

# Régression multiple

Soit la variable  $Y_i$  à expliquer par plusieurs variables quantitatives  $X_j$ :

$Y$ : variable à expliquer, quantitative

$X_1, X_2, \dots, X_k$ : variables explicatives, quantitatives

le Modèle s'écrit comme suit:

$$Y_{ij} = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_k X_{ik} + \underbrace{\varepsilon_i}_{\substack{\text{terme} \\ \text{aléatoire} \\ \text{les erreurs}}}$$

les Hypothèses à mettre en évidence pour estimer ce modèle:

$$\varepsilon_i \sim N(0, \sigma_i)$$

$\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ : ne sont pas corrélées entre elles.

on a:  $i$ : les individus  $j = 1 \dots k$  les variables

$$Y_1 = \beta_0 + \beta_1 X_{11} + \beta_2 X_{12} + \dots + \beta_k X_{1k} + \varepsilon_1$$

$$Y_2 = \beta_0 + \beta_1 X_{21} + \beta_2 X_{22} + \dots + \beta_k X_{2k} + \varepsilon_2$$

$$\vdots$$
$$Y_n = \beta_0 + \beta_1 X_{n1} + \beta_2 X_{n2} + \dots + \beta_k X_{nk} + \varepsilon_n$$

□

$$Y = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_m \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_k \end{pmatrix} \quad X = \begin{pmatrix} 1 & X_{11} & X_{12} & \dots & X_{1k} \\ & X_{21} & X_{22} & & X_{2k} \\ & & & & \\ & & & & \\ 1 & X_{m1} & X_{m2} & & X_{mk} \end{pmatrix}$$

$(m, 1) \qquad (k+1, 1) \qquad (k+1, m)$

$$\varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_m \end{pmatrix}$$

$(m, 1)$

le modèle s'écrit sous cette forme :

$$Y = \beta X + \varepsilon$$

on cherche l'ajustement qui minimise en  $\beta$  :

$$Q(\beta) = \sum \varepsilon_i^2 = \varepsilon' \varepsilon$$

donc pour  $\varepsilon \in \begin{pmatrix} \varepsilon \\ (m, 1) \end{pmatrix}$  existe il faut  $\varepsilon' \varepsilon$   
 $\begin{pmatrix} 1, m \end{pmatrix} \begin{pmatrix} m, 1 \end{pmatrix} =$

$$\varepsilon = Y - \beta X$$

$$\varepsilon' \varepsilon = (Y - \beta X)' (Y - \beta X)$$

$$Q(\beta) = Y'Y - Y'X\beta - \beta'X'Y + \beta'X'X\beta$$

on a  $Y'X\beta = \beta'X'Y$

$$Q(\beta) = Y'Y - 2Y'X\beta + \beta'X'X\beta$$



$$\frac{\partial Q(\beta)}{\partial \beta} = 0$$

$$\frac{\partial Y'Y}{\partial \beta} - \frac{\partial 2Y'X\beta}{\partial \beta} + \frac{\partial \beta'X'X\beta}{\partial \beta} = 0$$

$$0 - 2XY' + 2X'X\beta = 0$$

$$2X'X\beta = 2XY'$$

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

la commande R :  $\text{lm}(Y \sim X_1 + X_2 + X_3 + \dots + X_k)$

la fonction summary() : permet de produire les sorties suivantes :

- les coefficients estimés (associés à chaque variable explicative)
- leur écart-type (Std Error), la valeur de la statistique t de Student (t value) et la p-value ( $\text{Pr}(>|t|)$ ) : (la probabilité que le coefficient soit significativement différent de zéro) associées à chaque coefficient.

- $R^2$  ajusté (Adjusted R-squared) : Qualité d'ajustement du modèle

- la Statistique de Fisher (testant la significativité globales des variables), son degré de liberté et la p-value associé.

la fonction  $\text{coef}()$  : permet d'extraire les coefficients estimés

la fonction  $\text{fitted}()$  : permet d'extraire les valeurs prédites (estimées)  $\hat{y}_i$

la fonction  $\text{resid}()$  : permet d'extraire les résidus.  
Valeur prédite - valeur réelle =  $\hat{y}_i - y_i$

Plusieurs hypothèses sous-jacentes au modèle de régression linéaire portent sur les résidus de la régression :

- \* indépendance
- \* homoscedasticité (même variance)
- \* normalité (shapiro-test (résidus))