

Les régressions linéaires

Le problème linéaire

Soit x et y denote deux variables, comme le revenu et la consommation ou l'âge et la taille d'un enfant.

Avec un échantillon $\{(y_i, x_i) \in \mathbb{R}, i = 1, \dots, N\}$ on souhaite résoudre le problème suivant:

Trouver $\hat{a}, \hat{b} \in \mathbb{R}$ tel q $\hat{y}_i = \hat{a} + \hat{b}x_i$ soit le plus proche possible de y pour tout i .

Ce modèle est dit :

“Simple” pour un seul régresseur x .

“Multiple” si il y en a plusieurs.

Impossibilité d'une équation exacte :

En général il n'existe pas de couple (a, b) solution de :

$$y_1 = a + bx_1$$

$$y_2 = a + bx_2$$

...

$$y_N = a + bx_N$$

Si on a plus de deux équations : $(N \geq 2)$.

La régression linéaire est donc une approximation.

L'estimateur OLS (ou moindres carrés) est un co (\hat{a}, \hat{b}) qui minimise le problème suivant :

$$SSR(a, b) = \sum_{i=1}^N (y_i - a - bx_i)^2$$

$$SSR(\hat{a}, \hat{b}) = \min_{a, b} \sum_{i=1}^N (y_i - a - bx_i)^2$$

Les conditions de premier ordre(FOC) de ce problème (i.e. fixer la dérivée partielle à 0) sont :

$$\frac{\partial SSR(\hat{a}, \hat{b})}{\partial a} = -2 \sum_{i=1}^N (y_i - \hat{a} - \hat{b}x_i) = 0$$

$$\frac{\partial SSR(\hat{a}, \hat{b})}{\partial b} = -2 \sum_{i=1}^N x_i (y_i - \hat{a} - \hat{b}x_i) = 0$$

($SSR(a, b)$ est convexe ; donc les FOCs sont suffisantes.)

Soit :

$$\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i, \quad \bar{y} = \frac{1}{N} \sum_{i=1}^N y_i, \quad \overline{x^2} = \frac{1}{N} \sum_{i=1}^N x_i^2, \quad \overline{xy} = \frac{1}{N} \sum_{i=1}^N x_i y_i$$

On obtient :

$$\frac{1}{N} \sum_{i=1}^N (y_i - \hat{a} - \hat{b}x_i) = \bar{y} - \hat{a} - \hat{b}\bar{x}$$

$$\frac{1}{N} \sum_{i=1}^N x_i (y_i - \hat{a} - \hat{b}x_i) = \overline{xy} - \hat{a}\bar{x} - \hat{b}\overline{x^2}$$

Ce qui impl

$$(\hat{a}, \hat{b})$$

$$\begin{cases} \bar{y} &= \hat{a} + \hat{b}\bar{x} \\ \overline{xy} &= \hat{a}\bar{x} + \hat{b}\overline{x^2} \end{cases}$$

Ce système est appelé système normal.

On multiplie la première équation par \bar{X} et on soustrait la seconde :

$$\overline{xy} - \bar{x}\bar{y} = \hat{b}(\overline{x^2} - \bar{x}^2) \Leftrightarrow \hat{b} = \frac{\overline{xy} - \bar{x}\bar{y}}{\overline{x^2} - \bar{x}^2} = \frac{\text{Cov}_N(x_i, y_i)}{\text{Var}_N(x_i)}$$

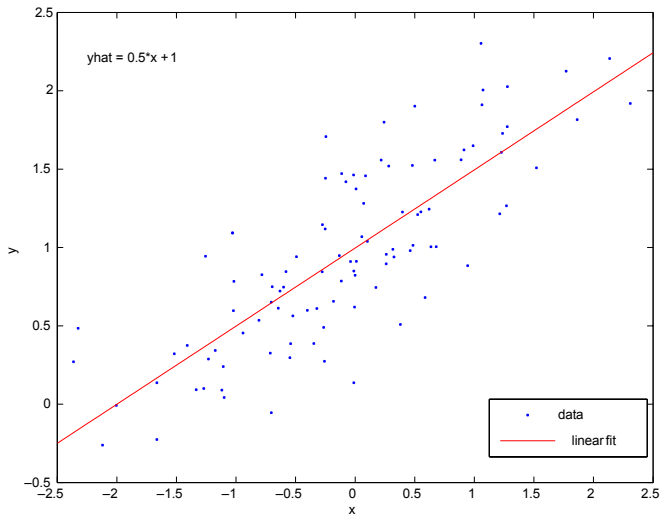
et

$$\hat{a} = \bar{y} - \hat{b}\bar{x}$$

Ceci est l'Estimateur « Ordinary Least Square » (OLS ou MCO en Français).

La ligne de régression : $y = ax + b$ est une droite passant par le barycentre du graphique (x, y)

Ligne de Regression



Pas de régresseur :

$$\min_a \sum_{i=1}^N (y_i - a)^2 \rightarrow \hat{a} = \frac{1}{N} \sum_{i=1}^N y_i = \bar{y}$$

Pas de constante (ordonnée à l'origine):

$$\min_b \sum_{i=1}^N (y_i - bx_i)^2 \rightarrow \hat{b} = \frac{\frac{1}{N} \sum_{i=1}^N x_i y_i}{\frac{1}{N} \sum_{i=1}^N x_i^2} = \frac{\overline{xy}}{\overline{x^2}}$$

On remplace chaque x et chaque y en leurs soustrayant leurs moyenne puis on régresse sans constante .

$$\hat{b} = \frac{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})(y_i - \bar{y})}{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2} = \frac{\text{Cov}_N(x_i, y_i)}{\text{Var}_N(x_i)}$$

Qui est l'estimateur de la pente de la droite de régression avec une constante.

Ajouter une constante revient au même que centrer toutes les variables de l'équation.

On remplace chaque x et chaque y en le divisant par son écart type :

$$\tilde{x}_i = \frac{x_i}{\sqrt{\text{Var}_N x_i}} \text{ and } \tilde{y}_i = \frac{y_i}{\sqrt{\text{Var}_N y_i}}$$

La variance de x et y vaut désormais 1

Cette méthode rend notre régression invariante aux changements d'échelles (par exemple changer la mesure entre des jours et des semaines, ou d'euros en dollars).

Dans ce cas :

$$\hat{b} = \frac{\text{Cov}_N(\tilde{x}_i, \tilde{y}_i)}{\text{Var}_N(\tilde{x}_i)} = \frac{\text{Cov}_N(\tilde{x}_i, \tilde{y}_i)}{1} = \frac{\text{Cov}_N(x_i, y_i)}{\sqrt{\text{Var}_N(x_i)}\sqrt{\text{Var}_N(y_i)}} = \text{Corr}_N(x_i, y_i)$$

L'OLS correspond à la corrélation.

On définit la prédiction par $\hat{y}_i = \hat{a} + \hat{b}x_i$ et le résidu $\hat{u}_i = y_i - \hat{y}_i$

Ainsi que : $SSR(\hat{a}, \hat{b}) = \sum_{i=1}^N \hat{u}_i^2$.

Le résidu moyen est 0 selon la première équation normale. Cela implique que la prédiction moyenne est égale à la moyenne des y_i :

$$\frac{1}{N} \sum_{i=1}^N \hat{y}_i = \hat{a} + \hat{b} \frac{1}{N} \sum_{i=1}^N x_i = \hat{a} + \hat{b}\bar{x} = \bar{y}$$

Corrélation avec les résidus

Les équations normales impliquent l'absence de corrélation des résidus:

$$\text{Cov}_N(x_i, \hat{u}_i) = \underbrace{\frac{1}{N} \sum_{i=1}^N x_i \hat{u}_i}_{=0} - \underbrace{\frac{1}{N} \sum_{i=1}^N x_i}_{=0} \frac{1}{N} \sum_{i=1}^N \hat{u}_i = 0$$

Ce qui implique que prédiction et résidus ne sont pas corrélés:

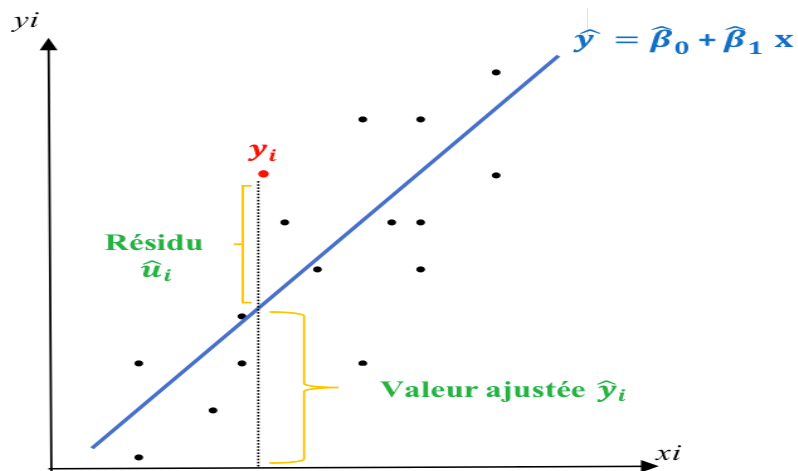
$$\begin{aligned} \text{Cov}_N(\hat{y}_i, \hat{u}_i) &= \frac{1}{N} \sum_{i=1}^N \hat{y}_i \hat{u}_i = \frac{1}{N} \sum_{i=1}^N (\hat{a} + \hat{b}x_i) \hat{u}_i \\ &= \hat{a} \frac{1}{N} \sum_{i=1}^N \hat{u}_i + \hat{b} \frac{1}{N} \sum_{i=1}^N x_i \hat{u}_i = 0 \end{aligned}$$

On dira que les résidus sont “orthogonaux” aux régresseurs et aux prédictions.

Formellement : le vecteur des résidus $(\hat{u}_1, \dots, \hat{u}_N)$ est orthogonal au vecteur des 1 $(1, \dots, 1)$ et (x_1, \dots, x_N) des x il le sera donc aussi au $(\hat{y}_1, \dots, \hat{y}_N)$ des y qui est une combinaison linéaire de x et de 1. La vérification se fait par le produit scalaire :

$$(x_1, \dots, x_N) \bullet (\hat{u}_1, \dots, \hat{u}_N) = \sum_{i=1}^N x_i \hat{u}_i = 0.$$

L'ajustement linéaire en graphique



Que nous apprennent les coefficients a et b :

Conséquences directes :

- Si x et y sont positivement corrélés alors b est positif
- Si x et y sont négativement corrélés alors b est négatif

Interprétation, 2 manières :

- A une valeur de x , on associe une valeur de y ($a+bx$).
- tcepa, si x varie de 1 unité, y varie de b unités.

Nb : ceteris paribus (ou tcepa) = suppose que tous que les autres facteurs (inclus dans le terme d'erreur) soient fixes.

- estimer l'effet de la variable explicative : ce qui explique les variations de y sont les variations de x

Quelle est la qualité de l'ajustement linéaire par rapport à y ?

Si $N = 2$ on peut résoudre parfaitement le système d'équations:

$$y_1 = a + bx_1$$

$$y_2 = a + bx_2$$

Mais avec $N \geq 3$ cela ne tient plus.

Comment évaluer la qualité de l'estimation ?

L'analyse de la variance (ANOVA)

Par définition : $y_i = \hat{y}_i + \hat{u}_i$ on va calculer la variance empirique de y_i :

$$\begin{aligned}\text{Var}_N y_i &= \text{Var}_N(\hat{y}_i + \hat{u}_i) \\ &= \text{Var}_N(\hat{y}_i) + \text{Var}_N(\hat{u}_i) + 2 \text{Cov}_N(\hat{y}_i, \hat{u}_i) \\ &= \text{Var}_N(\hat{y}_i) + \text{Var}_N(\hat{u}_i)\end{aligned}$$

Rappel : $\text{Cov}_N(\hat{y}_i, \hat{u}_i) = 0$.

On décompose ensuite cette variance en variance expliquée par le modèle et variance résiduelle :

- ▶ $\text{Var}_N(y_i) = \frac{1}{N} \sum_{i=1}^N (y_i - \bar{y})^2$
- ▶ $\text{Var}_N(\hat{y}_i) = \frac{1}{N} \sum_{i=1}^N (\hat{y}_i - \bar{y})^2$ because $\frac{1}{N} \sum_{i=1}^N \hat{y}_i = \frac{1}{N} \sum_{i=1}^N y_i = \bar{y}$
- ▶ $\text{Var}_N(\hat{u}_i) = \frac{1}{N} \sum_{i=1}^N \hat{u}_i^2$ because $\frac{1}{N} \sum_{i=1}^N \hat{u}_i = 0$

Decomposition en français :

Source de variation	Somme des carrés
Expliquée	$SCE = \sum_i (\hat{y}_i - \bar{y})^2$
Résiduelle	$SCR = \sum_i (y_i - \hat{y}_i)^2$
Totale	$SCT = \sum_i (y_i - \bar{y})^2$

$$SCT = SCR + SCE$$

- SCT : somme des carrés totaux
- SCR : somme des carrés résiduels, non expliqués par le modèle
- SCE : somme des carrés expliqués par le modèle

Le **coefficient de détermination** est

$$R^2 = \frac{\text{Variance expliquée}}{\text{total variance}} = \frac{SCE}{SCT} = 1 - \frac{\frac{SC}{R}}{\frac{SC}{T}} \in [0, 1]$$

Une bonne précision demande un R^2 élevé.

Le R^2 sera souvent plus élevé pour des séries temporelles que pour des données cross-sectionnelles parce que l'échantillon sera plus petit.

R^2 tends vers zero quand la taille de l'échantillon augmente, donc l'interprétation doit être éclairée.

Un R^2 faible ne veut pas forcément dire que le modèle est mauvais, on préférera comparer des modèles similaires.

- Deux types d'hypothèses
- Structurelles :
 - Sans lesquelles l'estimation n'est pas possible
 - Ou n'a pas de sens
- Comportementales:
 - Au sens du comportements des données
 - Les propriétés attendues des données
 - Qui permettent d'obtenir un estimateur désirable
- Ensemble = les hypothèses de Gauss-Markov



- Les hypothèses portent :
 - Sur les données, échantillon et variables (Y , X)
 - Sur le terme aléatoire
- Ces hypothèses pèsent :
 - Sur les propriétés des estimateurs (biais, convergence)
 - Sur l'inférence statistique (loi de distribution de l'estimateur)

- H1: Absence d'erreur de spécification
 - Pas d'erreur dans la liste des variables explicatives
 - La relation entre X et Y est linéaire
 - Sinon, pas d'estimation possible MCO
- H2: Echantillon aléatoire
 - Les probabilités d'inclusion/exclusion sont les mêmes pour tous les individus (sinon biais sélection ...)
- H3: Y est aléatoire
 - Via le terme d'erreur
 - La seule erreur que l'on commet provient des insuffisances dans les X à expliquer les valeurs de Y

Les hypothèses classiques

- H4: le nombre d'observations dans l'échantillon doit être supérieur au nombre de paramètres à estimer
- H5 : $V(x) > 0$
Les valeurs de X varient entre les individus,
Sinon on ne peut pas estimer l'impact de X sur Y
- H6: $E(x_i) = x_i$
Les variables explicatives sont certaines, exogènes (pas aléatoires)

- H7: $\text{Cov}(x_i, u_j) = 0$
 - Les erreurs sont indépendantes des variables explicatives
 - Pas de corrélation entre x et résidus, ce qui permet d'isoler l'effet des x_i
 - !!! Hypothèse d'exogénéité
- H8: $E(u_i) = 0$
 - La perturbation est d'espérance nulle
 - En moyenne les erreurs s'annulent (le modèle est bien spécifié)
 - Les aléas non pris en compte n'affectent pas la valeur moyenne de y

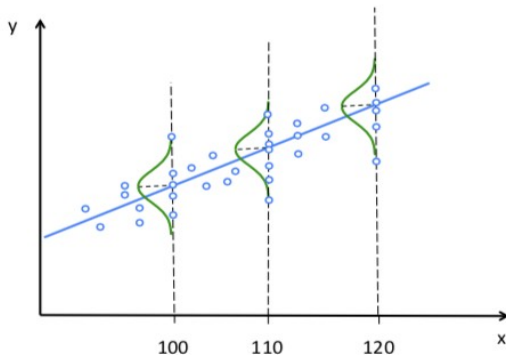
- H9a : $V(u_i) = \sigma^2$
 - Homoscédasticité
 - La variance des perturbations est constante, i.e. la même quelles que soient les valeurs des variables explicatives
- H9b : $\text{Cov}(u_i, u_j) = 0$
 - Non autocorrélation linéaire des perturbations
 - Une erreur faite sur une observation ne dépend linéairement pas d'une erreur faite sur une autre observation
 - Surtout pour les séries chronologiques

- Homoscédasticité
 - La variance des résidus doit être la même tout le long de la droite de régression,
 - Elle ne doit pas varier en fonction des valeurs de X :
 - L'écart moyen des points à la droite de régression est constant

$$Var(u|x_1, x_2, \dots, x_k) = \sigma^2$$

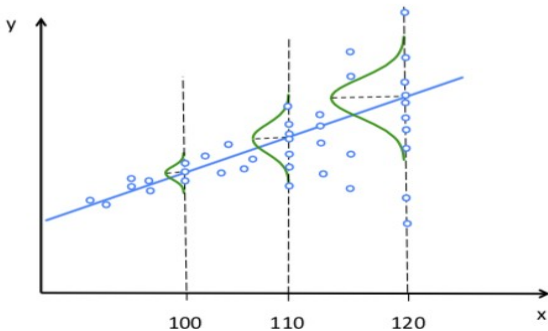
Les hypothèses classiques

$V(u)$ est constante pour toute valeur de x : les résidus sont **homoscédastiques**

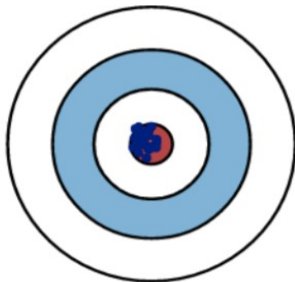


Les hypothèses classiques

$V(u)$ n'est pas constante pour toute valeur de x : les résidus sont **hétéroscédastiques**

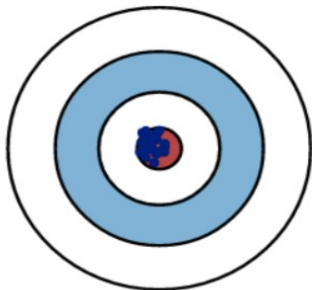


- Les hypothèses comportementales : des hypothèses à formuler en fonction des critères désirés pour l'estimateur
- Désirable =
 - Sans biais (espérance)
 - Précis (convergent, variance minimale)
 - De bonne famille (loi de distribution connue)



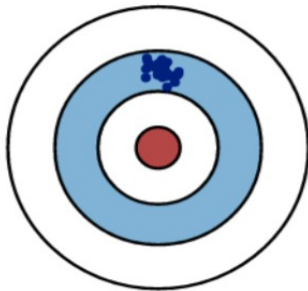
Légende

- Le centre de la cible représente la « vraie » valeur β du modèle en population
- Chaque point correspond à la valeur de $\hat{\beta}$ propre à un échantillon (l'étude_j)



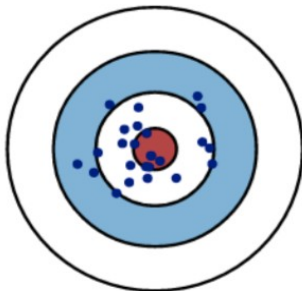
Faible biais, grande précision : l'estimateur idéal !

- Les estimations $\hat{\beta}$ sont proches de la « vraie » valeur β : $E(\hat{\beta}) \rightarrow \beta$
- Toutes ces estimations sont rapprochées : $V(\hat{\beta})$ est minimale



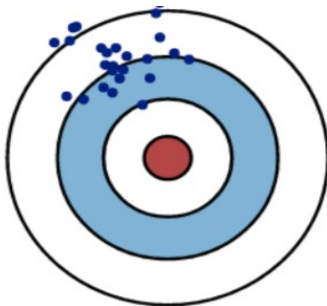
Fort biais, grande précision : l'estimateur faux que tout le monde trouve !

- Les estimations $\hat{\beta}$ sont loin de la « vraie » valeur β : $E(\hat{\beta}) \neq \beta$
- Mais toutes les estimations sont rapprochées : $V(\hat{\beta})$ est minimale



Faible biais, faible précision : l'estimateur est bon... dans une certaine mesure !

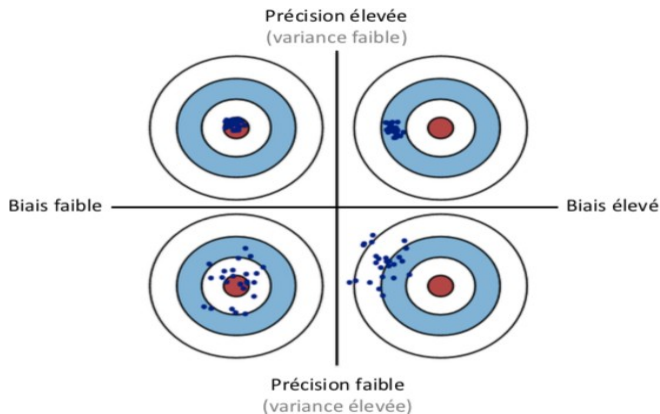
- Les estimations $\hat{\beta}$ sont proches de la « vraie » valeur β : $E(\hat{\beta}) \approx \beta$
- Mais toutes les estimations sont éloignées : $V(\hat{\beta})$ est grande



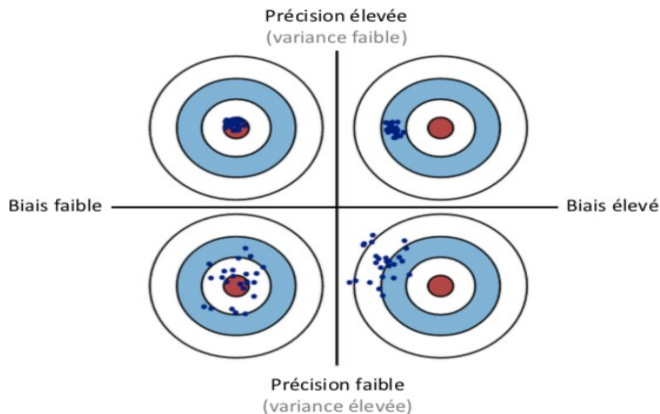
Fort biais, faible précision : le pire estimateur possible !

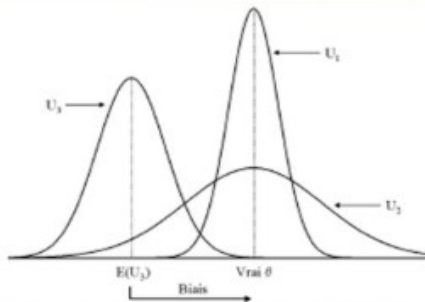
- Les estimations $\hat{\beta}$ sont éloignées de la « vraie » valeur β : $E(\hat{\beta}) \neq \beta$
- Et en plus, toutes les estimations sont éloignées : $V(\hat{\beta})$ est grande

Biais et précision



Biais et précision





Biais et variance pour 3 estimateurs d'un paramètre θ :
 U_1 et U_2 sont 2 estimateurs sans biais avec $Var(U_1) < Var(U_2)$
 U_3 est un estimateur biaisé

Quantifier la précision

La précision $V(\hat{\beta}_1) = \frac{\sigma^2}{nV(x)}$ depend
de :

La valeur de la variance des résidus = numérateur, Plus σ^2 est petite plus la précision est grande.

La taille de l'échantillon = dénominateur Plus n est grande plus $V\hat{\beta}_1$ est petite, plus la précision est grande
i.e. C'est bien un estimateur convergent

La taille de la variance de x :

- Plus V_x est grande (x dispersées), plus $V \hat{\beta}_1$ est petite et plus la précision est grande.

$$V(\hat{\beta}_1) = \frac{\sigma^2}{nV(x)}$$

- Repose sur H9
- Cette formule est fausse si l'hypothèse d'homoscédasticité n'est pas vérifiée
- L'homoscédasticité permet de simplifier la formule de la précision
et d'obtenir une variance minimale et donc l'estimateur le plus précis

- Best = variance minimale, le plus précis
 - Linear = estimateur linéaire
 - Unbiased = sans biais
 - Estimator = des MCO
- Il réunit 3 conditions : c'est un estimateur linéaire, sans biais, qui a la plus petite variance
- Sous certaines hypothèses : relation linéaire, X exogènes et variance >0 , espérance des perturbations nulle, homoscedasticité, etc.

- Sous certaines hypothèses ...
- Dans un modèle linéaire dans lequel les erreurs ont une espérance nulle, ne sont pas corrélées et de variance égales, le meilleur estimateur linéaire non biaisé des coefficients est l'estimateur MCO
 - Le meilleur estimateur non biaisé d'une combinaison linéaire des coefficients est son estimateur par les MCO
- Un des plus importants en économétrie
 - Avec les théorèmes asymptotiques (la LLN et le CLT)

Avec x_{N+1} comment prédire

$$\hat{y}_{N+1} = \hat{a} + \hat{b}x_{N+1}.$$

In-sample prediction: on choisit x_{N+1} qui est observé dans l'échantillon.

La capacité du modèle à estimer tout les x présents dans l'échantillon est une mesure de sa qualité, cependant cela n'est pas une mesure absolue mais plutôt relative. Par exemple si un trader peut prédire les mouvements de bourse un petit peu mieux que ses compétiteurs il gagnera beaucoup d'argent

Out-of-sample prediction: lorsque x_{N+1} n'a pas encore été observée. **Risqué!**

Une manière de faire est de couper le jeu de données en deux et de prédire la 2eme moitié avec la première et le modèle, puis de comparer.

https://fr.wikipedia.org/wiki/Conditions_d%27optimalit%C3%A9#:~:text=On%20parle%20de%20conditions%20du%20ordre%20des%20conditions%20font

<https://economictheoryblog.com/2014/11/05/proof/>