

# Régression linéaire simple

$Y$ : variable quantitative  
"à expliquer"  
endogène (dépendante)

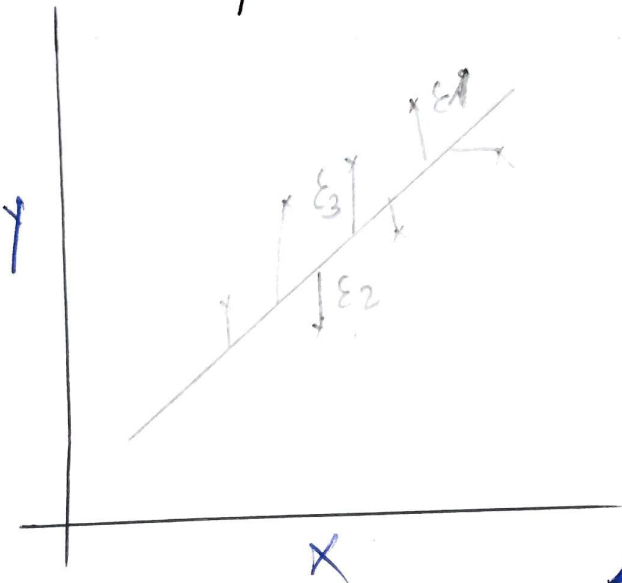
$X$ : variable quantitative  
"explicative"  
(exogène (effet fixe))

objectif

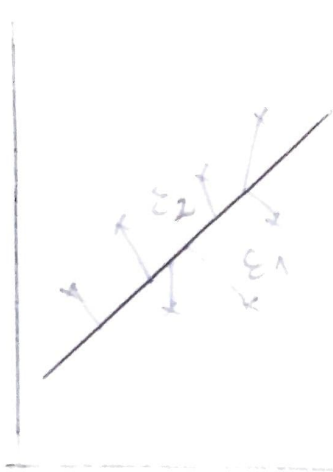
expliquer/exprimer une variable quantitative  $Y$  en fonction d'une variable quantitative  $X$ .

on suppose que la moyenne  $E[Y]$  est un fct linéaire.  
 $E[Y] = aX + b$

$$Y = aX + b + \varepsilon \quad (Y \text{ comme une fonction affine de } X).$$



4 "la droite affine ne passe pas nécessairement de tous les points"



$$Y_i$$

$$X_i$$

$$Y_i = b + a X_i + \varepsilon_i$$

modèle statistique

$\varepsilon_i$  : erreur qui représente les  
(Influence) autres facteurs qui peuvent influencer  
la variable  $Y$  autre que  $X$ .

↑ minimiser les  $\varepsilon_i$

## Hypothèses du modèle

1) la distribution de  $\varepsilon$  est  
indépendante de  $X$ .

[Pas de dépendance entre  $\varepsilon$  et  $X$ ]

$$E(\varepsilon_i | X_i) = 0 \quad H^0$$

2) Homoscedasticité : Erreur centrée  
et de variance  
constante.

$$E[\varepsilon_i] = 0, \quad V[\varepsilon_i] = \sigma^2$$

$H^0$  :  
↑  
s'annule au moyen  
- 2 -  
même  
variance.

3)

Pas de rupture du modèle :  
 $a$  et  $b$  sont deux constantes

4)

$\varepsilon_i$  indépendante.  
et  $\varepsilon_i \sim \mathcal{CR}(0, \sigma^2)$  H<sub>0</sub>.

$\Rightarrow$  Estimation des paramètres  
du modèle :  $a$ ?  $b$ ?

On cherche  $a$  et  $b$  ?.

Estimation par maximum  
de vraisemblance.

Estimation  
par les moindres  
carrés ordinaires

Plus facile  $\Rightarrow$  Moindres carrés.

minimiser la somme des carrés.

$$\min_{\vec{\theta}} \sum_{i=1}^n \varepsilon_i^2 = \min_{\vec{\theta}} \sum_{i=1}^n (y_i - a x_i - b)^2.$$

on cherche  $a$  et  $b$  tels que  $\sum_{i=1}^n \varepsilon_i^2$   
soit minimale.

On cherche a et b tq

- + la droite la "plus proche" du nuage des points
- + la somme des carrés des erreurs  $\varepsilon_i$  soit minimale.

$$\min_{(a,b)} \sum_{i=1}^n \varepsilon_i^2 = \min_{(a,b)} \sum_{i=1}^n (y_i - ax_i - b)^2.$$

On pose :  $F(a,b) = \sum_{i=1}^n (y_i - ax_i - b)^2$

$$\begin{cases} \frac{\partial F}{\partial a} = \sum_{i=1}^n -2(y_i - ax_i - b)x_i = -2 \sum_{i=1}^n x_i(y_i - ax_i - b) \\ \frac{\partial F}{\partial b} = \sum_{i=1}^n -2(y_i - ax_i - b) = -2 \sum_{i=1}^n (y_i - ax_i - b) \end{cases}$$

$$\begin{cases} \frac{\partial F}{\partial a} = 0 \\ \frac{\partial F}{\partial b} = 0 \end{cases} \Rightarrow \begin{cases} -2 \sum_{i=1}^n x_i(y_i - ax_i - b) = 0 \\ -2 \sum_{i=1}^n (y_i - ax_i - b) = 0 \end{cases}$$

$$\Rightarrow \begin{cases} \sum_{i=1}^n x_i(y_i - ax_i - b) = 0 \\ \sum_{i=1}^n (y_i - ax_i - b) = 0 \end{cases}$$



$$\begin{cases} \sum_{i=1}^n x_i y_i - a \sum_{i=1}^n x_i^2 - b \sum_{i=1}^n x_i = 0 \\ \sum_{i=1}^n y_i - a \sum_{i=1}^n x_i - nb = 0 \end{cases}$$

$$\begin{cases} \frac{1}{n} \sum_{i=1}^n x_i y_i - \frac{1}{n} a \sum_{i=1}^n x_i^2 - b \left( \frac{1}{n} \sum_{i=1}^n x_i \right) = 0 \\ \left( \frac{1}{n} \sum_{i=1}^n y_i \right) - a \left( \frac{1}{n} \sum_{i=1}^n x_i \right) - b = 0 \end{cases}$$

$\bar{y}$                        $\bar{x}$

$$\begin{cases} \frac{1}{n} \sum_{i=1}^n x_i y_i - a \frac{1}{n} \sum_{i=1}^n x_i^2 - b \bar{x} = 0 & (1) \\ \bar{y} - a \bar{x} - b = 0 & (2) \end{cases}$$

$$(2) \Rightarrow \boxed{b = \bar{y} - a\bar{x}}$$

$$(1) \quad \frac{1}{n} \sum_{i=1}^n x_i y_i - a \frac{1}{n} \sum_{i=1}^n x_i^2 - (\bar{y} - a\bar{x}) \bar{x} = 0$$

$$\frac{1}{n} \sum_{i=1}^n x_i y_i - a \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{y} \bar{x} + a \bar{x}^2 = 0$$

$$\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} = a \left( \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \right)$$

$$a = \frac{\frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y}}{\frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2} = \frac{\text{Cov}(X, Y)}{\sigma_x^2}$$

$$\left\{ \begin{aligned} \text{Cov}(X, Y) &= E[XY] - E[X]E[Y] \\ &= \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \bar{y} \\ \sigma_x^2 &= \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}^2 \end{aligned} \right.$$

a : pente de la droite

$$a = \frac{\text{Cov}(X, Y)}{\sigma_x^2}, \quad b = \bar{y} - a\bar{x}$$

$a$  et  $b$  sont les estimations.  
pour les réalisations.

$$\hat{y}_i = a x_i + b.$$

on définit :

Résidus

(résidus estimés)

$$e_i = y_i - \hat{y}_i$$

" la méthode des moindres carrés  
conduit  $\Rightarrow$  des résidus avec  
une somme des carrés minimales "

# Qualité d'ajustement

$y_i$  : observation réelle.

$\hat{y}_i$  : estimation par la méthode des moindres carrés.

$$e_i = y_i - \hat{y}_i$$


Une bonne qualité d'ajustement =

$$\sum_{i=1}^n e_i^2 = 0 \Rightarrow \sum_{i=1}^n (y_i - \hat{y}_i)^2 = 0$$

Somme des carrés résiduelle

$$SCR = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

Prédiction parfaite   $SCR = 0$

Mais  à partir de quelle valeur des SCR peut-on dire que la régression est mauvaise ?



On définit alors

{ la somme des carrés totale

$$SCT = SCR + SCE$$

←  
somme des  
carrés résiduelle

←  
somme  
des carrés  
expliquée

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

variabilité totale  
de  $y$ .

la différence  
entre les  $y_i$   
et la valeur de  
référence  $\bar{y}$  qui  
représente le centre  
de gravité du  
nuage de points

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

variabilité expliquée  
par le modèle, la  
variation de  $y$  expliquée  
par  $x$ .

la différence entre  
l'estimation et la  
valeur de référence

$$SCR = \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

différence entre les  
valeurs observées et les  
estimations.

variabilité non  
expliquée par  
le modèle :

Imaginant  
1) Situation extrême 1: Prédiction parfaite.

$$SCR = 0 \Rightarrow SCT = SCE.$$

"Les variations de  $y$  sont complètement  
expliquées par celles de  $X$ ."

→ modèle parfait 😊

2) Situation extrême 2:

$$SCE = 0 \rightarrow SCT = SCR.$$

→ \*

$M$  apporte aucune information  
sur  $y$ .

## "Coefficient de détermination"

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

"Variabilité totale expliquée par le modèle"

Plus que  $R^2 \rightarrow 1$ , plus que le modèle tend vers la perfection.

$$\circ) R^2 \rightarrow 1 \Rightarrow \frac{SCE}{SCT} \rightarrow 1$$
$$\Rightarrow SCE \rightarrow SCT$$

$$\circ) R^2 \rightarrow 1 \Rightarrow 1 - \frac{SCR}{SCT} \rightarrow 1 \rightarrow \frac{SCR}{SCT} \rightarrow 0$$
$$\Rightarrow SCR \rightarrow 0$$

$R^2$  proche de 1: modèle parfait  $\Rightarrow$  la connaissance de  $X$  nous permet de préciser  $Y$ .

$R^2$  proche de 0:  $X$  n'apporte pas d'information sur  $Y$ .

## Coefficient de Corrélation linéaire

$$R \text{ tq } a = R \frac{\sigma_y}{\sigma_x}.$$
$$R = \sqrt{R^2}.$$

$$a = \frac{\text{cov}(X, Y)}{\sigma_x^2} = R \frac{\sigma_y}{\sigma_x}.$$

$$\Rightarrow R = \frac{\text{cov}(X, Y)}{\sigma_x^2} \times \frac{\sigma_x}{\sigma_y}.$$

$$R = \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y}$$