

# ÉCONOMETRIE THÉORIQUE

## Chapitre 1

### L'estimation par moindres carrés ordinaires (MCO)

*Benoît Mulkey  
Université de Montpellier  
2023 - 2024*

1

## Chapitre I :

### L'estimation par moindres carrés ordinaires (MCO)

- I.1. Le modèle de régression linéaire classique.
- I.2. L'estimation par Moindres Carrés Ordinaires (MCO).
- I.3. L'estimateur de la variance de l'erreur.
- I.4. La qualité de l'ajustement.
- I.5. Interprétation géométrique des MCO.
- I.6. Le théorème de Frisch – Waugh.
- I.7. Les observations influentes

2

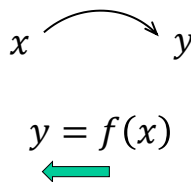
## I.1. Le modèle de régression linéaire classique

### a) Le modèle de régression multiple

HAYASHI [2000], Chapitre I.

L'analyse de régression :

- Modélisation d'une relation pour « expliquer » une **variable dépendante** (régressand), notée  $y$ ,
- par des **variables explicatives** (régresseurs), notées  $x$



#### Attention à la notion de causalité

supposée par hypothèse dans cette notation !!!

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

3

3

## L'échantillon ou les données

- C'est un ensemble d'observations sur plusieurs variables ou caractéristiques.
- Echantillon de  $N$  observations indicées :  $i = 1, 2, 3, \dots, N$ ;
  - Une **coupe instantanée** (données longitudinales) est observée sur différentes unités statistiques à une même période.
  - Elle est en général non ordonnée !
- Echantillon de  $T$  observations sur séries temporelles :  $t = 1, 2, 3, \dots, T$ .
  - Une **série temporelle** (ou chronologique) est observée sur une même unité statistique pour plusieurs périodes de temps (année, trimestres, jours,...)
  - Une série temporelle a un ordre naturel : passé  $\rightarrow$  présent  $\rightarrow$  futur

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

4

4

## Les observations

- Ici, une seule variable dépendante ( $y_i$ )
- et plusieurs variables explicatives ( $x_{k,i}$ ) :  $k = 1, 2, 3, \dots, K$  :

$$\mathbf{x}_i = (x_{1,i}, x_{2,i}, \dots, x_{k,i}, \dots, x_{K,i})'$$

- En général, la première variable explicative est constante :  $x_{1,i} = 1$
- Dans la régression simple :  $K = 2 \rightarrow \mathbf{x}_i = (1, x_{2,i})'$
- Dans la régression multiple :  $K > 2 \rightarrow \mathbf{x}_i = (1, x_{2,i}, \dots, x_{k,i}, \dots, x_{K,i})'$
- Régression\* de  $y$  sur les variables explicatives (régresseurs)  $x_k$

\* : Francis GALTON (1886) : "Regression Towards Mediocrity in Hereditary Stature", *The Journal of the Anthropological Institute of Great Britain and Ireland*, 15: 246–263.

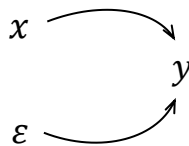
GALTON (1822 – 1911) était le cousin de Charles Darwin. Il a été un fondateur de la Société d'Eugénisme (Eugenics Society) qui voulait améliorer la qualité génétique des hommes.

## Le modèle

- La variable dépendante ( $y_i$ ) et les variables explicatives ( $x_{k,i}$ ) sont considérées comme des **variables aléatoires**.
- Les observations sont considérées comme des **réalisations** d'un tirage aléatoire !
- Ici les régresseurs sont aussi des variables aléatoires parce que :
  - On ne contrôle pas les variables explicatives.
  - En sciences sociales et plus particulièrement en économie, on est rarement dans le cadre d'expériences contrôlées (*randomisée*), mais plutôt avec des *données observationnelles*.
  - Le cas des régresseurs « fixes » est un cas particulier du cas des régresseurs aléatoires.

## Un modèle

- Un **modèle économétrique** est une ensemble de restrictions sur la distribution conjointe des variables dépendante et explicatives.
- C'est un ensemble de distributions conjointes qui satisfont un ensemble d'hypothèses.
- Comme la relation n'est pas parfaite, on introduit une variable supplémentaire pour tenir compte d'autres éléments → **l'erreur** :  $\varepsilon$  !



7

## Commentaires sur l'erreur : $\varepsilon_i$ (error , disturbance)

- L'erreur est aussi appelée aléas ou perturbation.
- Elle correspond à des éléments imprédictibles des aléas des comportements économiques, humains ou sociaux...
- Elle mesure les incertitudes du modèle.
- Elle incorpore des effets inobservables.
- Elle capture les effets d'un grand nombre de variables omises (*voir les hypothèses ci-après*).
- Elle prend en compte les erreurs de mesure sur la variable dépendante.
- L'erreur est une variable aléatoire avec des moments (espérance, variance,...) et une distribution à déterminer.

8

## HYPOTHÈSE 1 : Linéarité du modèle

$$y_i = \beta_1 + \beta_2 x_{2,i} + \beta_3 x_{3,i} + \dots + \beta_K x_{K,i} + \varepsilon_i \quad \text{pour } i = 1, 2, \dots, N.$$

- La relation entre les variables explicatives et la variable dépendante est **linéaire**.
- L'erreur est **additive**.
- Il n'y a pas d'erreurs de mesure sur les variables explicatives  $x$ .
- → **Le modèle est correctement spécifié !**

## HYPOTHÈSE 1 : Linéarité du modèle

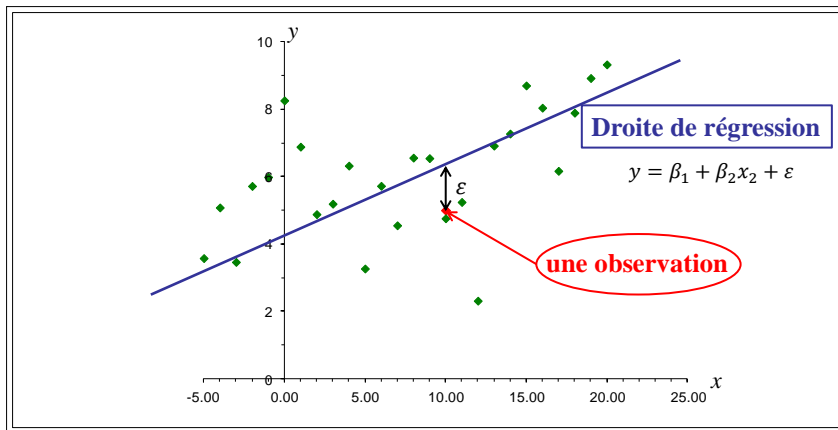
$$y_i = \beta_1 + \beta_2 x_{2,i} + \beta_3 x_{3,i} + \dots + \beta_K x_{K,i} + \varepsilon_i = \sum_{k=1}^K \beta_k x_{k,i} + \varepsilon_i$$

- $\beta_1 + \beta_2 x_{2,i} + \beta_3 x_{3,i} + \dots + \beta_K x_{K,i}$  est la **fonction de régression**.
- Les paramètres (ou coefficients)  $\beta_k$  sont les **coefficients de régression**.
- Ils représentent **l'effet marginal**, *toutes choses égales par ailleurs*, d'une variation de la variable  $x_{k,i}$  sur la variable dépendante  $y_i$  :

$$\beta_k = \frac{\partial y_i}{\partial x_{k,i}}$$

- La linéarité implique que les effets marginaux sont **constants** pour toutes les observations (individus).

### Cas de la régression simple : la droite de régression



### Cas de la régression multiple : un (hyper-)plan de régression

11

### Commentaires sur la forme linéaire

$$y_i = \beta_1 + \beta_2 x_{2,i} + \beta_3 x_{3,i} + \dots + \beta_K x_{K,i} + \varepsilon_i$$

Forme linéaire comme une **approximation locale** d'une forme non-linéaire inconnue...

ou comme une **projection linéaire** ...

Modèle linéaire dans les paramètres, pas nécessairement dans les variables.

Celles-ci peuvent être transformées pour retrouver la linéarité.

Exemple : **le modèle double log**

$$y = \lambda x^\beta \eta \rightarrow y = \exp(\alpha) \times x^\beta \times \exp(\varepsilon) \quad \text{avec} \begin{cases} \lambda = \exp(\alpha) \\ \eta = \exp(\varepsilon) \end{cases}$$

$$\rightarrow \log(y) = \alpha + \beta \log(x) + \varepsilon$$

$$\beta = \frac{\partial \log(y)}{\partial \log(x)} = \frac{\partial y/y}{\partial x/x} = \frac{\partial y}{\partial x} \times \frac{x}{y}$$

Effet marginal

→ élasticité de  $y$  (par rapport) à  $x$ .

12

Autre exemple : **le modèle polynomial** :

$$y = \alpha + \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \dots + \beta_p x^p + \varepsilon$$

$$\frac{\partial y}{\partial x} = \beta_1 + 2\beta_2 x + 3\beta_3 x^2 + \dots + p\beta_p x^{p-1}$$

Ici l'effet marginal dépend de la valeur de la variable explicative  $x$  !

Autre exemple : **le modèle à paramètres variables ou avec interaction** :

$$y = (\alpha_1 + \alpha_2 z) + (\beta_1 + \beta_2 z)x + \varepsilon \quad \rightarrow \quad y = \alpha_1 + \alpha_2 z + \beta_1 x + \beta_2 zx + \varepsilon$$

$$\frac{\partial y}{\partial x} = \beta_1 + \beta_2 z$$

Ici l'effet marginal dépend de la valeur de l'autre variable explicative  $z$  !

Choix théorique et empirique pour la spécification du modèle

Il faut faire attention à l'interprétation des paramètres : effets marginaux, élasticités, effets constants ou variables,...

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

13

**13**

**Extension aux modèles non linéaires :**

$$y_i = f(x_{1,i}, x_{2,i}, x_{3,i}, \dots, x_{K,i} ; \theta) + \varepsilon_i$$

$$g(y_i, x_{1,i}, x_{2,i}, x_{3,i}, \dots, x_{K,i} ; \theta) = \varepsilon_i$$

Modèles paramétriques avec  $f(\cdot)$  et / ou  $g(\cdot)$  connues

Modèles non paramétriques : on doit aussi déterminer la **forme** des fonctions  $f(\cdot)$  ou  $g(\cdot)$

→ voir un cours d'économétrie non paramétrique...

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

14

**14**

### **b) Notation matricielle du modèle de régression multiple**

Utilisation des vecteurs et des matrices

Simplification de la notation → on évite les sommes...

Les données sont des tableaux de nombres

Propriétés du calcul matriciel (algèbre linéaire)

→ simplification des calculs et des démonstrations...

Pour une observation  $i = 1, 2, \dots, N$  :

$$y_i = \beta_1 + \beta_2 x_{2,i} + \beta_3 x_{3,i} + \dots + \beta_K x_{K,i} + \varepsilon_i = \sum_{k=1}^K \beta_k x_{k,i} + \varepsilon_i$$

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i \quad \text{avec } \mathbf{x}_i \text{ et } \boldsymbol{\beta} \text{ des vecteurs (colonnes)} : K \times 1$$

$\mathbf{x}_i' \boldsymbol{\beta}$  → produit scalaire de 2 vecteurs de même dimension.

$$y_i = \mathbf{x}_i' \boldsymbol{\beta} + \varepsilon_i$$

$$\mathbf{x}_i : \text{un vecteur } (K \times 1) \text{ des } \underline{K \text{ variables explicatives}} : \mathbf{x}_i = \begin{pmatrix} 1 \\ x_{2,i} \\ x_{3,i} \\ \vdots \\ x_{K,i} \end{pmatrix}$$

$$\boldsymbol{\beta} : \text{un vecteur } (K \times 1) \text{ des } \underline{K \text{ paramètres inconnus}} : \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \vdots \\ \beta_K \end{pmatrix}$$

En empilant les  $N$  observations de l'échantillon, on obtient le modèle de régression multiple en notation matricielle :

$$\begin{cases} y_1 = \mathbf{x}_1' \boldsymbol{\beta} + \varepsilon_1 \\ y_2 = \mathbf{x}_2' \boldsymbol{\beta} + \varepsilon_2 \\ \vdots \\ y_N = \mathbf{x}_N' \boldsymbol{\beta} + \varepsilon_N \end{cases} \rightarrow \mathbf{y} = \mathbf{X} \boldsymbol{\beta} + \boldsymbol{\varepsilon}$$



$$\underset{(N \times 1)}{\mathbf{y}} = \underset{(N \times 1)}{\underbrace{\underset{(N \times K)}{\mathbf{X}} \underset{(K \times 1)}{\boldsymbol{\beta}}}} + \underset{(N \times 1)}{\boldsymbol{\varepsilon}}$$

$\mathbf{y}$  et  $\boldsymbol{\varepsilon}$  : des vecteurs  $(N \times 1)$  des variables dépendantes et des erreurs :

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} \quad \text{et} \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_N \end{pmatrix}$$

et  $\mathbf{X}$  : une matrice  $(N \times K)$  des  $K$  variables explicatives :

$$\mathbf{X} = \begin{pmatrix} x'_1 \\ x'_2 \\ \vdots \\ x'_N \end{pmatrix} = \left[ \begin{array}{cccccc} 1 & x_{2,1} & x_{3,1} & \cdots & x_{K,1} \\ 1 & x_{2,2} & x_{3,2} & \cdots & x_{K,2} \\ \vdots & \vdots & \ddots & & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{2,N} & x_{3,N} & \cdots & x_{K,N} \end{array} \right] \left. \vphantom{\begin{pmatrix} x'_1 \\ x'_2 \\ \vdots \\ x'_N \end{pmatrix}} \right\} \begin{array}{l} N \text{ observations} \\ \\ K \text{ variables} \end{array}$$

et toujours  $\boldsymbol{\beta}$  : un vecteur  $(K \times 1)$  des  $K$  paramètres inconnus :  $\boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \beta_3 \\ \vdots \\ \beta_K \end{pmatrix}$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

17

17

## **HYPOTHÈSE 2 : Stricte exogénéité des régresseurs (H2)**

$$\begin{array}{ll} E(\varepsilon_i | \mathbf{X}) = 0 & \text{pour tout } i = 1, 2, \dots, N. \\ E(\varepsilon_i | x_1, x_2, \dots, x_N) = 0 & \text{pour tout } i = 1, 2, \dots, N. \\ E(\varepsilon_i | x_j) = 0 & \text{pour tout } i \text{ et } j = 1, 2, \dots, N. \end{array}$$

L'erreur doit être **exogène** par rapport à **toutes** les observations de **toutes** les variables explicatives.

Si on considère la distribution conjointe des  $(NK + N)$  variables aléatoires :  $f(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N, x_1, x_2, \dots, x_N)$  pour l'ensemble de l'échantillon, la distribution conditionnelle pour chaque observation  $i$  :  $f(\varepsilon_i | x_1, x_2, \dots, x_N)$  aura une espérance nulle :  $E(\varepsilon_i | x_1, x_2, \dots, x_N) = E(\varepsilon_i | \mathbf{X}) = 0$ .

Cette hypothèse est **fondamentale** pour qu'on puisse trouver un estimateur des paramètres  $\boldsymbol{\beta}$  qui possède de bonnes propriétés statistiques.

On reviendra sur cette hypothèse dans le Chapitre VI.

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

18

18

### Implications de l'hypothèse d'exogénéité stricte :

- L'espérance inconditionnelle de l'erreur est nulle

$$E(\varepsilon_i) = 0 \quad \text{pour tout } i = 1, 2, \dots, N.$$

par la loi de l'espérance itérée :  $E(\varepsilon_i) = E_X[E(\varepsilon_i|\mathbf{X})] = 0$

- Les régresseurs sont orthogonaux aux termes d'erreurs pour toutes les observations.

*Deux variables aléatoires sont dites orthogonales si et seulement si l'espérance de leur produit (scalaire) est égale à zéro !*

$$E(x_{j,k}\varepsilon_i) = 0 \quad \text{pour tout } i, j = 1, 2, \dots, N \text{ et } k = 1, 2, \dots, K.$$

$$E(\mathbf{x}_j \varepsilon_i) = \begin{pmatrix} E(x_{j,1}\varepsilon_i) \\ E(x_{j,2}\varepsilon_i) \\ \vdots \\ E(x_{j,K}\varepsilon_i) \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{pmatrix} = \mathbf{0}_K \quad \text{pour tout } i, j = 1, 2, \dots, N.$$

### DÉMONSTRATION

Comme  $x_{j,k}$  est un élément de  $\mathbf{X}$ , la stricte exogénéité implique :

$$E(\varepsilon_i|x_{j,k}) = E[E(\varepsilon_i|\mathbf{X})|x_{j,k}] = 0$$

du fait de la loi des espérances itérées :  $E[E(y|x, z)|x] = E(y|x)$

On aura alors avec la loi de l'espérance totale :

$$E(x_{j,k}\varepsilon_i) = E[E(x_{j,k}\varepsilon_i|x_{j,k})] = E[x_{j,k}E(\varepsilon_i|x_{j,k})] = 0$$

L'avant-dernière égalité s'obtient avec la propriété de linéarité des espérances conditionnelles :  $E[f(x)y|x] = f(x)E(y|x)$

**CQFD**

On parlera alors de **conditions d'orthogonalité** pour cette hypothèse 2 :

$$E(X'\varepsilon) = \mathbf{0}_K$$

- **Les conditions d'orthogonalité sont équivalentes à des conditions d'absence de corrélation.**

Avec la définition de la covariance et en utilisant l'espérance inconditionnelle de  $\varepsilon_i$  :

$$\text{Cov}(\varepsilon_i, x_{j,k}) = E(x_{j,k}\varepsilon_i) - E(x_{j,k})E(\varepsilon_i) = E(x_{j,k}\varepsilon_i) = 0$$

du fait des conditions d'orthogonalité précédentes.

- **L'hypothèse de stricte exogénéité n'est pas souvent satisfaite pour les modèles sur séries temporelles...**

Il se peut que la variable  $y$  au temps  $t$  influence en retour les variables explicatives  $X$  dans les périodes suivantes.

$$\text{Corr}(y_t, x_{t+s}) \neq 0 \quad \text{pour } s > 0$$

**Exemple :** Considérons le modèle (très courant) autorégressif simple en série temporelle, c'est-à-dire qu'il y a un seul régresseur, la variable dépendante retardée :

$$y_t = \beta y_{t-1} + \varepsilon_t$$

Supposons que le régresseur soit orthogonal à l'erreur :  $E(y_{t-1}\varepsilon_t) = 0$ , on aura alors :

$$E(y_t \varepsilon_t) = E((\beta y_{t-1} + \varepsilon_t) \varepsilon_t) = \beta E(y_{t-1} \varepsilon_t) + E(\varepsilon_t^2) = E(\varepsilon_t^2)$$

qui n'est pas nulle, sinon toutes les erreurs  $\varepsilon_t$  seraient nulles.

Mais  $y_t$  est le régresseur pour l'observation  $t + 1$ . Celui-ci n'est donc pas orthogonal à l'erreur de la période précédente. Ce qui viole l'hypothèse de stricte exogénéité.

### **HYPOTHÈSE 3 : Absence de multicollinéarité parfaite (H3)**

Le rang de la matrice  $\mathbf{X}$  de dimension  $N \times K$  est  $K$  avec une probabilité 1 :

$$\Rightarrow \text{rang}(\mathbf{X}) = K$$

Cela veut dire qu'aucune colonne de la matrice  $\mathbf{X}$  (*aucune variable*) n'est combinaison linéaire parfaite des autres colonnes de cette matrice (*des autres variables*).

→  $\mathbf{X}$  est de rang-plein (colonne)

C'est une condition suffisante

Cette hypothèse permet l'identification des paramètres du modèle, c'est-à-dire cela permet d'obtenir une estimation unique des paramètres du modèle avec un critère donné.

Cette hypothèse (technique) rend possible l'estimation par moindres carrés ordinaires.

**23**

Une condition nécessaire mais pas suffisante est qu'il y ait au moins autant d'observations que de paramètres à estimer :  $N \geq K$ .

Les régresseurs sont dits « **parfaitement colinéaires** » si cette hypothèse n'est pas satisfaite. Cela se remarque souvent très facilement...

**24**

#### **HYPOTHÈSE 4 : Les erreurs ont une matrice de variance-covariance sphérique (H4)**

- a) **Homoscédasticité conditionnelle** : Les termes d'erreur  $\varepsilon_i$  de toutes les observations ont la même variance conditionnelle :

$$V(\varepsilon_i|\mathbf{X}) = \sigma^2 \quad \text{pour tout } i = 1, 2, \dots, N.$$

Cette variance est strictement positive et finie :  $0 < \sigma^2 < \infty$

On peut réécrire cette hypothèse comme :  $V(\varepsilon_i|\mathbf{X}) = E(\varepsilon_i^2|\mathbf{X}) = \sigma^2$

---

$$\begin{aligned} V(\varepsilon_i|\mathbf{X}) &= E\left[(\varepsilon_i - E(\varepsilon_i|\mathbf{X}))^2|\mathbf{X}\right] && \text{par définition de la variance} \\ &= E(\varepsilon_i^2|\mathbf{X}) - \underbrace{E(\varepsilon_i|\mathbf{X})^2}_{=0} && \text{parce que } E(\varepsilon_i|\mathbf{X}) = 0 \text{ (H2)} \\ &= E(\varepsilon_i^2|\mathbf{X}) = \sigma^2 \end{aligned}$$

---

- b) **Absence de corrélation entre les observations**

Il y a **indépendance** « statistique » entre les erreurs.

En fait, on a besoin d'une hypothèse **plus faible** d'absence de corrélation entre les erreurs :

$$E(\varepsilon_i \varepsilon_j|\mathbf{X}) = 0 \quad \text{pour tout } i, j = 1, 2, \dots, N \text{ et } i \neq j.$$

En conséquence les erreurs ne sont pas corrélées entre elles.

$$\begin{aligned} Cov(\varepsilon_i, \varepsilon_j|\mathbf{X}) &= E(\varepsilon_i \varepsilon_j|\mathbf{X}) - E(\varepsilon_i|\mathbf{X})E(\varepsilon_j|\mathbf{X}) = 0 \\ Corr(\varepsilon_i, \varepsilon_j|\mathbf{X}) &= \frac{Cov(\varepsilon_i, \varepsilon_j|\mathbf{X})}{\sqrt{V(\varepsilon_i|\mathbf{X}) \times V(\varepsilon_j|\mathbf{X})}} = \frac{E(\varepsilon_i \varepsilon_j|\mathbf{X})}{\sqrt{Var(\varepsilon_i|\mathbf{X})} \sqrt{Var(\varepsilon_j|\mathbf{X})}} = \frac{0}{\sigma^2} = 0 \end{aligned}$$

Conditionnellement à  $\mathbf{X}$ , il y a **non corrélation** entre les erreurs, et donc entre les observations !

→ **absence d'autocorrélation**

Cette dernière hypothèse (H4) peut se réécrire en langage matriciel, sous la forme d'une **matrice de variance-covariance** (carrée de dimension  $N \times N$ ) :

$$V(\boldsymbol{\varepsilon}|\mathbf{X}) = E[(\boldsymbol{\varepsilon} - E(\boldsymbol{\varepsilon}|\mathbf{X}))(\boldsymbol{\varepsilon} - E(\boldsymbol{\varepsilon}|\mathbf{X}))'|\mathbf{X}]$$

$$V(\boldsymbol{\varepsilon}|\mathbf{X}) = E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}'|\mathbf{X}] = \begin{bmatrix} \sigma^2 & 0 & 0 & \dots & 0 \\ 0 & \sigma^2 & 0 & \dots & 0 \\ 0 & 0 & \sigma^2 & \dots & \vdots \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \sigma^2 \end{bmatrix} = \sigma^2 \mathbf{I}_N$$

La matrice identité de dimension  $(N \times N)$ .

Cette matrice de variance-covariance des erreurs est alors dite « **sphérique** ».

Cette matrice de variance-covariance est non-singulière parce que  $\sigma^2 > 0$ , donc elle est inversible.

Elle est aussi définie positive avec un déterminant positif.

27

**DÉMONSTRATION :**  $V(\boldsymbol{\varepsilon}) = E[(\boldsymbol{\varepsilon} - E(\boldsymbol{\varepsilon}))(\boldsymbol{\varepsilon} - E(\boldsymbol{\varepsilon}))'] = E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}']$   
(en omettant le signe de conditionnalisation par rapport à  $\mathbf{X}$ )

$$V(\boldsymbol{\varepsilon}) = E[(\boldsymbol{\varepsilon} - E(\boldsymbol{\varepsilon}))(\boldsymbol{\varepsilon} - E(\boldsymbol{\varepsilon}))'] = E[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}']$$

$$= E \left[ \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_N \end{pmatrix} (\varepsilon_1 \quad \varepsilon_2 \quad \dots \quad \varepsilon_N)' \right] = E \begin{bmatrix} \varepsilon_1^2 & \varepsilon_1 \varepsilon_2 & \dots & \varepsilon_1 \varepsilon_N \\ \varepsilon_2 \varepsilon_1 & \varepsilon_2^2 & \dots & \varepsilon_2 \varepsilon_N \\ \vdots & \vdots & \ddots & \vdots \\ \varepsilon_N \varepsilon_1 & \varepsilon_N \varepsilon_2 & \dots & \varepsilon_N^2 \end{bmatrix}$$

$$= \begin{bmatrix} E(\varepsilon_1^2) & E(\varepsilon_1 \varepsilon_2) & \dots & E(\varepsilon_1 \varepsilon_N) \\ E(\varepsilon_2 \varepsilon_1) & E(\varepsilon_2^2) & \dots & E(\varepsilon_2 \varepsilon_N) \\ \vdots & \vdots & \ddots & \vdots \\ E(\varepsilon_N \varepsilon_1) & E(\varepsilon_N \varepsilon_2) & \dots & E(\varepsilon_N^2) \end{bmatrix}$$

$$= \begin{bmatrix} \sigma^2 & 0 & \dots & 0 \\ 0 & \sigma^2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \sigma^2 \end{bmatrix}$$

← Une matrice proportionnelle à la matrice identité de dimension  $(N \times N)$ .

**Avec les hypothèses H4a : Homoscédasticité :**  $V(\varepsilon_i) = E(\varepsilon_i^2) = \sigma^2$

**et H4b : Non –autocorrélation des erreurs :**  $Cov(\varepsilon_i, \varepsilon_j) = E(\varepsilon_i \varepsilon_j) = 0$

28

## HYPOTHÈSE 5 : Normalité des erreurs (H5)

Cette hypothèse n'est pas fondamentale pour l'obtention d'un bon estimateur, mais elle est **nécessaire** pour obtenir des propriétés en petits échantillons...

Si on suppose **les observations indépendantes**, la fonction de densité conjointe des erreurs est le produit des fonctions de densité marginale de chaque erreur :

$$f(\boldsymbol{\varepsilon}|\mathbf{X}) = f(\varepsilon_1, \varepsilon_2, \dots, \varepsilon_N|\mathbf{X}) = f(\varepsilon_1|\mathbf{X}) \times f(\varepsilon_2|\mathbf{X}) \times \dots \times f(\varepsilon_N|\mathbf{X}) = \prod_{i=1}^N f(\varepsilon_i|\mathbf{X})$$

avec  $f(\cdot)$  : la fonction de densité de la **loi normale univariée** :

$$\varepsilon_i|\mathbf{X} \approx \mathcal{N}(0, \sigma^2) \quad \rightarrow \quad f(\varepsilon_i|\mathbf{X}) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{\varepsilon_i^2}{2\sigma^2}\right\}$$

Dans ce cas, les erreurs suivent une loi normale multivariée de dimension  $N$ :

$$\boldsymbol{\varepsilon}|\mathbf{X} \approx \mathcal{N}_N(\mathbf{0}, \sigma^2 \mathbf{I}_N) \quad \rightarrow \quad f(\boldsymbol{\varepsilon}|\mathbf{X}) = (2\pi)^{-N/2} (\sigma^2)^{-N/2} \exp\left\{-\frac{\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}}{2\sigma^2}\right\}$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

29

29

La fonction de densité des erreurs de la loi normale multivariée s'écrit alors :

$$\begin{aligned} f(\boldsymbol{\varepsilon}|\mathbf{X}) &= \prod_{i=1}^N f(\varepsilon_i|\mathbf{X}) = \prod_{i=1}^N \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{\varepsilon_i^2}{2\sigma^2}\right\} \right] \\ &= \prod_{i=1}^N \left[ (2\pi)^{-1/2} (\sigma^2)^{-1/2} \exp\left\{-\frac{\varepsilon_i^2}{2\sigma^2}\right\} \right] \\ &= (2\pi)^{-N/2} (\sigma^2)^{-N/2} \prod_{i=1}^N \exp\left\{-\frac{\varepsilon_i^2}{2\sigma^2}\right\} \\ &= (2\pi)^{-N/2} (\sigma^2)^{-N/2} \exp\left\{\sum_{i=1}^N \left(-\frac{\varepsilon_i^2}{2\sigma^2}\right)\right\} \\ &= (2\pi)^{-N/2} (\sigma^2)^{-N/2} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^N \varepsilon_i^2\right\} \\ \text{(en langage matriciel)} \quad &= (2\pi)^{-N/2} (\sigma^2)^{-N/2} \exp\left\{-\frac{\boldsymbol{\varepsilon}'\boldsymbol{\varepsilon}}{2\sigma^2}\right\} \quad \text{avec} \quad \sum_{i=1}^N \varepsilon_i^2 = \boldsymbol{\varepsilon}'\boldsymbol{\varepsilon} \end{aligned}$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

30

30

## L2. L'estimation par Moindres Carrés Ordinaires (MCO).

### a) Le critère des moindres carrés

Comment estimer les paramètres inconnus  $\beta$  et  $\sigma^2$  ?

Utilisation de l'information de l'échantillon

Définir un critère ou un objectif à atteindre

Trouver une règle de décision ou de choix

Règle de décision → Un estimateur

Mise en œuvre – Faisabilité : L'estimation des paramètres du modèle

**Critère des moindres carrés (Least-squares) :**

→ *Minimisation de la somme des carrés des erreurs (Q) :*

$$\min_{\beta \in \mathbb{R}^k} Q(\beta) = \sum_{i=1}^N \varepsilon_i^2 = \sum_{i=1}^N (y_i - x_i' \beta)^2$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

31

31

$$\min_{\beta \in \mathbb{R}^k} Q(\beta) = \sum_{i=1}^N \varepsilon_i^2 = \sum_{i=1}^N (y_i - x_i' \beta)^2$$

### Avantages de ce critère :

- Donne les paramètres de la droite ou du plan de régression de  $y$  sur  $x$
- Facilité de mise en œuvre (règle de décision → estimateur linéaire)
- Bonnes propriétés statistiques
- Mais forte pénalité aux erreurs « importantes »
- Pas de contraintes sur les paramètres estimés !

D'autres critères sont possibles. Par exemple :

- Minimiser la valeur absolue
  - Minimiser le coefficient de Gini des erreurs
- voir des cours d'économétrie plus avancés...

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

32

32



### Propriétés du critère des moindres carrés $Q$ :

$Q(\beta)$  est une fonction scalaire qui dépend d'un vecteur ( $K \times 1$ ) de paramètres  $\beta$  :

$$Q(\beta) : \mathbb{R}^K \rightarrow \mathbb{R}^+$$

$Q(\beta)$  est une forme quadratique en  $\beta$  :

$$Q(\beta) = \sum_{i=1}^N \varepsilon_i^2 = \varepsilon' \varepsilon \quad \text{avec} \quad \varepsilon = y - X\beta$$

qui peut se réécrire comme :

$$\begin{aligned} Q(\beta) &= \varepsilon' \varepsilon \\ &= (y - X\beta)'(y - X\beta) = y'y - 2y'X\beta + \beta'X'X\beta \\ &= y'y - y'X\beta - \beta'X'y + \beta'X'X\beta \\ &= y'y - 2y'X\beta + \beta'X'X\beta \end{aligned}$$

parce que  $y'X\beta$  est un scalaire, et donc :  $y'X\beta = (y'X\beta)' = \beta'X'y$ .

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

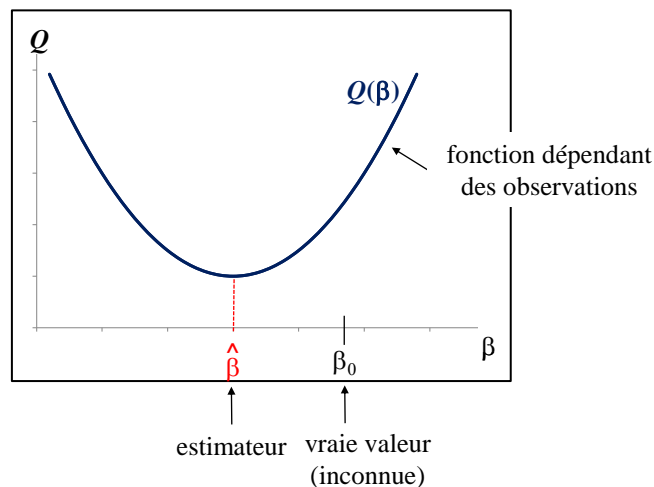
33

33

### Propriétés du critère des moindres carrés $Q$ :

$Q(\beta)$  a une forme en U avec un minimum unique (*voir plus loin*).

Dans le cas où il y a un seul paramètre :



Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

34

34

### b) L'estimateur des MCO

**Le critère** : *Minimisation de la somme des carrés des erreurs*

$$\min_{\beta \in \mathbb{R}^K} Q(\beta) = \sum_{i=1}^N \varepsilon_i^2 = \sum_{i=1}^N (y_i - x_i' \beta)^2$$

La solution de ce problème de minimisation est l'estimateur des moindres carrés :

$$\hat{\beta} = \arg \min_{\beta \in \mathbb{R}^K} Q(\beta)$$

On peut décomposer la forme quadratique  $Q(\beta)$  avec la notation de sommation :

$$\begin{aligned} \sum_{i=1}^N \varepsilon_i^2 &= \sum_{i=1}^N \left( y_i - \sum_{k=1}^K \beta_k x_{k,i} \right)^2 \\ &= \sum_{i=1}^N y_i^2 - 2 \sum_{i=1}^N \sum_{k=1}^K \beta_k x_{k,i} y_i + \sum_{i=1}^N \sum_{k=1}^K \sum_{l=1}^K \beta_k \beta_l x_{k,i} x_{l,i} \end{aligned}$$

(Voir démonstration ci-après...)

Il faut maintenant minimiser cette forme quadratique par rapport à tous les paramètres :  $\beta_1, \beta_2, \dots, \beta_k, \dots, \beta_K$  !!!

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

35

35

### DÉMONSTRATION :

Décomposition de la forme quadratique  $Q(\beta)$  :

$$Q(\beta) = \varepsilon' \varepsilon = (y - X\beta)'(y - X\beta) = \underbrace{y'y}_{\text{Décomposition 1}} - \underbrace{2y'X\beta}_{\text{Décomposition 2}} + \underbrace{\beta'X'X\beta}_{\text{Décomposition 3}}$$

Décomposition 1 :  $y'y$

$$y'y = (y_1 \ y_2 \ \dots \ y_N) \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix} = \sum_{i=1}^N y_i^2 \quad \begin{array}{l} \rightarrow \text{Scalaire (1} \times \text{1)} \\ \rightarrow \text{Ne dépend pas de } \beta \end{array}$$

Décomposition 2 :  $-2y'X\beta$

$$y'X = (y_1 \ y_2 \ \dots \ y_N) \begin{bmatrix} 1 & x_{2,1} & \dots & x_{K,1} \\ 1 & x_{2,2} & \dots & x_{K,2} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{2,N} & \dots & x_{K,N} \end{bmatrix} = \left( \sum_i y_i \quad \sum_i y_i x_{2,i} \quad \dots \quad \sum_i y_i x_{K,i} \right)$$

$\rightarrow$  Vecteur-ligne (1  $\times$  K)

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

36

36

$$\mathbf{y}'\mathbf{X} = \left( \sum_i y_i \quad \sum_i y_i x_{2,i} \quad \cdots \quad \sum_i y_i x_{K,i} \right) = (\mathbf{y}'\mathbf{x}_1 \quad \mathbf{y}'\mathbf{x}_2 \quad \cdots \quad \mathbf{y}'\mathbf{x}_K)$$

avec :  $\mathbf{y}'\mathbf{x}_k = \sum_{i=1}^N y_i x_{k,i}$

2<sup>ème</sup> étape :  $-2\mathbf{y}'\mathbf{X}\boldsymbol{\beta} = -2(\mathbf{y}'\mathbf{x}_1 \quad \mathbf{y}'\mathbf{x}_2 \quad \cdots \quad \mathbf{y}'\mathbf{x}_K) \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{pmatrix} = -2 \sum_{k=1}^K \beta_k \mathbf{y}'\mathbf{x}_k$

$$= -2 \sum_{k=1}^K \beta_k \sum_{i=1}^N y_i x_{k,i}$$

$$= -2 \sum_{i=1}^N \sum_{k=1}^K \beta_k y_i x_{k,i}$$

→ *Scalaire (1 × 1)*  
→ *Combinaison linéaire en  $\boldsymbol{\beta}$*

37

### Décomposition 3 : $\boldsymbol{\beta}'\mathbf{X}'\mathbf{X}\boldsymbol{\beta}$

$$\mathbf{X}'\mathbf{X} = \begin{bmatrix} 1 & 1 & \cdots & 1 \\ x_{2,1} & x_{2,2} & \cdots & x_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ x_{K,1} & x_{K,2} & \cdots & x_{K,N} \end{bmatrix} \begin{bmatrix} 1 & x_{2,1} & \cdots & x_{K,1} \\ 1 & x_{2,2} & \cdots & x_{K,2} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{2,N} & \cdots & x_{K,N} \end{bmatrix} = \begin{bmatrix} N & \sum_i x_{2,i} & \cdots & \sum_i x_{K,i} \\ \sum_i x_{2,i} & \sum_i x_{2,i}^2 & \cdots & \sum_i x_{2,i}x_{K,i} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_i x_{K,i} & \sum_i x_{K,i}x_{2,i} & \cdots & \sum_i x_{K,i}^2 \end{bmatrix}$$

→ *Matrice symétrique (K × K)*

$$\frac{1}{N}\mathbf{X}'\mathbf{X} = \frac{\mathbf{X}'\mathbf{X}}{N} = \begin{bmatrix} 1 & \overline{x_2} & \cdots & \overline{x_K} \\ \overline{x_2} & \frac{1}{N}\sum_i x_{2,i}^2 & \cdots & \frac{1}{N}\sum_i x_{2,i}x_{K,i} \\ \vdots & \vdots & \ddots & \vdots \\ \overline{x_K} & \frac{1}{N}\sum_i x_{K,i}x_{2,i} & \cdots & \frac{1}{N}\sum_i x_{K,i}^2 \end{bmatrix}$$

avec :  $\overline{x_k} = \frac{1}{N}\sum_{i=1}^N x_{k,i}$

→ *Matrice des moments des variables X*

Si on note  $\mathbf{x}_i$  le vecteur (K × 1) des K variables explicatives pour un individu, la matrice des moments se réécrit comme :

$$\frac{1}{N}\mathbf{X}'\mathbf{X} = \frac{1}{N}\sum_{i=1}^N \mathbf{x}_i\mathbf{x}_i'$$

38

Maintenant on calcule la forme quadratique :

$$\begin{aligned}
 \beta' X' X \beta &= (\beta_1 \quad \beta_2 \quad \dots \quad \beta_K) \begin{bmatrix} N & \sum_i x_{2,i} & \dots & \sum_i x_{K,i} \\ \sum_i x_{2,i} & \sum_i x_{2,i}^2 & \dots & \sum_i x_{2,i} x_{K,i} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_i x_{K,i} & \sum_i x_{K,i} x_{2,i} & \dots & \sum_i x_{K,i}^2 \end{bmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_K \end{pmatrix} \\
 \beta' X' X \beta &= (\beta_1 \quad \beta_2 \quad \dots \quad \beta_K) \begin{pmatrix} \beta_1 N + \beta_2 \sum_i x_{2,i} + \dots + \beta_K \sum_i x_{K,i} \\ \beta_1 \sum_i x_{2,i} + \beta_2 \sum_i x_{2,i}^2 + \dots + \beta_K \sum_i x_{K,i} x_{2,i} \\ \vdots \\ \beta_1 \sum_i x_{K,i} + \beta_2 \sum_i x_{K,i} x_{2,i} + \dots + \beta_K \sum_i x_{K,i}^2 \end{pmatrix} \\
 &= (\beta_1 \quad \beta_2 \quad \dots \quad \beta_K) \begin{pmatrix} \sum_{k=1}^K \beta_k \sum_i x_{k,i} \\ \sum_{k=1}^K \beta_k \sum_i x_{2,i} x_{k,i} \\ \vdots \\ \sum_{k=1}^K \beta_k \sum_i x_{K,i} x_{k,i} \end{pmatrix} \\
 &= \beta_1 \sum_{k=1}^K \beta_k \sum_i x_{k,i} + \beta_2 \sum_{k=1}^K \beta_k \sum_i x_{2,i} x_{k,i} + \dots + \beta_K \sum_{k=1}^K \beta_k \sum_i x_{K,i} x_{k,i}
 \end{aligned}$$

Benoît MULKAY

Université de Montpellier

Econométrie Théorique (M1 MBFA)

Chapitre 1 (2023 – 2024)

39

39

$$\beta' X' X \beta = \sum_{l=1}^K \beta_l \sum_{k=1}^K \beta_k \sum_i x_{l,i} x_{k,i}$$

Ce qui donne finalement :

$$\beta' X' X \beta = \sum_{i=1}^N \sum_{k=1}^K \sum_{l=1}^K \beta_k \beta_l x_{k,i} x_{l,i}$$

En reprenant les 3 décompositions, on obtient alors :

$$\begin{aligned}
 Q(\beta) &= y' y - 2 y' X \beta + \beta' X' X \beta \\
 &= \sum_{i=1}^N y_i^2 - 2 \sum_{i=1}^N \sum_{k=1}^K \beta_k y_i x_{k,i} + \sum_{i=1}^N \sum_{k=1}^K \sum_{l=1}^K \beta_k \beta_l x_{k,i} x_{l,i}
 \end{aligned}$$

Benoît MULKAY

Université de Montpellier

Econométrie Théorique (M1 MBFA)

Chapitre 1 (2023 – 2024)

40

40

### La minimisation du critère des moindres carrés en utilisant la notation matricielle :

$$\min_{\beta} Q(\beta) = \varepsilon' \varepsilon = (y - X\beta)'(y - X\beta) = y'y - 2y'X\beta + \beta'X'X\beta$$

Formule de la dérivée d'une **combinaison linéaire** (2<sup>ème</sup> terme) :

$$\frac{\partial z' \theta}{\partial \theta} = \frac{\partial}{\partial \theta} \left( \sum_{k=1}^K \theta_k z_k \right) = z \quad \Rightarrow \quad \frac{\partial (-2y'X\beta)}{\partial \beta} = -2 \frac{\partial (y'X)\beta}{\partial \beta} = -2X'y$$

Formule de la dérivée d'une **forme quadratique symétrique** (3<sup>ème</sup> terme) :

$$\text{si } Z \text{ est symétrique : } \frac{\partial \beta' Z \beta}{\partial \beta} = 2Z\beta \quad \Rightarrow \quad \frac{\partial (\beta'(X'X)\beta)}{\partial \beta} = 2X'X\beta$$

Dérivée première du critère  $Q(\beta)$  par rapport au vecteur des paramètres  $\beta$  :

$$\frac{\partial Q(\beta)}{\partial \beta} = \frac{\partial}{\partial \beta} (y'y - 2y'X\beta + \beta'X'X\beta) = -2X'y + 2X'X\beta = -2(X'y - X'X\beta)$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

41

41

Un extrémum (ici le minimum) s'obtient en égalisant cette dérivée à 0 :

$$\rightarrow \text{Condition du premier ordre : } \frac{\partial Q(\beta)}{\partial \beta} = 0_K \quad \Rightarrow \quad -2(X'y - X'X\beta) = 0_K$$

Il faut résoudre ce système d'équations pour obtenir **l'estimateur des moindres carrés** :  $\hat{\beta}$

$$X'y - X'X\hat{\beta} = 0_K \quad \Rightarrow \quad X'X\hat{\beta} = X'y$$

Ce système est appelé **les équations normales** du problème des moindres carrés.

Ce système comprend  $K$  équations linéaires avec  $K$  inconnues : les composantes du vecteur  $\beta$  des paramètres.

Le système des  $K$  équations normales  $X'X\hat{\beta} = X'y$  peuvent se réécrire comme :

$$\begin{bmatrix} N & \sum_i x_{2,i} & \dots & \sum_i x_{K,i} \\ \sum_i x_{2,i} & \sum_i x_{2,i}^2 & \dots & \sum_i x_{2,i} x_{K,i} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_i x_{K,i} & \sum_i x_{K,i} x_{2,i} & \dots & \sum_i x_{K,i}^2 \end{bmatrix} \begin{pmatrix} \hat{\beta}_1 \\ \hat{\beta}_2 \\ \vdots \\ \hat{\beta}_K \end{pmatrix} = \begin{pmatrix} \sum_i y_i \\ \sum_i y_i x_{2,i} \\ \vdots \\ \sum_i y_i x_{K,i} \end{pmatrix}$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

42

42

Pour résoudre ce problème, il faut que la matrice ( $K \times K$ )  $X'X$  soit de rang plein :

$$\text{rang}(X'X) = K \Leftrightarrow \det(X'X) = |X'X| \neq 0 \quad (\text{ici } \det(X'X) > 0)$$

Pour cela, il suffit que l'hypothèse d'absence de multicollinéarité parfaite (H3) soit satisfaite : Si  $\text{rang}(X) = K \leq N \Rightarrow \text{rang}(X'X) = K$

Dans ce cas, l'inverse de  $X'X$  existe :  $(X'X)^{-1}$

et on obtient une **solution unique** pour le système des équations normales :  $X'X\hat{\beta} = X'y$  :

**l'estimateur des moindres carrés ordinaire (MCO) :**

$$\hat{\beta} = (X'X)^{-1}X'y$$

*En anglais : Ordinary Least-Squares (OLS) Estimator*

**Remarque** : dans le livre de Hayashi, l'estimateur MCO est noté en lettre romaine minuscule : **b** !

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

43

43

### Minimum de la somme des carrés des erreurs ?

L'estimateur des MCO est bien **le minimum** du critère parce que la dérivée seconde (*la matrice Hessienne*) :

$$\frac{\partial^2 Q(\beta)}{\partial \beta \partial \beta'} = \frac{\partial}{\partial \beta'} (-2X'y + 2X'X\beta) = 2X'X$$

est une **matrice symétrique définie-positive** si la condition de rang précédente est satisfaite (Hypothèse H3) :

$$\text{Si } \text{rang}(X) = K \leq N \Rightarrow \text{rang}(X'X) = K$$

La parabole du critère des moindres carrés est orientée vers le haut, elle a une forme en U.

Le minimum obtenu (dans le modèle linéaire) est **unique** (*identification*) !

On utilise souvent la notation :  $\frac{\partial^2 Q(\beta)}{\partial \beta \partial \beta'} = 2X'X > 0$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

44

44

## Dérivation alternative de l'estimateur des MCO : la méthode des moments

On part de l'hypothèse 2 de stricte exogénéité

$$E(\mathbf{x}_i \varepsilon_i) = \mathbf{0}_K \quad \text{pour tout } i = 1, 2, \dots, N.$$

On remplace  $\varepsilon_i$  par le vrai modèle:  $y_i - \mathbf{x}_i' \boldsymbol{\beta}$

$$E(\mathbf{x}_i (y_i - \mathbf{x}_i' \boldsymbol{\beta})) = \mathbf{0}_K$$

Le principe d'analogie de la méthode des moments est de remplacer l'espérance sur la population (des individus) par la moyenne sur un échantillon observé :

$$\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i (y_i - \mathbf{x}_i' \mathbf{b}) = \mathbf{0}_K$$

On obtient alors un système de  $K$  équations à  $K$  inconnues (les  $K$  éléments du vecteur  $\boldsymbol{\beta}$ ) pour calculer l'estimateur de la méthode des moments  $\mathbf{b}$ .

$$\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i (y_i - \mathbf{x}_i' \mathbf{b}) = \mathbf{0}_K$$

On peut réécrire ce système d'équation sous la forme :

$$\frac{1}{N} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i' \mathbf{b} = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i y_i$$

qui correspond aux équations normales de l'estimateur des MCO :

$$\sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i' \mathbf{b} = \sum_{i=1}^N \mathbf{x}_i y_i \quad \Leftrightarrow \quad \mathbf{X}' \mathbf{X} \mathbf{b} = \mathbf{X}' \mathbf{y}$$

En effet :  $\mathbf{X}' \mathbf{X} = \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i'$  et  $\mathbf{X}' \mathbf{y} = \sum_{i=1}^N \mathbf{x}_i y_i$ .

L'estimateur de la méthode des moments :

$$\mathbf{b} = (\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y} = \hat{\boldsymbol{\beta}}$$

est identique à l'estimateur des MCO.

### c) Remarques sur l'estimation par MCO

#### DÉFINITION :

La **valeur calculée** (*fitted*) de  $\mathbf{y}$  :  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$

Le **résidu** (*residuals*) de l'estimation est l'écart entre la valeur observée et la valeur calculée de  $\mathbf{y}$  :  $\mathbf{e} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}$

avec  $\mathbf{e}$  le vecteur ( $N \times 1$ ) des **résidus**, composé d'éléments  $e_i = y_i - \mathbf{x}'_i\hat{\boldsymbol{\beta}}$ .

Il ne faut pas confondre les **résidus**  $\mathbf{e}$  (calculables avec les observations) avec les **erreurs**  $\boldsymbol{\varepsilon}$  (inconnues par définition).

#### PROPRIÉTÉ DES MCO :

La minimisation du critère des moindres carrés implique **toujours** que le **vecteur des résidus soit orthogonal à la matrice des variables explicatives**.

$$\mathbf{X}'\mathbf{e} = \mathbf{0}_K$$

DÉMONSTRATION : A partir des  **$K$  équations normales**, on a :

$$\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{X}'\mathbf{y} = \mathbf{0}_K$$

$$\mathbf{X}'(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) = \mathbf{0}_K$$

$$\mathbf{X}'\mathbf{e} = \mathbf{0}_K$$

Attention à ne pas confondre, cette propriété de l'estimation par MCO :  $\mathbf{X}'\mathbf{e} = \mathbf{0}_K$  qui est toujours vérifiée avec l'hypothèse H2 :  $E(\mathbf{X}'\boldsymbol{\varepsilon}) = \mathbf{0}_K$

Pourquoi ?

#### CONSÉQUENCE 1 :

Il y a une **corrélation nulle** entre les variables explicatives et les résidus :

$$\mathbf{x}'_k\mathbf{e} = \sum_i x_{k,i}e_i = 0 \quad \text{pour tout } k = 1, 2, \dots, K.$$



### CONSÉQUENCE 2 :

**La somme des résidus est toujours nulle s'il y a une constante** dans la matrice de variables explicatives : (soit  $\mathbf{x}_1 = \mathbf{i}$  un vecteur de 1)

$$\mathbf{x}'_1 \mathbf{e} = \mathbf{i}' \mathbf{e} = \sum_i e_i = 0$$

Donc il est **inutile** de vérifier cette propriété quand il y a une constante dans la régression !!!

### CONSÉQUENCE 3 :

**Le plan de régression passe par le point moyen du nuage de points s'il y a une constante dans la régression :**

En prenant la première ligne des équations normales  $\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} = \mathbf{X}'\mathbf{y}$ , on a :

$$\begin{aligned} \left( N \quad \sum_i x_{2,i} \quad \cdots \quad \sum_i x_{K,i} \right) \hat{\boldsymbol{\beta}} &= \sum_i y_i \\ \left( 1 \quad \frac{1}{N} \sum_i x_{2,i} \quad \cdots \quad \frac{1}{N} \sum_i x_{K,i} \right) \hat{\boldsymbol{\beta}} &= \frac{1}{N} \sum_i y_i \\ (1 \quad \bar{x}_2 \quad \cdots \quad \bar{x}_K) \hat{\boldsymbol{\beta}} &= \bar{y} \quad \Rightarrow \quad \bar{y} = \bar{\mathbf{x}}' \hat{\boldsymbol{\beta}} \end{aligned}$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

49

49

**Dans la pratique**, la majorité des logiciels économétriques n'utilisent pas ces formules qui nécessitent d'inverser la matrice  $\mathbf{X}'\mathbf{X}$  ...

Ils utilisent plutôt une méthode basée sur la décomposition en base orthogonale et triangulaire (dite décomposition **QR**) beaucoup moins coûteuse à calculer en temps et en nombre d'opérations élémentaires (*voir Davidson-McKinnon [1993], section I.5*)

Soit la matrice  $N \times K$  des variables explicatives  $\mathbf{X}$ , on peut la réécrire comme :

$$\underset{N \times K}{\mathbf{X}} = \underset{N \times K}{\mathbf{Q}} \underset{K \times K}{\mathbf{R}} \quad \text{avec} \quad \mathbf{Q}'\mathbf{Q} = \mathbf{I}_K$$

Les colonnes de  $\mathbf{Q}$  sont donc orthonormales.

$\mathbf{R}$  est une matrice carrée de dimension  $K \times K$  triangulaire supérieure telle que :

$$\mathbf{X}'\mathbf{X} = \mathbf{R}'\mathbf{R}$$

En effet :  $\mathbf{X}'\mathbf{X} = (\mathbf{QR})'(\mathbf{QR}) = \mathbf{R}'\mathbf{Q}'\mathbf{Q}\mathbf{R} = \mathbf{R}'\mathbf{R}$  parce que  $\mathbf{Q}'\mathbf{Q} = \mathbf{I}_K$

On peut montrer que l'estimateur des moindres carrés se calcule alors comme :

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{R}^{-1}\mathbf{Q}'\mathbf{y}$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

50

50

En effet :

$$\hat{\beta} = (X'X)^{-1}X'y = (R'R)^{-1}(QR)'y = R^{-1}R'^{-1}R'Q'y = R^{-1}Q'y$$

Cette formulation alternative nécessite seulement l'inversion de la matrice triangulaire supérieure  $R$ , ce qui est facile récursivement ...

mais au prix du calcul relativement simple (récursivement) de  $Q$  et de  $R$  à partir de la matrice  $X$  !

Le nombre d'opérations élémentaires est réduit et les arrondis de calculs sont moindres avec cette méthode alternative !

En plus, cette méthode permet d'identifier la (les) variable(s) qui seraient parfaitement co-linéaires !

### I.3. L'estimateur de la variance de l'erreur

On peut également estimer un second paramètre : **la variance de l'erreur** :  $\sigma^2$

Comme dans le cas de la régression simple, on va utiliser la somme des carrés des résidus.

Le vecteur ( $N \times 1$ ) des résidus  $e$  est défini par :  $e = y - X\hat{\beta}$

En remplaçant l'estimateur MCO par sa formule :

$$\begin{aligned} e &= y - X(X'X)^{-1}X'y \\ &= (I_N - X(X'X)^{-1}X')y = My \quad \text{avec : } M = I_N - X(X'X)^{-1}X' \end{aligned}$$

La matrice  $M$  de dimension ( $N \times N$ ) est une matrice symétrique et idempotente :

$$M' = M \quad \text{et} \quad MM = M \quad (\text{à démontrer ?})$$

Cette matrice  $M$  est appelée « *residual – maker* » par W. Greene.

De plus, si on pré-multiplie  $X$  par cette matrice  $M$ , on obtient une matrice nulle :

$$MX = (I_N - X(X'X)^{-1}X')X = X - X(X'X)^{-1}X'X = X - X = 0_{N \times K}$$

On aura avec le résultat précédent :

$$e = My = M(X\beta + \varepsilon) = MX\beta + M\varepsilon = M\varepsilon$$

Cependant cette expression est uniquement théorique parce qu'on ne connaît pas  $\varepsilon$ . On ne peut pas trouver le vecteur des erreurs  $\varepsilon$  à partir du vecteur des résidus  $e$  en calculant :  $\varepsilon = M^{-1}e$ , parce que la matrice  $M$  n'est pas inversible ! (Pourquoi ?)

Il n'est pas d'un intérêt pratique de calculer (ni d'imprimer) cette matrice  $M$ , parce qu'elle peut être très grande :  $N \times N$  (même si elle est symétrique).

Imaginez un échantillon de 1 000 observations : elle serait de dimension  $1\,000 \times 1\,000$ , contenant 1 000 000 de nombres, soit 8 Mo !  
(en général : 1 nombre = 8 octets)

**DÉFINITION** : La somme des carrés des résidus (SCR) :

$$SCR = e'e = \sum_{i=1}^N e_i^2 \quad \text{avec} \quad e_i = y_i - x_i'\hat{\beta}$$

Cette somme des carrés des résidus peut s'exprimer théoriquement en fonction du terme d'erreurs :

$$SCR = e'e = (M\varepsilon)'(M\varepsilon) = \varepsilon'M'M\varepsilon = \varepsilon'M\varepsilon$$

On peut aussi réécrire la SCR de la manière suivante :

$$SCR = e'e = (y - X\hat{\beta})'(y - X\hat{\beta}) = y'y - 2y'X\hat{\beta} + \hat{\beta}'X'X\hat{\beta}$$

Ce qui correspond à la **valeur minimisée** du critère des MCO !

On remplace dans le dernier terme, le premier estimateur MCO :  $\hat{\beta}' = y'X(X'X)^{-1}$  pour obtenir :

$$SCR = y'y - 2y'X\hat{\beta} + y'X(X'X)^{-1}X'X\hat{\beta} = y'y - 2y'X\hat{\beta} + y'X\hat{\beta}$$

Et finalement :  $SCR = y'y - y'X\hat{\beta} = y'(y - X\hat{\beta}) = y'e$

On montrera plus loin (Section II.1.d) que :

$$E(SCR) = E(\mathbf{e}'\mathbf{e}) = E(\boldsymbol{\varepsilon}'\mathbf{M}\boldsymbol{\varepsilon}) = (N - K)\sigma^2$$

On suppose maintenant que le nombre d'observations est strictement supérieur au nombre de variables :  $N > K$ .

**DÉFINITION** : L'estimateur des moindres carrés de la **variance de l'erreur** sera donné par :

$$\widehat{\sigma}^2 = \frac{SCR}{N - K} = \frac{1}{N - K} \sum_{i=1}^N e_i^2$$

On a ainsi estimé la dispersion de la variable dépendante autour de la droite de régression.

Attention, ce n'est pas la moyenne des carrés des résidus :  $\frac{1}{N} \sum_{i=1}^N e_i^2$  !  
ni la variance des résidus :  $\frac{1}{N-1} \sum_{i=1}^N e_i^2$  !

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

55

55

**DÉFINITION** : Le facteur  $(N - K)$  est appelé le nombre de **degrés de liberté** (**ddl**) de la régression (ou tout simplement) les degrés de liberté...  
(*en anglais : degrees of freedom - df*)

$$ddl = (N - K)$$

Remarquez que cet **estimateur de la variance** de l'erreur est mesuré dans le **carré** des unités de  $\mathbf{e}$  (ou de  $\mathbf{y}$ ) !

**DÉFINITION** : **L'écart-type de la régression** est la racine carrée de cet estimateur de la variance.  
(*En anglais : standard-error of the regression ou encore : Root mean-squared error – RMSE*)

$$\hat{\sigma} = \sqrt{\widehat{\sigma}^2} = \sqrt{\frac{SCR}{N - K}}$$

L'écart-type de la régression est mesuré avec les unités de  $\mathbf{e}$  (ou de  $\mathbf{y}$ ) !

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

56

56

## I.4. La qualité de l'ajustement

### a) Le coefficient de détermination : $R^2$

F. HAYASHI, *Econometrics* [2000], Section I.2. (p. 20-21)

W. GREENE, *Econométrie* [2018], Section III.4.

On a vu que l'estimation par MCO du modèle de régression impliquait par définition :

$$\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e} = \hat{\mathbf{y}} + \mathbf{e}$$

Pré-multiplions cette expression par sa transposée :

$$\mathbf{y}'\mathbf{y} = (\mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e})'(\mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e})$$

$$\mathbf{y}'\mathbf{y} = \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} + \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{e} + \mathbf{e}'\mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e}'\mathbf{e}$$

Nous avons vu plus haut que les équations normales impliquaient :  $\mathbf{X}'\mathbf{e} = \mathbf{0}_K$ . En conséquence les deux termes centraux sont nuls ...

$$\mathbf{y}'\mathbf{y} = \hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e}'\mathbf{e}$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

57

57

Soustrayons maintenant le carré de la moyenne de  $\mathbf{y}$  multiplié par le nombre d'observations de chaque côté de l'égalité :

$$(\mathbf{y}'\mathbf{y} - N\bar{y}^2) = (\hat{\boldsymbol{\beta}}'\mathbf{X}'\mathbf{X}\hat{\boldsymbol{\beta}} - N\bar{y}^2) + \mathbf{e}'\mathbf{e}$$

Comme  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$ , il est évident que la moyenne de  $\hat{\mathbf{y}}$  est identique à la moyenne de  $\mathbf{y}$ .

$$\underbrace{(\mathbf{y}'\mathbf{y} - N\bar{y}^2)}_{SCT} = \underbrace{(\hat{\mathbf{y}}'\hat{\mathbf{y}} - N\bar{y}^2)}_{SCE} + \underbrace{\mathbf{e}'\mathbf{e}}_{SCR}$$

avec :

- $SCT = \mathbf{y}'\mathbf{y} - N\bar{y}^2$  : la somme des carrés totaux (*de la variable dépendante*)
- $SCE = \hat{\mathbf{y}}'\hat{\mathbf{y}} - N\bar{y}^2$  : la somme des carrés expliqués (*par la régression*)
- $SCR = \mathbf{e}'\mathbf{e}$  : la somme des carrés des résidus

La **somme des carrés totaux** ( $SCT$ ) est définie par :

$$SCT = \sum_{i=1}^N (y_i - \bar{y})^2 = \mathbf{y}'\mathbf{y} - N\bar{y}^2$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

58

58

La **somme des carrés expliqués** ( $SCE$ ) par la régression est définie par :

$$SCE = \sum_{i=1}^N (\hat{y}_i - \bar{y})^2 = \mathbf{\hat{y}}' \mathbf{\hat{y}} - N \bar{y}^2$$

La **somme des carrés des résidus** ( $SCR$ ) est définie par :

$$SCR = \sum_{i=1}^N e_i^2 = \mathbf{e}' \mathbf{e}$$

### En notation matricielle

La moyenne se calcule comme :  $\bar{y} = \frac{\mathbf{j}_N' \mathbf{y}}{N}$  avec  $\mathbf{j}_N = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$

Dès lors :  $N \bar{y}^2 = N \left( \frac{\mathbf{j}_N' \mathbf{y}}{N} \right)' \left( \frac{\mathbf{j}_N' \mathbf{y}}{N} \right) = \mathbf{y}' \frac{\mathbf{j}_N \mathbf{j}_N'}{N} \mathbf{y} = \mathbf{y}' \frac{\mathbf{j}_N}{N} \mathbf{y}$  avec  $\mathbf{j}_N = \mathbf{j}_N \mathbf{j}_N'$

La matrice  $\mathbf{J}_N$  est une matrice carrée de dimension  $N \times N$  qui ne contient que des valeurs 1.

En conséquence :  $SCT = \mathbf{y}' \mathbf{y} - N \bar{y}^2 = \mathbf{y}' \mathbf{y} - \mathbf{y}' \frac{\mathbf{j}_N}{N} \mathbf{y} = \mathbf{y}' \left( \mathbf{I}_N - \frac{\mathbf{j}_N}{N} \right) \mathbf{y} = \mathbf{y}' \mathbf{M}_0 \mathbf{y}$

La matrice  $\mathbf{M}_0$  effectue le centrage des observations  $\mathbf{y}$  par rapport à leur moyenne :

$$\mathbf{M}_0 \mathbf{y} = \left( \mathbf{I}_N - \frac{\mathbf{j}_N}{N} \right) \mathbf{y} = \mathbf{y} - \mathbf{j}_N \frac{\mathbf{j}_N' \mathbf{y}}{N} = \mathbf{y} - \mathbf{j}_N \bar{y}$$

La matrice  $\mathbf{M}_0$  est une matrice symétrique et idempotente.

Remarque : elle donne le vecteur du résidu d'une régression où la seule variable explicative est la constante !

Si on pré-multiplie la régression par la matrice  $M_0$ , on obtient :

$$y = X\hat{\beta} + e \rightarrow M_0 y = M_0 X\hat{\beta} + M_0 e = M_0 X\hat{\beta} + e$$

parce que les résidus ont une moyenne nulle s'il y a une constante parmi les variables explicatives  $X$  :

$$M_0 e = \left(I_N - \frac{J_N}{N}\right) e = e - j_N \frac{j_N' e}{N} = e - j_N \times 0 = e$$

On effectue alors le produit scalaire de ce vecteur :

$$(M_0 y)'(M_0 y) = (M_0 X\hat{\beta} + e)'(M_0 X\hat{\beta} + e)$$

$$y' M_0 y = \hat{\beta}' X' M_0 X \hat{\beta} + \hat{\beta}' X' M_0 e + e' M_0 X \hat{\beta} + e' e$$

$$y' M_0 y = \hat{\beta}' X' M_0 X \hat{\beta} + e' e \quad \text{parce que } X' M_0 e = X' e = 0$$

$$SCT = SCE + SCR$$

du fait des définitions de  $SCT$ ,  $SCE$  et  $SCR$ .

Comme la somme des carrés totaux (la variabilité) de la variable  $y$  se décompose dans la somme de ce qui est expliqué par la régression ( $SCE$ ) et de ce qui est inexpliqué ( $SCR$ ) :  $SCT = SCE + SCR$ , un bon ajustement voudrait que la part de la  $SCE$  soit importante, et la part de la  $SCR$  soit faible.

Si on divise partout par la  $SCT$ , on aura :  $\frac{SCE}{SCT} + \frac{SCR}{SCT} = 1$

**DEFINITION** : On mesure la qualité de l'ajustement de la régression avec le coefficient de détermination :  $R^2$  qui est défini comme :

$$R^2 = \frac{SCE}{SCT} = 1 - \frac{SCR}{SCT}$$

$$R^2 = \frac{\sum_{i=1}^N (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^N (y_i - \bar{y})^2} = 1 - \frac{\sum_{i=1}^N e_i^2}{\sum_{i=1}^N (y_i - \bar{y})^2}$$

On peut aussi réécrire le  $R^2$  sous la forme :

$$R^2 = 1 - \frac{SCR}{SCT} = 1 - \frac{\mathbf{y}'\mathbf{M}\mathbf{y}}{\mathbf{y}'\mathbf{M}_0\mathbf{y}} \quad \left\{ \begin{array}{l} \mathbf{M} = \mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}' \\ \mathbf{M}_0 = \mathbf{I}_N - \frac{\mathbf{J}_N}{N} = \mathbf{I}_N - \mathbf{j}_N(\mathbf{j}_N'\mathbf{j}_N)^{-1}\mathbf{j}_N' \end{array} \right.$$

**PROPRIÉTÉ 1 :** Pour autant qu'il y ait une constante dans le modèle, le coefficient de détermination  $R^2$  est compris entre 0 et 1.

**PROPRIÉTÉ 2 :** Le coefficient de détermination  $R^2$  est aussi le carré du coefficient de corrélation entre la valeur observée et la valeur calculée de la variable dépendante :

$$R^2 = \text{Corr}^2(y, \hat{y}) = \frac{\text{Cov}^2(y, \hat{y})}{V(y) \times V(\hat{y})}$$

(Voir démonstration ci-après ...)

### Démonstration :

D'après la définition des variances :

$$V(y) = \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})^2 = \frac{SCT}{N-1}$$

$$V(\hat{y}) = \frac{1}{N-1} \sum_{i=1}^N (\hat{y}_i - \bar{y})^2 = \frac{SCE}{N-1}$$

La covariance entre la variable dépendante  $\mathbf{y}$  et la variable calculée  $\hat{\mathbf{y}} = \mathbf{X}\hat{\boldsymbol{\beta}}$  est par définition :

$$\begin{aligned} \text{Cov}(y, \hat{y}) &= \frac{1}{N-1} \sum_{i=1}^N (y_i - \bar{y})(\hat{y}_i - \bar{y}) = \frac{1}{N-1} \left[ \sum_{i=1}^N y_i \hat{y}_i - N\bar{y}^2 \right] \\ &= \frac{1}{N-1} (\mathbf{y}'\hat{\mathbf{y}} - N\bar{y}^2) = \frac{1}{N-1} \left( \mathbf{y}'\hat{\mathbf{y}} - \mathbf{y}'\frac{\mathbf{J}_N}{N}\hat{\mathbf{y}} \right) = \frac{1}{N-1} \mathbf{y}'\mathbf{M}_0\hat{\mathbf{y}} \end{aligned}$$



On peut aussi réécrire la covariance comme :

$$\begin{aligned}
 Cov(y, \hat{y}) &= \frac{1}{N-1} \mathbf{y}' \mathbf{M}_0 \hat{\mathbf{y}} \\
 &= \frac{1}{N-1} (\mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e})' \mathbf{M}_0 \mathbf{X}\hat{\boldsymbol{\beta}} && \text{parce que : } \mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} + \mathbf{e} \\
 &= \frac{1}{N-1} (\hat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{M}_0 \mathbf{X} \hat{\boldsymbol{\beta}} + \mathbf{e}' \mathbf{M}_0 \mathbf{X} \hat{\boldsymbol{\beta}}) \\
 &= \frac{1}{N-1} \hat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{M}_0 \mathbf{X} \hat{\boldsymbol{\beta}} && \text{parce que : } \mathbf{e}' \mathbf{M}_0 \mathbf{X} = \mathbf{e}' \mathbf{X} = \mathbf{0} \\
 &= \frac{SCE}{N-1}
 \end{aligned}$$

En conséquence, le carré de la corrélation entre variable dépendante et variable prédite devient le coefficient de détermination:

$$Corr^2(y, \hat{y}) = \frac{Cov^2(y, \hat{y})}{V(y) \times V(\hat{y})} = \frac{\left(\frac{SCE}{N-1}\right)^2}{\left(\frac{SCT}{N-1}\right) \left(\frac{SCE}{N-1}\right)} = \frac{SCE}{SCT} = R^2$$

**CQFD.**

### **b) Le $R^2$ non-centré**

**Problème du  $R^2$**  : Si il n'y a pas de constante dans le modèle,  
le  $R^2$  peut devenir négatif !!!  
En effet on aura :

Si  $SCR > SCT$  ou si  $\sum_{i=1}^N e_i^2 > \sum_{i=1}^N (y_i - \bar{y})^2$  parce que  $\sum_{i=1}^N e_i \neq 0$

alors  $R^2 = 1 - \frac{SCR}{SCT} < 0$

Pour éviter cette valeur négative dans ce cas, on propose d'utiliser un **coefficient de détermination non centré** :

$$R_{nc}^2 = 1 - \frac{\sum_{i=1}^N e_i^2}{\sum_{i=1}^N y_i^2} = 1 - \frac{\mathbf{e}' \mathbf{e}}{\mathbf{y}' \mathbf{y}}$$

En effet, on part de la décomposition de la variable dépendante :  $y = \hat{y} + e$ .

Calculons son produit scalaire :

$$\begin{aligned} y'y &= (\hat{y} + e)'(\hat{y} + e) = \hat{y}'\hat{y} + 2\hat{y}'e + e'e \\ &= \hat{y}'\hat{y} + 2\hat{\beta}'X'e + e'e && \text{parce que : } \hat{y} = X\hat{\beta} \\ &= \hat{y}'\hat{y} + e'e && \text{parce que : } X'e = 0 \end{aligned}$$

Le coefficient de détermination non centré est alors :

$$R_{nc}^2 = \frac{\hat{y}'\hat{y}}{y'y} = 1 - \frac{e'e}{y'y} \quad \text{avec } 0 \leq R_{nc}^2 \leq 1$$

Il est toujours non-négatif et supérieur au coefficient de détermination centré :  
 $R_{nc}^2 \geq R^2$

*Pourquoi ? Démontrez ...*

Si il n'y a pas de constante dans le modèle, certains logiciels (*Stata par exemple*) calculent automatiquement le coefficient de détermination non centré

$$R_{nc}^2 = 1 - \frac{e'e}{y'y}$$

plutôt que le coefficient de détermination centré :

$$R^2 = 1 - \frac{\sum_{i=1}^N e_i^2}{\sum_{i=1}^N (y_i - \bar{y})^2} = 1 - \frac{e'e}{y'y - N\bar{y}^2}$$

afin d'éviter d'obtenir une valeur négative pour ce dernier.

Mais attention au cas où il n'y a pas de constante, mais où une somme de variables explicatives est constante ...!!!

### c) Le $R^2$ ajusté pour les degrés de liberté

**Problème du  $R^2$**  : Il augmente toujours lorsque l'on ajoute une variable explicative supplémentaire...  
parce que la somme des carrés des résidus diminue...

**Correction du  $R^2$**  : Ajustement pour le nombre de variables explicatives ( $K$ )

$$\bar{R}^2 = 1 - \frac{SCR/(N - K)}{SCT/(N - 1)}$$

$$\bar{R}^2 = 1 - \frac{\mathbf{e}'\mathbf{e}/(N - K)}{\mathbf{y}'\mathbf{M}_0\mathbf{y}/(N - 1)} = 1 - \frac{\hat{\sigma}^2}{V(\mathbf{y})}$$

On aura un **coefficient de détermination ajusté** (corrigé) pour les degrés de liberté.

### **Relation entre le $R^2$ et le $R^2$ ajusté :**

A partir des définitions des coefficients de détermination, il est facile de montrer que :

$$\bar{R}^2 = 1 - \frac{N - 1}{N - K} (1 - R^2)$$

De même le  $\bar{R}^2$  ajusté sera toujours inférieur ou égal au  $R^2$  classique :

$$\bar{R}^2 \leq R^2$$

L'égalité s'obtient lorsque l'ajustement est parfait :  $\bar{R}^2 = R^2 = 1$  !

Attention le  $\bar{R}^2$  ajusté peut être négatif !

• supposons que le  $R^2$  soit nul  $\rightarrow \bar{R}^2 = \frac{1 - K}{N - K} \leq 0$

• Si le  $R^2 \leq \frac{K-1}{N-1}$ , alors  $\bar{R}^2 \leq 0$

Le  $\bar{R}^2$  ajusté peut diminuer lorsque l'on ajoute une variable explicative supplémentaire...

... si celle-ci n'apporte pas d'information suffisante pour un meilleur ajustement de la variable dépendante  $y$ .

Cette variable explicative supplémentaire serait alors non pertinente !

(voir la démonstration dans la Section I.6)

#### **d) Remarque sur l'utilisation du $R^2$**

Le  $R^2$  est souvent utilisé pour distinguer le pouvoir explicatif de deux régressions.

On peut choisir le modèle qui a le coefficient de détermination le plus élevé.

Cette stratégie de modélisation soulève cependant plusieurs remarques...

- 1) Comme le  $R^2$  classique **augmente toujours** lorsque l'on ajoute une variable explicative à la régression, *même si elle est non pertinente (non significative)*, il est préférable d'utiliser le  $\bar{R}^2$  ajusté pour les degrés de liberté.
- 2) Pour choisir entre des régresseurs, on doit comparer des régressions comparables : c'est-à-dire faites sur le **même échantillon**.  
Donc il faut **vérifier que la taille et la composition de l'échantillon** ne changent pas entre les 2 régressions.

- 3) Le **choix d'un modèle ne peut pas se faire sur la seule base du coefficient de détermination  $R^2$** . D'autres critères doivent être pris en compte :
- Significativité des paramètres estimés
  - Interprétation économique des paramètres estimés
  - Parcimonie
  - Autres tests statistiques (*voir suite du cours*)
- 4) Il faut aussi que la **variable dépendante soit la même** dans les 2 modèles. En effet elle apparaît dans la **SCT** au dénominateur du  $R^2$ .

## I.5. Interprétation géométrique des MCO

On peut considérer chaque variable comme un point (ou un vecteur) dans un espace vectoriel à  $N$  dimensions :  $\mathbb{R}^N$ .

Donc la variable dépendante  $y$  représente un vecteur dans cet espace vectoriel, et les  $K$  variables explicatives  $X$  représentent  $K$  vecteurs dans ce même espace.

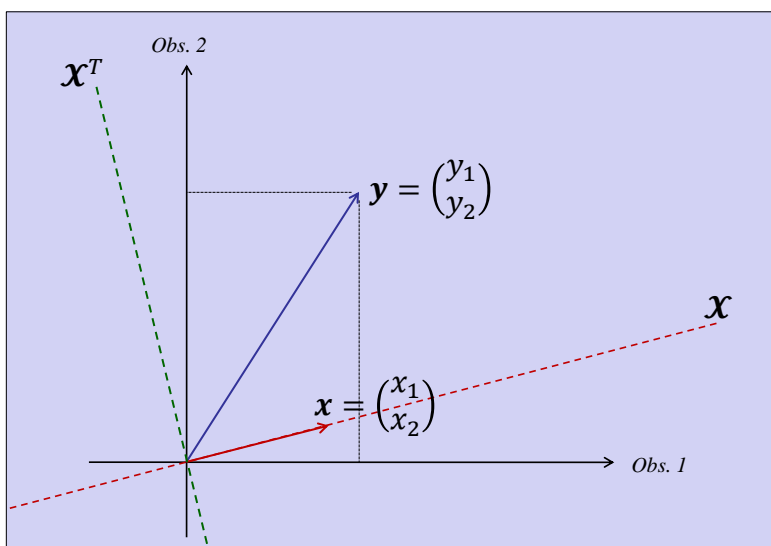
Ces  $K$  vecteurs définissent un sous-espace vectoriel (un hyperplan)  $\mathcal{X}$  de dimension  $K$  avec :  $\mathcal{X} \subset \mathbb{R}^N$  avec  $\dim(\mathcal{X}) = K$ .

De même dans  $\mathbb{R}^N$ , on peut définir un espace vectoriel (un hyperplan) orthogonal à  $\mathcal{X}$  de dimension  $N - K$ , appelé  $\mathcal{X}^T$  :

$$\mathcal{X} \perp \mathcal{X}^T \text{ et } \mathcal{X}^T \subset \mathbb{R}^N \text{ avec } \dim(\mathcal{X}^T) = N - K$$

$$\text{tel que } \mathcal{X}^T \cup \mathcal{X} = \mathbb{R}^N$$

**Exemple** :  $N = 2$  observations et  $K = 1$ , une seule variable explicative



Benoît MULKAY  
Université de Montpellier

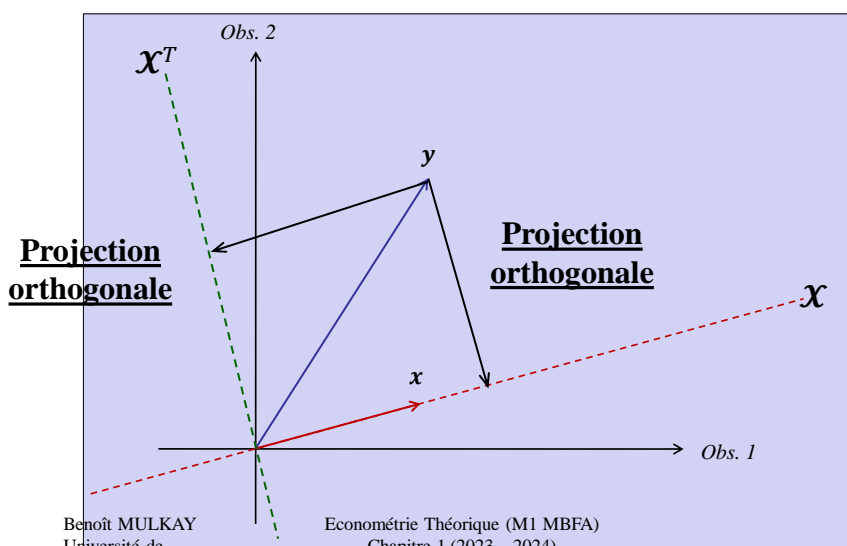
Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

75

75

La régression consiste à effectuer deux **projections orthogonales** du vecteur  $y$

- dans le sous-espace vectoriel  $\mathcal{X}$  engendré par les variables explicatives  $X$ ,
- dans le sous-espace vectoriel orthogonal  $\mathcal{X}^T \perp \mathcal{X}$ .



Benoît MULKAY  
Université de  
Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

76

76

Mathématiquement, en géométrie analytique, une **projection orthogonale** est réalisée par des opérations matricielles sur les vecteurs.

Une matrice  $\mathbf{P}$  de dimension  $N \times N$  est une **matrice de projection** si elle est symétrique ( $\mathbf{P}' = \mathbf{P}$ ) et idempotente ( $\mathbf{P}\mathbf{P} = \mathbf{P}$ ).

La dimension du sous-espace de projection est donnée par :  $\dim(\mathcal{X}) = \text{rang}(\mathbf{P})$

Pour effectuer une projection dans le sous-espace  $\mathcal{X}$  engendré par les variables explicatives  $\mathbf{X}$ , on utilisera la **matrice de projection** :

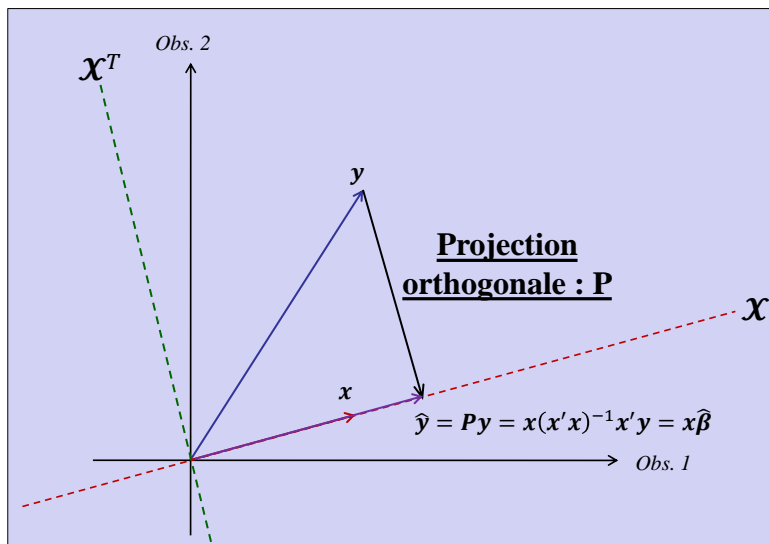
$$\mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

Vérifiez que cette matrice est de dimension  $N \times N$ , symétrique et idempotente et que  $\dim(\mathcal{X}) = \text{rang}(\mathbf{P}) = K$ .

Le vecteur projeté dans le sous-espace  $\mathcal{X}$  de la variable dépendante  $\mathbf{y}$  sera alors :

$$\mathbf{P}\mathbf{y} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} = \hat{\mathbf{y}}$$

**Exemple** :  $N = 2$  observations et  $K = 1$ , une seule variable explicative



La **matrice de projection** dans le sous-espace orthogonal  $\mathcal{X}^T$  à celui engendré par les variables explicatives  $\mathbf{X}$ , est la matrice  $\mathbf{M}$  précédente :

$$\mathbf{M} = \mathbf{I}_N - \mathbf{P} = (\mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}')$$

Cette matrice  $N \times N$  est symétrique et idempotente et orthogonale à  $\mathbf{P}$  :

$$\mathbf{PM} = \mathbf{P}(\mathbf{I}_N - \mathbf{P})\mathbf{I}_N = \mathbf{P} - \mathbf{PP} = \mathbf{P} - \mathbf{P} = \mathbf{0}$$

La dimension du sous-espace  $\mathcal{X}^T$  est :  $\dim(\mathcal{X}^T) = \text{rang}(\mathbf{M}) = N - K$

Si on applique cette projection  $\mathbf{M}$  à la matrice des variables explicatives  $\mathbf{X}$ , on obtient le vecteur nul :

$$\mathbf{MX} = (\mathbf{I}_N - \mathbf{P})\mathbf{X} = \mathbf{X} - \mathbf{PX} = \mathbf{X} - \mathbf{X} = \mathbf{0}$$

En effet le sous-espace  $\mathcal{X}^T$  est par construction orthogonal au sous-espace  $\mathcal{X}$  qui est supporté par les vecteurs de  $\mathbf{X}$ .

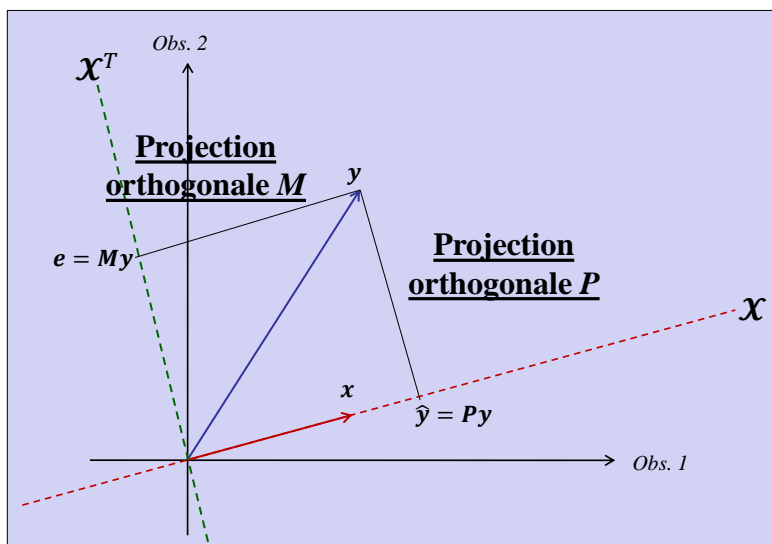
En revanche si on applique cette projection à la variable dépendante  $\mathbf{y}$ , on obtient le vecteur du résidu  $\mathbf{e}$  :

$$\mathbf{My} = (\mathbf{I}_N - \mathbf{P})\mathbf{y} = \mathbf{y} - \mathbf{Py} = \mathbf{y} - \hat{\mathbf{y}} = \mathbf{e}$$

Donc le résidu  $\mathbf{e}$  est le vecteur  $\mathbf{y}$  projeté orthogonalement dans le sous-espace vectoriel  $\mathcal{X}^T$  à celui engendré par les variables explicatives  $\mathbf{X}$ .



**Exemple** :  $N = 2$  observations et  $K = 1$ , une seule variable explicative



Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

81

81

On a ainsi décomposé la variable explicative  $y$  en deux composantes vectorielles orthogonales :

$$y = \hat{y} + e \rightarrow \begin{cases} \hat{y} \in \mathcal{X} \\ e \in \mathcal{X}^\perp \end{cases} \text{ avec } \hat{y} \perp e$$

On a utilisé toute l'information contenue dans les variables explicatives  $X$  pour obtenir un  $y$  calculé.

En conséquence, il n'y a plus d'information disponible dans  $X$ , contenue dans le résidu  $e$ , afin d'améliorer la prévision de  $y$  **sans hypothèse supplémentaire**.

L'estimateur des MCO utilise ainsi l'ensemble de l'information disponible dans les vecteurs des variables explicatives...

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

82

82

### **b) Interprétation géométrique du $R^2$**

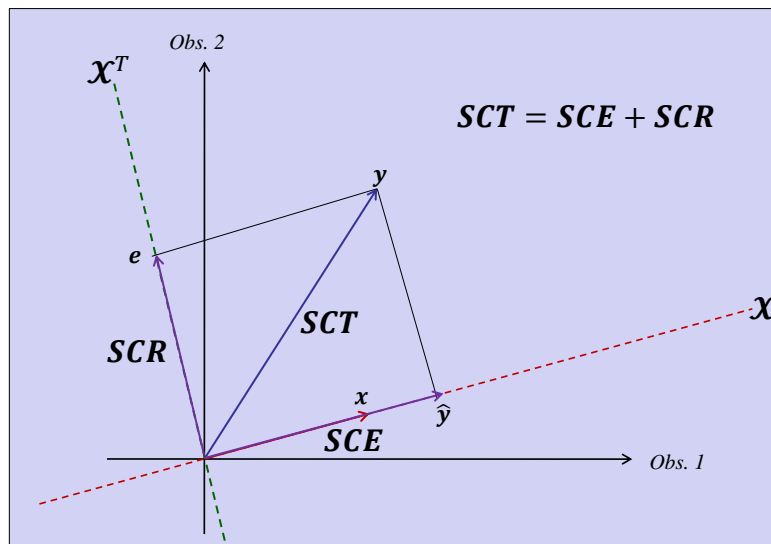
Retour sur l'interprétation géométrique de la régression par MCO  
(cfr. Section I.4).

La décomposition  $SCT = SCE + SCE$  provient aussi du théorème de Pythagore :

*Dans un triangle rectangle, le carré de la longueur de l'hypoténuse est égal à la somme des carrés des longueurs des deux autres côtés.*

avec le fait qu'on a une décomposition orthogonale du vecteur  $y$  :

$$y = \hat{y} + e \quad \text{avec} \quad \hat{y} \perp e$$



La régression consiste à faire une projection orthogonale du vecteur  $\mathbf{y}$  dans le plan des régresseurs avec la matrice de projection  $\mathbf{P}$ .

On obtient le vecteur des variables calculées :

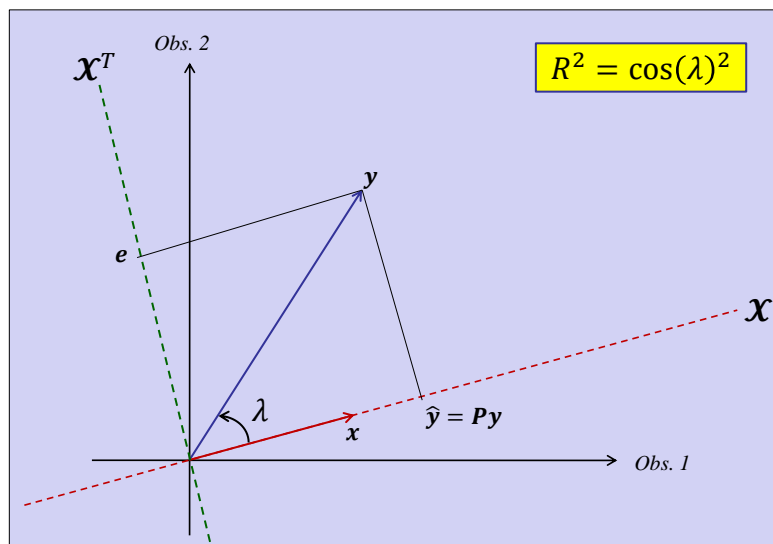
$$\mathbf{Py} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} = \mathbf{X}\hat{\boldsymbol{\beta}} = \hat{\mathbf{y}}$$

Le coefficient de détermination  $R^2$  est relié à l'angle  $\lambda$  entre le vecteur des variables dépendantes  $\mathbf{y}$  et le vecteur des variables calculées  $\hat{\mathbf{y}}$  :

$$R^2 = \cos(\lambda)^2$$

Plus le vecteur de la variable dépendante sera proche du plan de régression (engendré par les variables explicatives), plus le coefficient de détermination sera important.

**Exemple** :  $N = 2$  observations et  $K = 1$ , une seule variable explicative



## L6. Le théorème de Frisch – Waugh – Lovell

Estimation séparée d'un sous-ensemble de paramètres :

- les paramètres d'intérêt
- les paramètres de « nuisance »

Comment enlever une tendance déterministe des données ?

Comment éviter de faire des inversions matricielles trop complexes ?

Voir : Ragnar FRISCH et Frederick WAUGH (1933) : « Partial Time Regressions as Compared to Individual Trends », *Econometrica*, I(4), pp. 387-401.

Michael C. LOVELL (1963) : « Seasonal Adjustment of Economic Time Series and Multiple Regression Analysis », *Journal of American Statistical Association*, Vol. 58, N°304, pp. 996-1010.

L'intérêt pratique de ce théorème est assez limité de nos jours !

Mais il est très utile pour beaucoup de démonstration des propriétés de l'estimateur des MCO.

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

87

87

### a) La régression partielle

Modèle de régression multiple avec  $K$  variables explicatives :  $y = X\beta + \varepsilon$

Les  $K$  variables explicatives sont **partitionnées** en 2 groupes de variables :

$$y = X\beta + \varepsilon = (X_1 \quad X_2) \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \varepsilon$$

avec les matrices :  $X_1 : N \times K_1$  et  $X_2 : N \times K_2$  avec  $K_1 + K_2 = K$   
et les vecteurs :  $\beta_1 : K_1 \times 1$  et  $\beta_2 : K_2 \times 1$ .

On s'intéresse plus particulièrement aux paramètres  $\beta_1$  (les paramètres d'intérêt).

Cela peut être le paramètre d'une seule variable explicative si :  $K_1 = 1$  et  $K_2 = K - 1$  ( $X_1$  devient alors un vecteur).

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

88

88

Voir : W. GREENE, Econométrie [2018], Théorème 3.3

### Théorème de Frisch-Waugh-Lovell

Dans la régression linéaire par MCO de  $y$  sur deux groupes de variables explicatives  $X_1$  et  $X_2$ ,

$$y = X_1\beta_1 + X_2\beta_2 + \varepsilon$$

le sous-vecteur  $\widehat{\beta}_1$  des paramètres est obtenu quand les résidus ( $u$ ) d'une régression de  $y$  sur  $X_2$  sont régressés sur l'ensemble des résidus ( $V$ ) d'une régression de chaque colonne de  $X_1$  sur  $X_2$ .

$$\begin{cases} y = X_2\gamma + v & \rightarrow & u = y - X_2\hat{\gamma} \\ X_1 = X_2\Gamma + Y & \rightarrow & V = X_1 - X_2\hat{\Gamma} \end{cases}$$

$$\Rightarrow u = V\beta_1 + \eta \quad \rightarrow \quad \widehat{\beta}_1 = (V'V)^{-1}V'u$$

qui peut être réécrit comme :

$$\widehat{\beta}_1 = (X_1'M_2X_1)^{-1}X_1'M_2y \quad \text{avec} \quad M_2 = I_N - X_2(X_2'X_2)^{-1}X_2'$$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

89

89

### DEMONSTRATION :

Considérons d'abord cette deuxième méthode.

On régresse la variable dépendante  $y$  sur les variables  $X_2$  :

$$y = X_2\gamma + v \quad \rightarrow \quad \hat{\gamma} = (X_2'X_2)^{-1}X_2'y$$

Le résidu est alors égal à :  $u = y - X_2\hat{\gamma} = y - X_2(X_2'X_2)^{-1}X_2'y = M_2y$

$$\text{avec } M_2 = I_N - X_2(X_2'X_2)^{-1}X_2'$$

On effectue aussi la régression (multivariée) de  $X_1$  sur  $X_2$  :

$$X_1 = X_2\Gamma + Y \quad \rightarrow \quad \hat{\Gamma} = (X_2'X_2)^{-1}X_2'X_1$$

ce qui donne pour les résidus :  $V = X_1 - X_2\hat{\Gamma} = X_1 - X_2(X_2'X_2)^{-1}X_2'X_1 = M_2X_1$

En fait, on projette  $y$  et  $X_1$  sur le plan orthogonal à celui engendré par  $X_2$ .

On va maintenant régresser les résidus de la première régression ( $u$ ) sur les résidus de la seconde régression ( $V$ ) :  $u = V\beta_1 + \eta$

Ce qui donne :  $\widehat{\beta}_1 = (V'V)^{-1}V'u = (X_1'M_2X_1)^{-1}X_1'M_2y$

On a enlevé de la variable dépendante  $y$  et des régresseurs  $X_1$ , l'effet des autres variables explicatives  $X_2$  (*partialling-out*).

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

90

90

Calculons maintenant la régression globale :

$$y = X\beta + \varepsilon = \begin{pmatrix} X_1 & X_2 \end{pmatrix} \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \varepsilon$$

Si on reprend les équations normales des moindres carrés, on aura :

$$\begin{pmatrix} X_1'X_1 & X_1'X_2 \\ X_2'X_1 & X_2'X_2 \end{pmatrix} \begin{pmatrix} \widehat{\beta}_1 \\ \widehat{\beta}_2 \end{pmatrix} = \begin{pmatrix} X_1'y \\ X_2'y \end{pmatrix}$$

Pour résoudre ce système, on peut utiliser les résultats de calcul matriciel sur les matrices partitionnées, ou procéder par substitution :

Prenons le second groupe d'équations :  $(X_2'X_1)\widehat{\beta}_1 + (X_2'X_2)\widehat{\beta}_2 = X_2'y$

dont la solution en fonction de  $\widehat{\beta}_1$  est :  $\widehat{\beta}_2 = (X_2'X_2)^{-1}X_2'y - (X_2'X_2)^{-1}(X_2'X_1)\widehat{\beta}_1$

$\uparrow$   $\uparrow$   
*Régression* *Correction dépendante*  
*de y sur X<sub>2</sub>* *de la corrélation*  
*entre X<sub>1</sub> et X<sub>2</sub>*

91

On substitue alors ce résultat dans le premier groupe d'équations normales :

$$(X_1'X_1)\widehat{\beta}_1 + (X_1'X_2)\widehat{\beta}_2 = X_1'y$$

pour donner :

$$\begin{aligned} (X_1'X_1)\widehat{\beta}_1 + (X_1'X_2)[(X_2'X_2)^{-1}X_2'y - (X_2'X_2)^{-1}(X_2'X_1)\widehat{\beta}_1] &= X_1'y \\ (X_1'X_1)\widehat{\beta}_1 + (X_1'X_2)(X_2'X_2)^{-1}X_2'y + (X_1'X_2)(X_2'X_2)^{-1}(X_2'X_1)\widehat{\beta}_1 &= X_1'y \\ (X_1'X_1)\widehat{\beta}_1 + (X_1'X_2)(X_2'X_2)^{-1}(X_2'X_1)\widehat{\beta}_1 &= X_1'y - (X_1'X_2)(X_2'X_2)^{-1}X_2'y \\ X_1'(I_N - X_2(X_2'X_2)^{-1}X_2')X_1\widehat{\beta}_1 &= X_1'(I_N - X_2(X_2'X_2)^{-1}X_2')y \end{aligned}$$

Du fait de la définition de la matrice de projection :  $M_2 = I_N - X_2(X_2'X_2)^{-1}X_2'$ , on obtient alors :

$$(X_1'M_2X_1)\widehat{\beta}_1 = X_1'M_2y$$

Ce qui donne finalement exactement le même estimateur que précédemment, du fait que la matrice  $(X_1'M_2X_1)$  est inversible :

$$\widehat{\beta}_1 = (X_1'M_2X_1)^{-1}X_1'M_2y$$

**→ Les deux méthodes sont équivalentes. CQFD**

92

Le théorème de Frisch – Waugh permet de neutraliser l'influence d'un groupe de variables explicatives  $X_2$ , pour se concentrer uniquement sur les paramètres d'intérêt  $\beta_1$ .

On a enlevé de la variable dépendante  $y$  et des régresseurs  $X_1$ , l'effet des autres variables explicatives  $X_2$

En anglais on appelle cette opération : *partialling – out* ou *netting – out*.

Ce théorème de Frisch – Waugh a une portée plus théorique que pratique... mais il permet d'effectuer de nombreuses démonstrations des propriétés des moindres carrés !

### **b) Exemple 1 : Le centrage des variables par rapport à leur moyenne**

Soit  $X_2$  le vecteur de la constante :  $X_2 = j_N = \begin{pmatrix} 1 \\ 1 \\ \vdots \\ 1 \end{pmatrix}$

on aura pour la matrice de projection  $M_2$  sur l'espace orthogonal engendré par ce vecteur de constante :

$$M_2 = I_N - j_N(j_N'j_N)^{-1}j_N' = I_N - \frac{j_N j_N'}{N} = I_N - \frac{J_N}{N}$$

avec  $J_N$  : une matrice carrée  $N \times N$  composée entièrement de 1.

Cette matrice  $M_2$  s'écrit alors :

$$M_2 = I_N - \frac{J_N}{N} = \begin{bmatrix} 1 - \frac{1}{N} & -\frac{1}{N} & -\frac{1}{N} & \dots & -\frac{1}{N} \\ -\frac{1}{N} & 1 - \frac{1}{N} & -\frac{1}{N} & \dots & -\frac{1}{N} \\ -\frac{1}{N} & -\frac{1}{N} & 1 - \frac{1}{N} & \dots & -\frac{1}{N} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ -\frac{1}{N} & -\frac{1}{N} & -\frac{1}{N} & \dots & 1 - \frac{1}{N} \end{bmatrix}$$

En prémultipliant un vecteur  $\mathbf{y}$  par cette matrice  $\mathbf{M}_2$ , on **centre ce vecteur par rapport à sa moyenne**  $\bar{y} = \mathbf{j}'_N \mathbf{y} / N$  :

$$\mathbf{M}_2 \mathbf{y} = \left( \mathbf{I}_N - \frac{\mathbf{J}_N}{N} \right) \mathbf{y} = \mathbf{y} - \mathbf{j}_N \frac{\mathbf{j}'_N \mathbf{y}}{N} = \mathbf{y} - \bar{y} \mathbf{j}_N = \begin{pmatrix} y_1 - \bar{y} \\ y_2 - \bar{y} \\ \vdots \\ y_N - \bar{y} \end{pmatrix} = \tilde{\mathbf{y}}$$

L'estimateur des  $K - 1$  paramètres d'intérêt d'un modèle où les variables sont centrées  $\tilde{\mathbf{y}} = \tilde{\mathbf{X}}_1 \boldsymbol{\beta}_1 + \tilde{\mathbf{v}} \rightarrow \mathbf{M}_2 \mathbf{y} = \mathbf{M}_2 \mathbf{X}_1 \boldsymbol{\beta}_1 + \mathbf{M}_2 \mathbf{v}$  est par le théorème de Frisch-Waugh :

$$\widehat{\boldsymbol{\beta}}_1 = (\tilde{\mathbf{X}}_1' \tilde{\mathbf{X}}_1)^{-1} \tilde{\mathbf{X}}_1' \tilde{\mathbf{y}} = (\mathbf{X}_1' \mathbf{M}_2 \mathbf{X}_1)^{-1} \mathbf{X}_1' \mathbf{M}_2 \mathbf{y}$$

**Donc l'estimation des paramètres de pente n'est pas modifiée si on centre toutes les variables par rapport à leurs moyennes.**

*Remarquez qu'il ne serait pas nécessaire de « centrer » la variable dépendante  $y$  si on centre uniquement les variables explicatives. Pourquoi ?*

Autres exemples : élimination d'une tendance,  
ou d'effets saisonniers additifs,...

Il suffit d'éliminer ces effets de chacune des variables explicatives d'intérêt.

Cependant : Intérêt pratique du théorème quasiment nul actuellement !  
*Sauf modèle à effets fixes sur données de panel ...*

L'intérêt principal de ce théorème est d'ordre théorique...



**c) Exemple 2 : L'introduction d'une variable supplémentaire**

On reprend la régression partielle avec une seule variable supplémentaire pour la variable  $X_1 = \mathbf{z}$  et  $X_2 = \mathbf{X}$  :

$$\mathbf{y} = (\mathbf{X}_1 \quad \mathbf{X}_2) \begin{pmatrix} \beta_1 \\ \boldsymbol{\beta} \end{pmatrix} + \varepsilon \quad \rightarrow \quad \mathbf{y} = (\mathbf{z} \quad \mathbf{X}) \begin{pmatrix} \gamma \\ \boldsymbol{\beta} \end{pmatrix} + \varepsilon$$

L'utilisation du théorème de Frisch-Waugh permet d'estimer  $\gamma$  séparément de  $\boldsymbol{\beta}$  :

$$\hat{\gamma} = (\mathbf{z}' \mathbf{M}_X \mathbf{z})^{-1} \mathbf{z}' \mathbf{M}_X \mathbf{y} \quad \text{avec : } \mathbf{M}_X = \mathbf{I}_N - \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$$

Ce qui donne en indiquant par une étoile (\*) les variables pré-multipliées par  $\mathbf{M}_X$  :

$$\hat{\gamma} = (\mathbf{z}'_* \mathbf{z}_*)^{-1} \mathbf{z}'_* \mathbf{y}_* = \frac{\mathbf{z}'_* \mathbf{y}_*}{\mathbf{z}'_* \mathbf{z}_*} \quad \text{avec } \mathbf{z}_* = \mathbf{M}_X \mathbf{z} \text{ et } \mathbf{y}_* = \mathbf{M}_X \mathbf{y}$$

C'est l'estimateur des MCO d'une régression simple de  $\mathbf{y}_*$  sur la seule variable explicative  $\mathbf{z}_*$  :

$$\mathbf{y}_* = \gamma \mathbf{z}_* + \mathbf{v} \quad \rightarrow \quad \mathbf{y}_* = \hat{\gamma} \mathbf{z}_*$$

$\mathbf{y}_* = \mathbf{M}_X \mathbf{y}$  est le résidu de la régression de  $\mathbf{y}$  sur les variables explicatives  $\mathbf{X}$ .

$\mathbf{z}_* = \mathbf{M}_X \mathbf{z}$  est le résidu de la régression de  $\mathbf{z}$  sur les variables explicatives  $\mathbf{X}$ .

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

97

97

$$\mathbf{y}_* = \gamma \mathbf{z}_* + \mathbf{v}$$

C'est la régression partielle après avoir éliminé l'influence des variables explicative  $\mathbf{X}$  de la variable dépendante  $\mathbf{y}$  et de la variable explicative d'intérêt  $\mathbf{z}$ .

$$\hat{\gamma} = (\mathbf{z}'_* \mathbf{z}_*)^{-1} \mathbf{z}'_* \mathbf{y}_*$$

De même, selon le théorème de Frisch-Waugh, l'estimateur MCO de  $\boldsymbol{\beta}$  dans le modèle complet  $\hat{\boldsymbol{\beta}}$  peut se réécrire en fonction de  $\boldsymbol{\beta}^*$  (l'estimateur MCO de  $\boldsymbol{\beta}$  sans la variable  $\mathbf{z}$ ) :

$$\boldsymbol{\beta}^* = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y}$$

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{y} - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\hat{\gamma} = \boldsymbol{\beta}^* - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\hat{\gamma}$$

Remarquez que si  $\hat{\gamma} \neq 0$ , les estimateurs sont égaux ( $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}^*$ ) seulement si  $\mathbf{X}'\mathbf{z} = \mathbf{0}$ , c'est-à-dire s'il n'y a pas de covariance (corrélation) entre les variables  $\mathbf{X}$  et la variable  $\mathbf{z}$  !

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

98

98

Notons le résidu de cette **régression « complète »** avec les variables  $\mathbf{X}$  et  $\mathbf{z}$  :

$$\mathbf{e} = \mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}} - \mathbf{z}\hat{\gamma}$$

alors que le résidu de la **régression « simplifiée »** avec seulement les variables  $\mathbf{X}$  :

$$\mathbf{e}_* = \mathbf{y} - \mathbf{X}\boldsymbol{\beta}^* = \mathbf{M}_X \mathbf{y} = \mathbf{y}_*$$

En introduisant l'estimateur  $\hat{\boldsymbol{\beta}} = \boldsymbol{\beta}^* - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\hat{\gamma}$  dans l'expression du résidu  $\mathbf{e}$  :

$$\begin{aligned}\mathbf{e} &= \mathbf{y} - \mathbf{X}(\boldsymbol{\beta}^* - (\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\hat{\gamma}) - \mathbf{z}\hat{\gamma} \\ &= \mathbf{y} - \mathbf{X}\boldsymbol{\beta}^* + \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'\mathbf{z}\hat{\gamma} - \mathbf{z}\hat{\gamma} \\ &= \mathbf{e}_* - \mathbf{M}_X \mathbf{z}\hat{\gamma} = \mathbf{e}_* - \mathbf{z}_* \hat{\gamma}\end{aligned}$$

La somme des carrés des résidus dans la régression « complète » devient alors :

$$\begin{aligned}SCR &= \mathbf{e}'\mathbf{e} = (\mathbf{e}_* - \mathbf{z}_* \hat{\gamma})'(\mathbf{e}_* - \mathbf{z}_* \hat{\gamma}) \\ &= \mathbf{e}_*'\mathbf{e}_* + \hat{\gamma}^2 \mathbf{z}_*'\mathbf{z}_* - 2\hat{\gamma} \mathbf{z}_*'\mathbf{e}_* \\ &= SCR^* - \hat{\gamma}^2 \mathbf{z}_*'\mathbf{z}_*\end{aligned}$$

parce que :  $\mathbf{z}_*'\mathbf{e}_* = \mathbf{z}_*'\mathbf{y}_* = \mathbf{z}_*'\hat{\gamma}\mathbf{z}_* = \hat{\gamma}\mathbf{z}_*'\mathbf{z}_*$

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

99

99

**En conséquence, la SCR du modèle « complet » est inférieure ou égale à la SCR\* du modèle « simplifié » :  $SCR = SCR^* - \hat{\gamma}^2 \mathbf{z}_*'\mathbf{z}_* \leq SCR^*$**

En insérant  $\hat{\gamma} = \mathbf{z}_*'\mathbf{y}_* / \mathbf{z}_*'\mathbf{z}_*$ , on obtient finalement :

$$SCR^* - SCR = \hat{\gamma}^2 \mathbf{z}_*'\mathbf{z}_* = \frac{(\mathbf{z}_*'\mathbf{y}_*)^2}{(\mathbf{z}_*'\mathbf{z}_*)} = r_*^2 (\mathbf{y}_*'\mathbf{y}_*)$$

où  $r_*$  est la **corrélation partielle** entre la variable dépendante  $\mathbf{y}$  et la variable explicative  $\mathbf{z}$ , **en neutralisant les effets des autres variables explicatives  $\mathbf{X}$**  :

$$r_* = \text{Corr}(\mathbf{y}_*, \mathbf{z}_*) = \frac{(\mathbf{z}_*'\mathbf{y}_*)}{\sqrt{(\mathbf{z}_*'\mathbf{z}_*)(\mathbf{y}_*'\mathbf{y}_*)}}$$

Quand on ajoute une variable explicative supplémentaire  $\mathbf{z}$  à un modèle, Le gain en termes d'ajustement, ou la diminution de la SCR dépend de la corrélation partielle entre  $\mathbf{z}$  et  $\mathbf{y}$ , compte tenu des autres variables explicatives  $\mathbf{X}$ .

Il faut que la variable supplémentaire  $\mathbf{z}$  apporte une information nouvelle qui n'est pas encore contenue dans les variables explicatives  $\mathbf{X}$ .

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

100

100

Dans le modèle « simplifié », le coefficient de détermination est :

$$R_*^2 = 1 - \frac{SCR_*}{SCT}$$

et dans le modèle « complet » :  $R^2 = 1 - \frac{SCR}{SCT}$

Comme  $SCR \leq SCR_*$ , on aura :  $1 - R^2 = \frac{SCR}{SCT} \leq \frac{SCR_*}{SCT} = 1 - R_*^2$   
 $-R^2 \leq -R_*^2$   
 $R^2 \geq R_*^2$

Lorsqu'on ajoute une variable explicative supplémentaire, le  $R^2$  du nouveau modèle « complet » augmente par rapport au  $R_*^2$  du modèle « simplifié ».

$$R^2 = R_*^2 + r_*^2 \frac{\mathbf{y}'_* \mathbf{y}_*}{SCT} = R_*^2 + r_*^2 \frac{\mathbf{y}' \mathbf{M}_X \mathbf{y}}{SCT}$$

### **d) Corollaire : la décomposition du $R^2$**

(Voir Henri THEIL (1971) : *Principles of Econometrics*, John Wiley & Sons, Chapitre IV.)

On a vu que le coefficient de détermination pouvait s'écrire :

$$R^2 = 1 - \frac{SCR}{SCT} = 1 - \frac{\mathbf{y}' \mathbf{M} \mathbf{y}}{\mathbf{y}' \mathbf{M}_0 \mathbf{y}}$$

$$R^2 = \frac{SCE}{SCT} = \frac{\widehat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{M}_0 \mathbf{X} \widehat{\boldsymbol{\beta}}}{\mathbf{y}' \mathbf{M}_0 \mathbf{y}}$$

Supposons que les variables soient centrées (il suffit d'avoir une constante dans le modèle pour obtenir la même propriété), alors :

$$\mathbf{y}' \mathbf{M}_0 \mathbf{y} = \mathbf{y}' \mathbf{y} \quad \text{et} \quad \widehat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{M}_0 \mathbf{X} \widehat{\boldsymbol{\beta}} = \widehat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{X} \widehat{\boldsymbol{\beta}} = \widehat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{y}$$

On aura alors pour le  $R^2$  :  $R^2 = \frac{\widehat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{X} \widehat{\boldsymbol{\beta}}}{\mathbf{y}' \mathbf{y}} = \frac{\widehat{\boldsymbol{\beta}}' \mathbf{X}' \mathbf{y}}{\mathbf{y}' \mathbf{y}}$

qui peut s'exprimer sous la forme :  $R^2 = \frac{\sum_{k=1}^K \widehat{\beta}_k \sum_{i=1}^N x_{k,i} y_i}{\sum_{i=1}^N y_i^2}$

$$R^2 = \frac{\sum_{k=1}^K \widehat{\beta}_k \sum_{i=1}^N x_{k,i} y_i}{\sum_{i=1}^N y_i^2} = \sum_{k=1}^K \widehat{\beta}_k \frac{\sum_{i=1}^N x_{k,i} y_i}{\sum_{i=1}^N y_i^2}$$

Pour une variable explicative, le ratio  $\widehat{\beta}_k \frac{\sum_{i=1}^N x_{k,i} y_i}{\sum_{i=1}^N y_i^2}$  pourrait être considéré comme la contribution de la  $k^{ième}$  variable explicative au  $R^2$ .

**MAIS** ce ratio peut être négatif pour une variable explicative !

Ce qui empêche toute utilisation de cette mesure comme contribution individuelle de la  $k^{ième}$  variable explicative au  $R^2$ .

### Explorons une autre voie...

On a vu plus haut que lorsqu'on ajoute une variable explicative supplémentaire, le  $R^2$  du nouveau modèle « complet » augmente par rapport au  $R_*^2$  du modèle « simplifié » :

$$R^2 = R_*^2 + r_*^2 \frac{\mathbf{y}' \mathbf{y}_*}{SCT} = R_*^2 + r_*^2 \frac{\mathbf{y}' \mathbf{M}_X \mathbf{y}}{SCT}$$

Avec :  $R_*^2$  : le  $R^2$  du modèle « simplifié » sans la  $k^{ième}$  variable explicative,  
 $r_*^2$  : le coefficient de corrélation partielle entre cette  $k^{ième}$  variable explicative et la variable dépendante  $y$ ,  
 $SCT$  : la somme des carrés totaux  $SCT = \mathbf{y}' \mathbf{y}$ ,  
 $\mathbf{M}_X = \mathbf{I}_N - \mathbf{X}(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}'$  (sans la  $k^{ième}$  variable explicative)

On obtient alors :

$$R^2 = R_*^2 + r_*^2 \frac{\mathbf{y}' \mathbf{M}_X \mathbf{y}}{\mathbf{y}' \mathbf{y}} = R_*^2 + r_*^2 \frac{\mathbf{y}' \mathbf{y} - \mathbf{y}' \mathbf{X}(\mathbf{X}' \mathbf{X})^{-1} \mathbf{X}' \mathbf{y}}{\mathbf{y}' \mathbf{y}} = R_*^2 + r_*^2 \left( 1 - \frac{\mathbf{y}' \mathbf{X} \widehat{\boldsymbol{\beta}}}{\mathbf{y}' \mathbf{y}} \right)$$

$$R^2 = R_*^2 + r_*^2 (1 - R_*^2)$$

Finalement :  $R^2 - R_*^2 = r_*^2(1 - R_*^2)$

*Voir Theil, formule 3.12, p.175.*

Ceci représente la contribution incrémentale (marginale) au  $R^2$  de la  $k^{ième}$  variable explicative à l'ajustement de la régression.

C'est le produit de deux facteurs positifs au plus égal à 1 :

- le carré de la corrélation partielle de la  $k^{ième}$  variable explicative avec la variable dépendante  $y$ ,
- la part inexpliquée (par les variables de la régression « simplifiée ») de la variation de la variable dépendante  $y$ .

On peut aussi ré-écrire cela sous la forme (*Theil, formule 3.1, p.172*) :

$$1 - R^2 = (1 - r_*^2)(1 - R_*^2)$$

## I.7. Les observations influentes

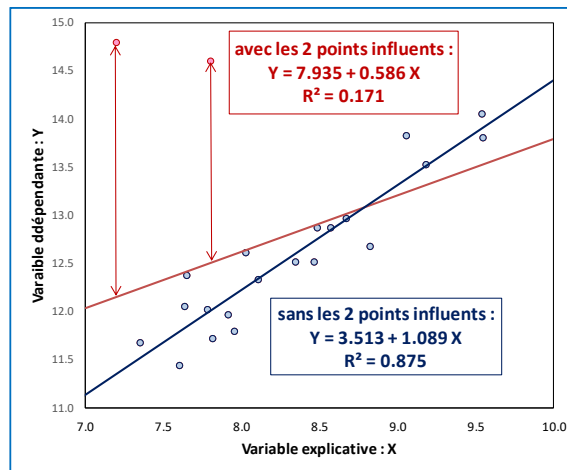
### a) L'importance d'une observation : le levier

Dans la méthode d'estimation des MCO, plus l'erreur d'une observation est grande, plus elle influence les paramètres estimés.

$$\hat{\beta} = \arg \min_{\beta} \sum_{i=1}^N \varepsilon_i^2$$

La méthode des MCO est sensible aux **observations extrêmes** (outliers) avec une forte erreur au détriment des observations avec une erreur faible ...

### Exemple



Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

111

111

### Peut-on détecter les observations les plus influentes dans les MCO ?

Dans quelle mesure l'estimation par MCO est-elle modifiée lorsqu'on retire ou on ajoute une observation ?

Voir le livre de David BELSLEY, Edwin KUH, et Roy WELSH (1980) : "Regression diagnostics: identifying influential data and sources of collinearity". New York: John Wiley & Sons.

#### Remarque :

- Si on ajoute une observation qui correspond à la moyenne de l'échantillon ( $\bar{y}, \bar{X}$ ), celle-ci est déjà sur le plan de régression (à condition d'avoir une constante).
- Cette observation n'influence pas la régression parce que son résidu sera nul !
- Donc il est important de savoir si une observation est éloignée de la moyenne de l'échantillon...

Benoît MULKAY  
Université de Montpellier

Econométrie Théorique (M1 MBFA)  
Chapitre 1 (2023 – 2024)

112

112

Estimateur *MCO* sur l'ensemble de l'échantillon :  $\hat{\beta} = (X'X)^{-1}X'y$

Si on retire la  $i^{\text{ème}}$  observation, on aura un estimateur des *MCO* :

$$\hat{\beta}_{(i)} = \hat{\beta} - \left( \frac{1}{1 - h_i} \right) (X'X)^{-1}x_i'e_i \quad \text{avec } e_i = y_i - x_i'\hat{\beta}$$

*(voir démonstration ci-dessous)*

**DÉFINITION** : le scalaire  $h_i$  est appelé le « levier » (**leverage**) de l'observation  $i$ .

Il est calculé comme le  $i^{\text{ème}}$  élément sur la diagonale principale de la matrice  $H$  :

$$H = X(X'X)^{-1}X' \rightarrow h_i = x_i'(X'X)^{-1}x_i$$

## DÉMONSTRATION

Si on enlève la  $i^{\text{ème}}$  observation de l'échantillon, on définit  $y_{(i)}$  le vecteur de la variable dépendante, et  $X_{(i)}$  la matrice des variables explicatives sans l'observation  $y_i$  et  $x_i$ .

En conséquence :

$$X'X = X'_{(i)}X_{(i)} + x_i x_i' \quad \text{et} \quad X'y = X'_{(i)}y_{(i)} + x_i y_i$$

On veut calculer l'inverse de  $X'_{(i)}X_{(i)} = X'X - x_i x_i'$  en utilisant la formule matricielle (*voir démonstration page suivante*) :

$$\text{si } W = A - BDC \quad \text{alors } W^{-1} = A^{-1} + A^{-1}B(D^{-1} - CA^{-1}B)^{-1}CA^{-1}$$

Si on note ici :  $W = X'_{(i)}X_{(i)}$ ,  $A = X'X$ ,  $B = x_i$ ,  $C = x_i'$ , et  $D = I_N$ , on aura :

$$\begin{aligned} (X'_{(i)}X_{(i)})^{-1} &= (X'X)^{-1} + (X'X)^{-1}x_i(1 - x_i'(X'X)^{-1}x_i)^{-1}x_i'(X'X)^{-1} \\ &= (X'X)^{-1} + \frac{1}{1 - h_i} (X'X)^{-1}x_i x_i' (X'X)^{-1} \end{aligned}$$

parce que  $h_i = x_i'(X'X)^{-1}x_i$

**Démonstration :** l'inverse de la matrice  $W = A - BDC'$  est :

$$W^{-1} = A^{-1} + A^{-1}B(D^{-1} - CA^{-1}B)^{-1}CA^{-1}$$

En effet :

$$\begin{aligned} WW^{-1} &= (A - BDC)(A^{-1} + A^{-1}B(D^{-1} - CA^{-1}B)^{-1}CA^{-1}) \\ &= I + B(D^{-1} - CA^{-1}B)^{-1}CA^{-1} - BDCA^{-1} - BDCA^{-1}B(D^{-1} - CA^{-1}B)^{-1}CA^{-1} \\ &= I - B[(D^{-1} - CA^{-1}B)^{-1} + D + DCA^{-1}B(D^{-1} - CA^{-1}B)^{-1}]CA^{-1} \\ &= I - B[D - (I - DCA^{-1}B)(D^{-1} - CA^{-1}B)^{-1}]CA^{-1} \\ &= I - B[D - D(D^{-1} - CA^{-1}B)(D^{-1} - CA^{-1}B)^{-1}]CA^{-1} \\ &= I - B[D - D]CA^{-1} \\ &= I \end{aligned}$$

**115**

Maintenant comme  $X'_{(i)}y_{(i)} = X'y - x_iy_i$ , l'estimateur des MCO sur l'échantillon sans la  $i^{\text{ème}}$  observation de l'échantillon sera :

$$\begin{aligned} \widehat{\beta}_{(i)} &= (X'_{(i)}X_{(i)})^{-1}X'_{(i)}y_{(i)} = \left[ (X'X)^{-1} + \frac{1}{1-h_i} (X'X)^{-1}x_ix'_i(X'X)^{-1} \right] (X'y - x_iy_i) \\ &= (X'X)^{-1}X'y - (X'X)^{-1}x_iy_i + \frac{1}{1-h_i} (X'X)^{-1}x_ix'_i(X'X)^{-1}X'y \\ &\quad - \frac{1}{1-h_i} (X'X)^{-1}x_ix'_i(X'X)^{-1}x_iy_i \\ \widehat{\beta}_{(i)} &= \widehat{\beta} - (X'X)^{-1}x_iy_i + \frac{1}{1-h_i} (X'X)^{-1}x_ix'_i\widehat{\beta} - \frac{1}{1-h_i} (X'X)^{-1}x_ix'_i(X'X)^{-1}x_iy_i \\ \widehat{\beta}_{(i)} &= \widehat{\beta} + \frac{1}{1-h_i} (X'X)^{-1}x_ix'_i\widehat{\beta} - (X'X)^{-1}x_iy_i - \frac{1}{1-h_i} (X'X)^{-1}x_ih_iy_i \\ \widehat{\beta}_{(i)} &= \widehat{\beta} + \frac{1}{1-h_i} (X'X)^{-1}x_ix'_i\widehat{\beta} - \left( 1 + \frac{h_i}{1-h_i} \right) (X'X)^{-1}x_iy_i \\ \widehat{\beta}_{(i)} &= \widehat{\beta} + \frac{1}{1-h_i} (X'X)^{-1}x_ix'_i\widehat{\beta} - \frac{1}{1-h_i} (X'X)^{-1}x_iy_i \\ \widehat{\beta}_{(i)} &= \widehat{\beta} - \frac{1}{1-h_i} (X'X)^{-1}x_i(y_i - x'_i\widehat{\beta}) = \widehat{\beta} - \frac{1}{1-h_i} (X'X)^{-1}x_ie_i \end{aligned}$$

**CQFD.**

**116**



La matrice de projection  $H = X(X'X)^{-1}X'$  est appelée la « **hat matrix** » en anglais...

En fait c'est notre matrice de projection  $P$  sur le plan engendré par les variables explicatives  $X$ .

Par analogie à la valeur calculée (*fitted*) de la variable dépendante :  $\hat{y}$

$$\hat{y} = X\hat{\beta} = X(X'X)^{-1}X'y = Hy$$

Le levier d'une observation sera élevé si celle-ci diffère de la majorité des observations...

Prenons le cas d'une régression simple :  $y_i = \alpha + \beta x_i + \varepsilon_i$ , le levier sera :

$$\begin{aligned} h_i &= x_i'(X'X)^{-1}x_i = (1 \quad x_i) \begin{pmatrix} N & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ x_i \end{pmatrix} \\ &= (1 \quad x_i) \frac{1}{N \sum_i x_i^2 - (\sum_i x_i)^2} \begin{pmatrix} \sum_i x_i^2 & -\sum_i x_i \\ -\sum_i x_i & N \end{pmatrix} \begin{pmatrix} 1 \\ x_i \end{pmatrix} \\ &= \frac{\sum_i x_i^2 - 2x_i \sum_i x_i + Nx_i^2}{N \sum_i (x_i - \bar{x})^2} \\ &= \frac{\sum_i x_i^2 - \frac{1}{N} (\sum_i x_i)^2 + \frac{1}{N} (\sum_i x_i)^2 - 2x_i \sum_i x_i + Nx_i^2}{N \sum_i (x_i - \bar{x})^2} \\ &= \frac{\sum_i (x_i - \bar{x})^2 + N(x_i^2 - 2x_i\bar{x} + \bar{x}^2)}{N \sum_i (x_i - \bar{x})^2} = \frac{1}{N} + \frac{(x_i - \bar{x})^2}{\sum_i (x_i - \bar{x})^2} \end{aligned}$$

Le levier d'une observation sera élevé si celle-ci diffère de la moyenne des observations.

Si  $x_i = \bar{x}$ , le deuxième terme est nul, et le levier sera :  $h_i = 1/N$ .

## Propriétés des leviers $h_i = \mathbf{x}_i'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i$

### 1) Les leviers sont compris entre 0 et 1 : $0 < h_i \leq 1$

C'est une forme quadratique avec une matrice définie positive  $(\mathbf{X}'\mathbf{X})^{-1}$ . En conséquence :  $h_i > 0$ .

De même,  $h_i \leq 1$ . Considérons le vecteur unitaire  $\mathbf{t}_i$  composé de zéros, sauf une valeur 1 pour le  $i^{\text{ème}}$  élément, et calculons :

$$h_i = \mathbf{x}_i'(\mathbf{X}'\mathbf{X})^{-1}\mathbf{x}_i = \mathbf{t}_i'\mathbf{P}\mathbf{t}_i > 0$$

avec  $\mathbf{H} = \mathbf{P} = \mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'$ . Mais on a vu que  $\mathbf{M} = \mathbf{I} - \mathbf{P}$ . Dès lors :

$$h_i = \mathbf{t}_i'\mathbf{P}\mathbf{t}_i = \mathbf{t}_i'(\mathbf{I} - \mathbf{M})\mathbf{t}_i = \mathbf{t}_i'\mathbf{t}_i - \mathbf{t}_i'\mathbf{M}\mathbf{t}_i = 1 - \mathbf{t}_i'\mathbf{M}\mathbf{t}_i$$

parce que  $\mathbf{t}_i'\mathbf{t}_i = 1$  par définition. Finalement comme  $\mathbf{t}_i'\mathbf{M}\mathbf{t}_i \geq 0$  parce que  $\mathbf{M}$  est une matrice semi-définie positive, on aura :

$$h_i = 1 - \mathbf{t}_i'\mathbf{M}\mathbf{t}_i \leq 1$$

### 2) La somme des leviers est aussi égale au nombre de variables explicatives :

$$\sum_{i=1}^N h_i = K.$$

Cette somme est la trace de la matrice  $\mathbf{H}$  :

$$\sum_{i=1}^N h_i = \text{tr}(\mathbf{H}) = \text{tr}[\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}\mathbf{X}'] = \text{tr}[\mathbf{X}'\mathbf{X}(\mathbf{X}'\mathbf{X})^{-1}]$$

par la propriété de circularité de la trace. En conséquence :

$$\text{tr}(\mathbf{H}) = \text{tr}[\mathbf{I}_K] = K.$$

Donc la moyenne des leviers est :  $\bar{h} = \frac{1}{N} \sum_{i=1}^N h_i = \frac{K}{N}$

### 3) Les leviers sont supérieurs ou égaux à l'inverse du nombre d'observations.

$$\sum_{i=1}^N h_i = K.$$

On partitionne :  $\mathbf{x}'_i = (1, \mathbf{z}'_i)$ . Sans perte de généralité, on peut remplacer  $\mathbf{z}_i$  par les variables centrées autour de leur moyenne :  $\tilde{\mathbf{z}}_i = \mathbf{z}_i - \bar{\mathbf{z}}$ .  
Dès lors, comme  $\tilde{\mathbf{z}}_i$  est orthogonal à la constante :

$$\begin{aligned} h_i &= \mathbf{x}'_i (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i = (1 \quad \tilde{\mathbf{z}}'_i) \begin{pmatrix} N & 0 \\ 0 & \tilde{\mathbf{Z}}'\tilde{\mathbf{Z}} \end{pmatrix}^{-1} \begin{pmatrix} 1 \\ \tilde{\mathbf{z}}_i \end{pmatrix} \\ &= (1 \quad \tilde{\mathbf{z}}'_i) \begin{pmatrix} 1/N & 0 \\ 0 & (\tilde{\mathbf{Z}}'\tilde{\mathbf{Z}})^{-1} \end{pmatrix} \begin{pmatrix} 1 \\ \tilde{\mathbf{z}}_i \end{pmatrix} \\ &= \frac{1}{N} + \tilde{\mathbf{z}}'_i (\tilde{\mathbf{Z}}'\tilde{\mathbf{Z}})^{-1} \tilde{\mathbf{z}}_i \geq \frac{1}{N} \end{aligned}$$

Dans la formule de l'estimateur *MCO* sans la  $i^{ème}$  observation :

$$\hat{\boldsymbol{\beta}}_{(i)} = \hat{\boldsymbol{\beta}} - \left( \frac{1}{1 - h_i} \right) (\mathbf{X}'\mathbf{X})^{-1} \mathbf{x}_i' e_i$$

Plus le levier de la  $i^{ème}$  observation sera proche de 1, plus l'observation aura du poids dans les MCO... et son absence modifiera fortement l'estimateur !

Elle sera donc très influente sur la régression.

Donc on peut calculer, pour chacune des  $N$  observations de l'échantillon, leurs leviers (leverage) respectifs  $h_i$  pour savoir quelles sont les observations les plus influentes...

On peut considérer que les observations avec un levier  $h_i > 2\bar{h} = 2K/N$  sont « influentes » et méritent une investigation plus poussée...

**Que faire** si on détecte une ou plusieurs observations extrêmes (**outliers**) ?

1. Eliminer l'observation ...
2. Corriger l'observation ...
3. Mettre une indicatrice pour cette observation ...  
*... cela équivaut à l'éliminer !*