# Régression Linéaire

- **La régression simple**

$$Y, X : \Omega \to \mathbb{R}$$

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

- **La régression multiple**

$$Y, \underbrace{X_1, X_2, \dots X_p}$$

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon$$
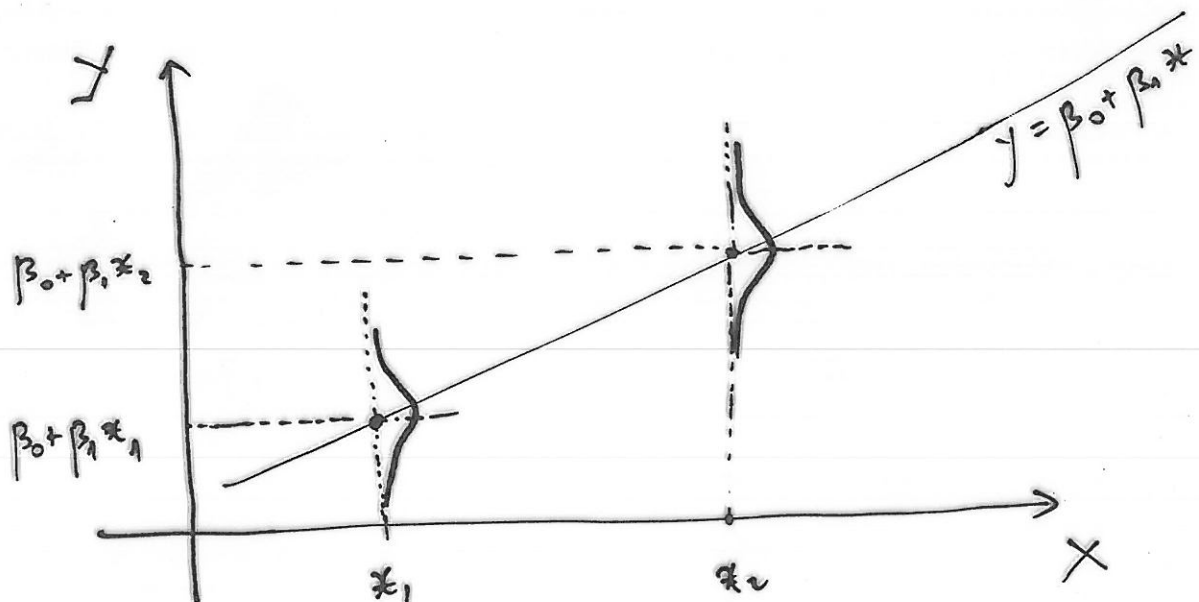
But :
  - estimation des $\{\beta_i\}_{i \geq 0}$
  - qualité du modèle

# La régression simple

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

Pour $X = x$ fixe

$$\left[ \begin{array}{l} Y \big|_{X=x} = \beta_0 + \beta_1 x + \underline{\varepsilon} \\[2mm] \qquad \begin{cases} E(\varepsilon) = 0 \\[1mm] Var(\varepsilon) = \sigma^2 \end{cases} \\[3mm] E(Y \big| X=x) = \beta_0 + \beta_1 x \end{array} \right.$$

Notons par $\hat{y} = \beta_0 + \beta_1 X$

Alors $y = \hat{y} + \varepsilon$

modèle $\qquad$ alea

## Estimation de $\beta_0, \beta_1, \sigma^2$

$(y_1, x_1), \quad \ldots \quad (y_n, x_n) \quad , \quad n \geq 1$

$$
\begin{cases}
y_1 = \beta_0 + \beta_1 x_1 + e_1 \\
y_2 = \beta_0 + \beta_1 x_2 + e_2 \\
\quad \vdots \\
y_n = \beta_0 + \beta_1 x_n + e_n
\end{cases}
$$

Moindres carrés :

chercher $\beta_0, \beta_1$ t.q :

$$
L(\beta_0, \beta_1) = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} \left[ y_i - (\beta_0 + \beta_1 x_i) \right]^2 = \sum_{i=1}^{n} (y_i - \hat{y}_i)^2
$$

$$\frac{\partial L}{\partial \beta_0} = 0 \qquad \frac{\partial L}{\partial \beta_1} = 0$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x} \qquad avec \begin{cases} \bar{y} = \frac{1}{n} \sum y_i \\ \bar{x} = \frac{1}{n} \sum x_i \end{cases}$$

$$\hat{\beta}_1 = \frac{\sum_{i=1}^{n} (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^{n} (x_i - \bar{x})^2}$$

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^{n} e_i^2$$

Th. : $\underline{\hat{\beta}_0, \hat{\beta}_1, \text{ et } \hat{\sigma}^2}$ sont estimateurs

sans biais de $\underline{\beta_0, \beta_1, \sigma^2}$.

$$\begin{cases} E(\hat{\beta}_0) = \beta_0 \\ \\ \text{Var}(\hat{\beta}_0) = \sigma^2 \cdot \left( \dfrac{1}{n} + \dfrac{\overline{x}^2}{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2} \right) \end{cases}$$

$$\begin{cases} E(\hat{\beta}_1) = \beta_1 \\ \\ \text{Var}(\hat{\beta}_1) = \sigma^2 \left( \dfrac{1}{\sum\limits_{i=1}^{n}(x_i - \overline{x})^2} \right) \end{cases}$$

## Remarque

Soit $\Delta^2_y = \dfrac{1}{n} \sum\limits_{i=1}^{n}(y_i - \overline{y})^2$

$\Delta^2_x = \dfrac{1}{n} \sum\limits_{i=1}^{n}(x_i - \overline{x})^2$

$\Delta^2_e = \dfrac{1}{n} \sum\limits_{i=1}^{n}(y_i - \hat{y}_i)^2 = \dfrac{1}{n} \sum\limits_{i=1}^{n} e_i^2$

$r = \dfrac{1/n \sum(x_i - \overline{x})(y_i - \overline{y})}{\frac{1}{n}\sqrt{\sum(x_i - \overline{x}_i)^2 \cdot \sum(y_i - \overline{y})^2}} = \dfrac{\text{cov}(x,y)}{\Delta_x \cdot \Delta_y}$

$$S_e^2 = (1 - r^2) S_y^2$$

ou

$$S_y^2 = S_e^2 + r^2 S_y^2$$

var totale $\quad$ var résiduelle $\quad$ variance expliquée par le modèle !

On a également :

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

# Inférence

$$\frac{\hat{\beta}_1 - \beta_1}{\sqrt{Var(\hat{\beta}_1)}} \sim T_{n-2}$$

$$\frac{\hat{\beta}_0 - \beta_0}{\sqrt{Var(\hat{\beta}_0)}} \sim T_{n-2}$$

---

Rem : Si $\rho = \dfrac{Cov(x,y)}{\sqrt{V(x) \cdot V(y)}} = 0$ alors $\beta_1 = 0$

$$\frac{\hat{\beta}_1}{\sqrt{Var(\hat{\beta}_1)}} \sim T_{n-2}$$

$$\Downarrow$$

$$\boxed{\frac{r}{\sqrt{1-r^2}} \sqrt{n-2} \sim T_{n-2}}$$

Test pour $r = 0$ !

# Tests dans le modèle linéaire

$$\begin{cases} H_0 : & \beta_1 = 0 \\[2mm] H_1 : & \beta_1 \neq 0 \end{cases}$$

- à l'aide du Student

- utilisant la décomposition de la variance.

$$\frac{1}{\sigma^2} \times \left| \quad \frac{1}{n} \sum_{i=1}^{n} (y_i - \bar{y})^2 = \frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2 + \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - \bar{y})^2 \right.$$

$$\downarrow \qquad\qquad\qquad \downarrow \qquad\qquad\qquad \downarrow$$

$$\chi^2_{n-1} \qquad\qquad \begin{array}{c} \text{si } \beta_1 = 0 \\ \chi^2 \end{array} \qquad\qquad \chi^2_{n-2}$$

$$\boxed{\frac{\displaystyle\sum_{i} (\hat{y}_i - \bar{y})^2}{\displaystyle\sum_{i=1}^{n} (y_i - \hat{y}_i)^2} \, (n-2) \;\sim\; F(1, n-2)}$$

$$\beta_1 = 0$$

$$\|$$

$$\frac{r^2}{1 - r^2} \, (n-2)$$

## Etude des résidus

$e_1, e_2, \dots e_n$

$e_1 + e_2 + \dots + e_n = 0$ \qquad (pas d'indépendance)

On teste généralement une tendance ou une dépendance entre $e_i$ et $e_{i+1}$

## Test de Durbin-Watson :

$$\begin{cases} H_0 : \text{il n'y a pas corrélations entre } \varepsilon_i \text{ et } \varepsilon_{i+1} \\ H_1 : \varepsilon_{i+1} = \rho \, \varepsilon_i + u_i \end{cases}$$

Statistique de test :

$$d = \frac{\sum_{i=2}^{n}(e_i - e_{i-1})^2}{\sum_{i=1}^{n} e_i^2} \underset{H_0}{\simeq} \text{près de } \underline{\underline{2}}$$

$$\downarrow$$

table statistique.

On vérifie que $0 \leq d \leq 4$

# Prévision

Soit $X = x_0$.

Alors

$$\hat{y}_0 = \hat{\beta}_0 + \hat{\beta}_1 x_0$$

$$\text{Var}(\hat{y}_0) = \sigma^2 \left( \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_i (x_i - \bar{x})^2} \right) \quad (\text{Démo}).$$

Mais $Y_0$ et $\hat{y}_0$ sont indépendantes !

$\Rightarrow Y_0 - \hat{y}_0$ :

$$\text{Var}(\underbrace{Y_0 - \hat{y}_0}_{\varepsilon}) = \sigma^2 \left( 1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum_i (x_i - \bar{x})^2} \right)$$

$$\boxed{ \frac{Y_0 - \hat{y}_0}{\hat{\sigma} \sqrt{1 + \frac{1}{n} + \frac{(x_0 - \bar{x})^2}{\sum (x_i - \bar{x})^2}}} \sim T_{n-2} }$$

$\underline{IC^{**}(Y_0)}$ : $\nearrow$ grand si $|x_0 - \bar{x}|$ grand.