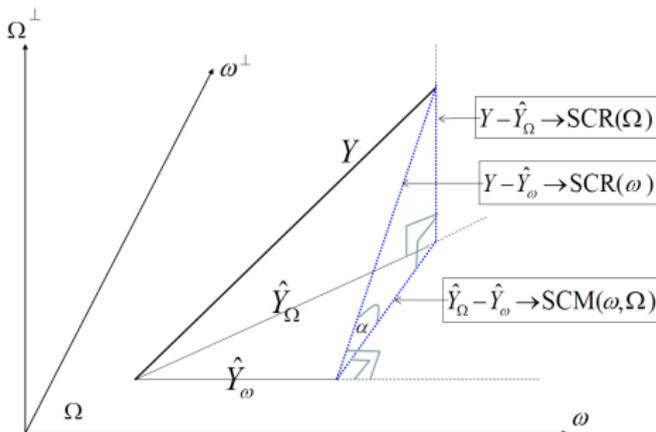
On note $\hat{Y}_\Omega = H_\Omega(Y)$. On a :

$$\|Y - \hat{Y}_\omega\|^2 = \|\hat{Y}_\Omega - \hat{Y}_\omega\|^2 + \|Y - \hat{Y}_\Omega\|^2$$

$$\text{SCR}(\omega) = \text{SCM}(\omega, \Omega) + \text{SCR}(\Omega).$$

$$F = \frac{(\text{SCR}(\omega) - \text{SCR}(\Omega))/(p - q)}{\text{SCR}(\Omega)/(n - p)} = \frac{\text{SCM}(\omega, \Omega)/(p - q)}{\text{SCR}(\Omega)/(n - p)}$$

$$\sim_\omega \mathcal{F}(p - q, n - p).$$



Théorème (Test de Fisher d'un sous-modèle linéaire)

Soit le modèle linéaire gaussien $Y = m + \varepsilon$, $m \in \Omega$,
 $\varepsilon \sim \mathcal{N}(0, I_n)$, de dimension $\dim(\Omega) < n$. Soit ω un sev de Ω
tel que $\dim(\omega) = q < p$. Le test de Fisher

$$H_0 : m \in \omega \text{ contre } H_1 : m \in \Omega \setminus \omega.$$

de région de rejet

$$\mathcal{R}_\alpha = \{F > f_{p-q, n-p, 1-\alpha}\},$$

où $f_{p-q, n-p, 1-\alpha}$ est le quantile de la loi de Fisher
 $\mathcal{F}(p - q, n - p)$, est de niveau α .

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Application : tableau d'analyse de la variance

$(H_0) : \omega$ contre $(H_1) : \Omega \setminus \omega$ tq. $\omega \subset \Omega$, $\dim(\omega) = q$, $\dim(\Omega) = p$

Source	Res.Df	RSS	Df	Sum of Sq	F	Prob > F
ω	$n - q$	$SCR(\omega)$				
Ω	$n - p$	$SCR(\Omega)$	$p - q$	$SCM(\omega, \Omega)$	$f_{obs} = \frac{SCM/(p-q)}{SCR(\Omega)/(n-p)}$	p-value

```
> anova(lm(Y~1,data=df), lm(Y~.,data=df))
```

```
Model 1: Y ~ 1
```

```
Model 2: Y ~ V1+V2+V3 #équivalent à 1+V1+V2+V3
```

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	18	30654				
2	15	260	3	30394	584.55	9.386e-16 ***

```
#F-statistic: 584.6 on 3 and 15 DF, p-value: 9.386e-16
```

Introduction

Préférences

A-Définition et
hypothèses

Exemples

Identifiabilité

B-Estimation

EMC

Propriétés

EMV

C-Lois des
estimateurs

RC

D-Tests

de Student

de Wald

de Fisher

E-Validation

Tableau d'analyse de la variance

Une autre utilisation de la fonction anova

```
> anova(lm(Y~.,data=df))
```

Response: Y

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
V1	1	18252.3	18252.3	1053.11	2.603e-15	***
V2	1	7647.5	7647.5	441.24	1.548e-12	***
V3	1	4494.2	4494.2	259.31	7.101e-11	***
Residuals	15	260.0	17.3			

(H_0)	(H_1)
iid: $Y \sim 1$	$Y \sim V1$
$Y \sim V1$	$Y \sim V1+V2$
$Y \sim V1+V2$	$Y \sim V1+V2+V3$

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Exemple (suite)

```
> anova(lm(Y~.,data=df))
```

```
Response: Y
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
V1	1	18252.3	18252.3	1053.11	2.603e-15	***
V2	1	7647.5	7647.5	441.24	1.548e-12	***
V3	1	4494.2	4494.2	259.31	7.101e-11	***
Residuals	15	260.0	17.3			

```
> anova(lm(Y~V3+V2+V1,data=df))
```

```
Response: y
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)	
V3	1	30377.9	30377.9	1752.7292	<2e-16	***
V2	1	2.3	2.3	0.1349	0.7186	
V1	1	13.8	13.8	0.7976	0.3859	
Residuals	15	260.0	17.3			

Attention à l'ordre des termes quand les variables explicatives ne sont pas orthogonales

Sommaire

Introduction

A-Définition et hypothèses

B-Estimation

C-Lois des estimateurs

D-Tests

E-Validation

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

- ▶ L'étape de validation fait partie intégrante de la démarche statistique
- ▶ Elle peut remettre en cause le modèle initialement choisi, donner des idées de modifications pour repartir pour un tour...
- ▶ Outils : tests, critères, étude graphique
↔ étude des résidus

Introduction

Références

A-Définition et hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Critère du R^2

Si ω est le modèle iid (constant) :

- Coefficient de **détermination**

$$R^2 = \frac{SCM}{SCT} = \frac{\|\hat{Y} - \bar{y}\mathbf{1}\|^2}{\|Y - \bar{y}\mathbf{1}\|^2} = \cos^2 \alpha$$

Part de la variabilité expliquée par le modèle sur la variabilité totale

- Indication de la **qualité d'ajustement**, mais pas forcément de la qualité de prévision
- Pour compte de la dimension de $Im(X)$: $\hookrightarrow R^2$ ajusté :

$$R_a^2 = 1 - \frac{n-1}{n-p} \frac{\|\hat{\epsilon}\|^2}{\|Y - \bar{y}\mathbf{1}\|^2}$$

Multiple R-squared: 0.9915, Adjusted R-squared: 0.9898

Introduction

Références

A-Définition et hypothèses

Exemples

Identifiabilité

B-Estimation

EMC

Propriétés

EMV

C-Lois des estimateurs

RC

D-Tests

de Student

de Wald

de Fisher

E-Validation

- ▶ **Significativité globale** : test de Fisher du modèle iid contre le modèle d'étude.
- ▶ **Significativité d'un régresseur** : \Leftrightarrow des stratégies de sélection de variables.
- ▶ **Adéquation** : comparer le modèle d'étude Ω à un surmodèle Ω_S peu contestable qui contient Ω
- ▶ Hypothèse **gaussienne** : test de Kolmogorov-Smirnov ou test de Shapiro-Wilks.

Introduction

Références

A-Définition et hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des estimateurs

RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Proposition

Dans le modèle linéaire gaussien, les résidus $\hat{\varepsilon} = Y - \hat{Y}$ possèdent les propriétés suivantes :

- ▶ *centrés* : $\mathbf{E}(\hat{\varepsilon}) = 0$.
- ▶ *hétéroscédastiques* : $\text{var}(\hat{\varepsilon}) = \sigma^2(I - H)$.
- ▶ *décorrélés avec les valeurs estimées* : $\text{cov}(\hat{Y}, \hat{\varepsilon}) = 0$.
- ▶ Si $\mathbb{I} \in \text{Im}(X)$, les résidus estimés sont *linéairement dépendants*.

Introduction

Références

A-Définition et
hypothèses

Exemples

Identifiabilité

B-Estimation

EMC

Propriétés

EMV

C-Lois des
estimateurs

RC

D-Tests

de Student

de Wald

de Fisher

E-Validation

Différents résidus

- ▶ Les **résidus normalisés** éliminent l'hétéroscédasticité

$$r_i = \frac{\hat{\varepsilon}_i}{\sigma\sqrt{1-h_{ii}}}.$$

- ▶ **résidus standardisés**, parfois également appelés **studentisés** : de variance 1

$$t_i = \frac{\hat{\varepsilon}_i}{\hat{\sigma}\sqrt{1-h_{ii}}}.$$

- ▶ Les résidus **studentisés par validation croisée** sont définis par

$$t_i^* = \frac{\hat{\varepsilon}_i}{\hat{\sigma}_{(i)}\sqrt{1-h_{ii}}}.$$

où $\hat{\sigma}_{(i)}$ est calculé dans un échantillon privé de l'observation i . Dans le cas gaussien, $t_i^* \sim \mathcal{T}(n-1-p)$.

↪ points **aberrants**

Différents types de graphiques permettent une étude visuelle :

- ▶ valeurs estimées en fonction de valeurs observées,
- ▶ résidus en fonction des valeurs estimées,
- ▶ résidus en fonction d'une valeur de covariable,
- ▶ résidus de l'observation $i + 1$ en fonction du résidu de l'observation i ,
- ▶ graphe quantile/quantile des résidus.

Introduction

Références

A-Définition et
hypothèses

Exemples
Identifiabilité

B-Estimation

EMC
Propriétés
EMV

C-Lois des
estimateurs

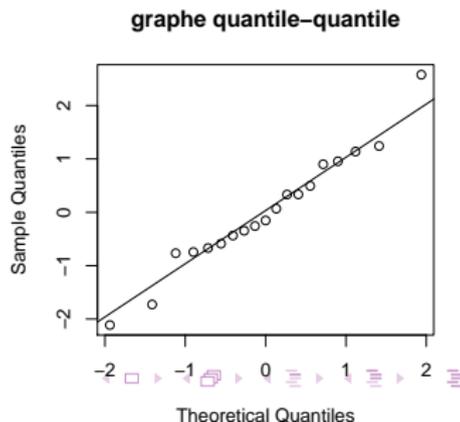
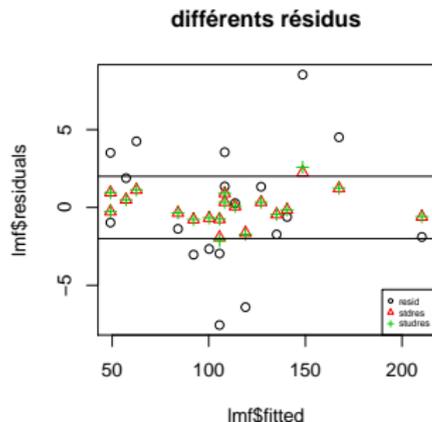
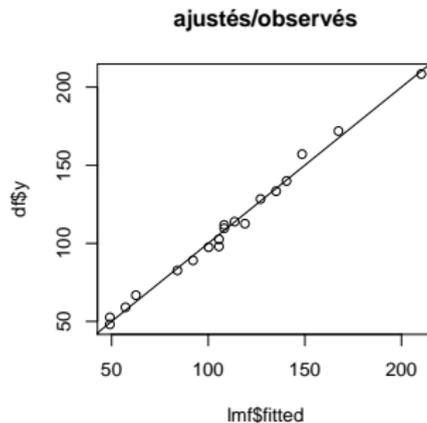
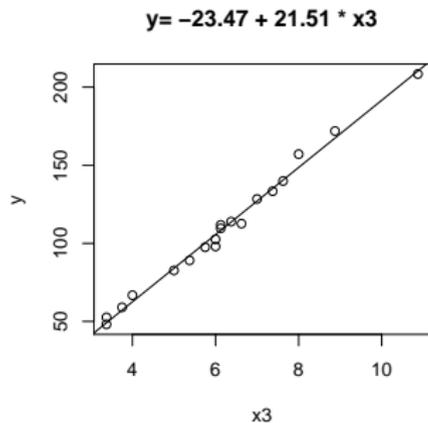
RC

D-Tests

de Student
de Wald
de Fisher

E-Validation

Exemple de diagnostics visuels



Introduction

Références

A-Définition et
hypotheses

Exemples

Identifiabilité

B-Estimation

EMC

Propriétés

EMV

C-Lois des
estimateurs

RC

D-Tests

de Student

de Wald

de Fisher

E-Validation