

Modèles linéaires - Statistical Analysis System (SAS)

Chapitre 2 : Hypothèses

Ordinary Least Squares Estimator (OLS)

Emmanuelle Gautherat^(a)

(a) Crieg-Regards, Université de Reims Champagne Ardenne

Second semestre - 4 ECTS

Outline

1. Objets génériques

Introduction

Méthode OLS fondée sur la corrélation, la moyenne empirique et utilise pleinement l'ensemble des hypothèses

→ très sensible à l'écart à ces hypothèses.

Importance de la [détection des écarts aux hypothèses](#).

→ très sensibles aux valeurs extrêmes ou atypiques.

Importance de la [détection des valeurs atypiques](#) : plusieurs sens donné à 'atypique'.

Ici : cadre du modèle OLS standard avec constante.

Par défaut, l'ensemble des hypothèses est supposé vérifié.

Objets génériques : Matrice H

On note dorénavant $H = (h_{ij})_{i,j=1,\dots,n}$ (prononcer "matrice hat") la matrice de projection $\Pi_{(x)}$.

$$\hat{Y} = \Pi_{(x)} Y = HY.$$

En notant \tilde{x} la matrice dont les colonnes sont les variables $x^{(k)}$ centrées :

$$\tilde{x} = \begin{pmatrix} x_1^{(1)} - \overline{x^{(1)}} & \dots & x_1^{(K)} - \overline{x^{(K)}} \\ \vdots & & \vdots \\ x_n^{(1)} - \overline{x^{(1)}} & \dots & x_n^{(K)} - \overline{x^{(K)}} \end{pmatrix},$$

on obtient

$$\begin{aligned} \hat{Y}_i &= (HY)_i = \sum_{j=1}^n h_{ij} Y_j = \sum_{j=1}^n \left(\frac{1}{n} + \tilde{x} (\tilde{x}^t \tilde{x})^{-1} \tilde{x}^t \right)_{ij} Y_j \\ h_{ii} &= \frac{1}{n} + \sum_{k=1}^K \left(\frac{(v^k)^t (x_i - \bar{x})}{\sqrt{\lambda_k}} \right)^2 \end{aligned}$$

avec λ_k val. propre associée au vect. propre normé v^k de $\tilde{x}^t \tilde{x}$.

R : Comme $H = \Pi_{(x)}$, on obtient la diagonale de la matrice H avec

`diag(x%*% solve(t(x) %*% x) %*% t(x))`

Ou encore

`hat(x).`

Estimateur du niveau de bruit $\hat{\sigma}_{(-i)}$

On rappelle que \hat{u}_i , le **résidu** pour l'observation i , est défini par

$$\begin{aligned}\hat{u} = \Pi_{x^\perp} Y &= (Id_n - H)Y \sim \mathcal{N}_n(0 ; \sigma^2(Id - H)) \\ \hat{u}_i &\sim \mathcal{N}_1(0 ; \sigma^2(1 - h_{ii})).\end{aligned}$$

Estimateur du niveau de bruit sans l'observation i : $\hat{\sigma}_{(-i)}^2$.

•

$$\hat{\sigma}^2 = \frac{\|\hat{u}\|_2^2}{n - K - 1} = \frac{\hat{u}^t \hat{u}}{n - K - 1}.$$

- Soit $\hat{\sigma}_{(-i)}^2$ un estimateur du niveau de bruit calculé sans utiliser l'observation i .
C'est-à-dire :

$$\begin{aligned}\hat{\sigma}_{(-i)}^2 &= \frac{\hat{u}_{(-i)}^t \hat{u}_{(-i)}}{n - K - 2} \\ \hat{\sigma}_{(-i)}^2 &\sim \chi^2((n - 1) - (K + 1)) = \chi^2(n - K - 2)\end{aligned}$$

.

Résidus standardisés \hat{e}_i , résidus studentisés \hat{t}_i

Rappel : Pour tout i , $\frac{\hat{u}_i}{\sqrt{\sigma^2(1-h_{ii})}} \sim \mathcal{N}_1(0; 1)$.

Résidus standardisés pour l'observation i : \hat{e}_i

On définit \hat{e}_i par :

$$\hat{e}_i = \frac{\hat{u}_i}{\sqrt{\hat{\sigma}^2(1-h_{ii})}} \sim T(n-K-1).$$

Résidus studentisés pour l'observation i : \hat{t}_i

On définit \hat{t}_i par :

$$\hat{t}_i = \frac{\hat{u}_i}{\sqrt{\hat{\sigma}_{(-i)}^2(1-h_{ii})}} \sim T(n-K-2).$$

Avec R :

Pour `reg<-lm(.....)`

`\hat{u}` \longrightarrow `residuals(reg)`

`\hat{e}` \longrightarrow `rstandard(reg)`

`\hat{t}` \longrightarrow `rstudent(reg)`

`$\hat{\sigma}_{(-i)}^2$` \longrightarrow `Influ2<-lm.influence(reg); Influ2$sigma`

Prédiction sans l'individu "i" : $\hat{Y}_{(-i)}$ et $\hat{\beta}_{(-i)}$

On désigne par $\hat{Y}_{(-i)}$ la prédiction de Y lorsque l'individu i a été omis du calcul de $\hat{\beta}$.

Soit $x_{(-i)}$, la matrice x à laquelle on a retiré la ligne i : $x_{(-i)} \in \mathcal{M}_{n-1, K+1}(\mathbb{R})$.

Soit $Y_{(-i)}$ le vecteur des endogènes auquel on a retiré l'observation i : $Y_{(-i)} \in \mathcal{M}_{n-1, 1}(\mathbb{R})$

On remarque que $Y_{(-i)}$ a bien n valeurs : la valeur de Y pour i est prédite à partir des autres individus.

On définit :

$$\hat{\beta}_{(-i)} = (x_{(-i)}^t x_{(-i)})^{-1} x_{(-i)}^t Y_{(-i)}$$

$$\hat{Y}_{(-i)} = x \hat{\beta}_{(-i)}.$$

On a bien

$$\hat{\beta}_{(-i)} \in \mathcal{M}_{K+1, 1}(\mathbb{R}).$$