

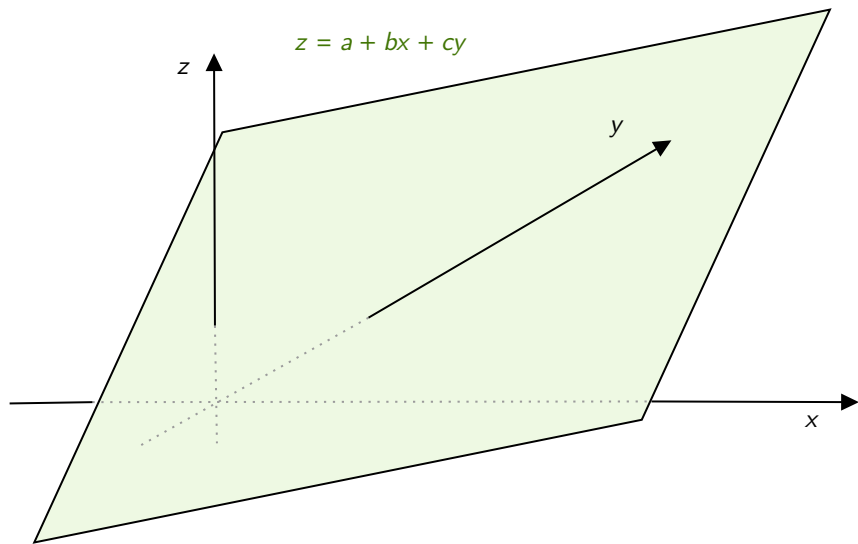
# Statistiques pour les sciences (MAT-4681)

Arthur Charpentier

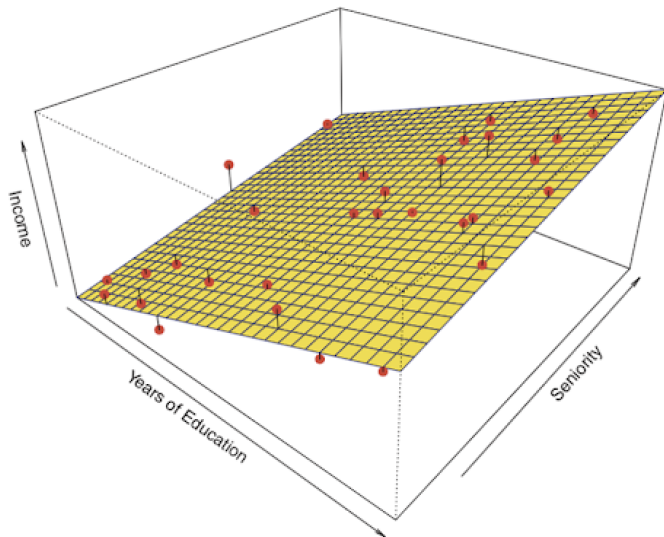
# 16 - Régression multiple

été 2022

## Plan (dans l'espace, $\mathbb{R}^3$ )



# Plan (dans l'espace, $\mathbb{R}^3$ )



# Moindres carrés

Soit  $(\mathbf{x}, \mathbf{y}) = \{(x_{1,1}, x_{2,1}, y_1), \dots, (x_{1,n}, x_{2,n}, y_n)\}$  un échantillon de trois variables. On suppose que

$$y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \varepsilon_i$$

- ▶  $y$  est la variable d'intérêt (que l'on veut prédire)
- ▶  $x_1$  et  $x_2$  sont deux variables explicatives (possibles)

On va chercher le plan qui passe au mieux dans le nuage de points,

$$\min_{\alpha, \beta} \left\{ \sum_{i=1}^n \varepsilon_i^2 \right\} = \min_{\beta_0, \beta_1, \beta_2} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1,i} - \beta_2 x_{2,i})^2 \right\}$$

# Moindres carrés

## Plan de régression, moindres carrées (OLS)

Soit  $(\mathbf{x}, \mathbf{y}) = \{(x_{1,1}, x_{2,1}, y_1), \dots, (x_{1,n}, x_{2,n}, y_n)\}$  un échantillon. Le plan de régression qui minimise la somme des carrés des erreurs est  $y = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2$  où

$$(\hat{\beta}_0, \hat{\beta}_1, \hat{\beta}_2) = \underset{\beta_0, \beta_1, \beta_2}{\operatorname{argmin}} \left\{ \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_{1,i} - \beta_2 x_{2,i})^2 \right\}$$

**Note** il existe une unique solution à ce programme d'optimisation

**Note** pour donner les valeurs des paramètres  $(\hat{\beta}_j)$  on va devoir passer par une représentation matricielle (cf MAT105 - 201-NYC)

# Matrices et vecteurs

Soient  $m, n \geq 1$ . Une matrice  $\mathbf{A}$  de taille  $(m, n)$  à coefficients réels est un tableau de nombres réels ayant  $m$  lignes et  $n$  colonnes. On note également par  $(\mathbf{A})_{ij}$  ou plus simplement  $A_{ij}$  l'élément sur la ligne  $i$  et sur la colonne  $j$  de  $\mathbf{A}$ .

**Example:**

$$\mathbf{A} = \begin{pmatrix} 1.5 & 2 & 3.1 & 8 \\ -1 & 4 & 5 & 6.5 \end{pmatrix}$$

$\mathbf{A}$  est de taille  $(2 \times 4)$  et par exemple  $A_{13} = 3.1$ .

Une matrice ne contenant qu'une colonne est appelée un vecteur et une matrice ne contenant qu'une ligne est un vecteur ligne. Par exemple  $\mathbf{x} = \begin{pmatrix} 1.5 \\ -1 \end{pmatrix}$  et  $\mathbf{y} = (1.5 \ 2 \ 3.1 \ 8)$  sont respectivement de taille  $(2, 1)$  et  $(1, 4)$ .

# Transposée

Soit  $\mathbf{A}$  une matrice réelle de taille  $(m, n)$ . La matrice transposée notée  $\mathbf{A}^T$  de taille  $(n, m)$  est définie par  $(\mathbf{A}^T)_{ij} = A_{ji}$  pour  $i = 1, \dots, n$  et  $j = 1, \dots, m$ .  
Et  $(\mathbf{A}^T)^T = \mathbf{A}$ .

$$\begin{pmatrix} 1 & 2 \\ 3 & 4 \\ 5 & 6 \end{pmatrix}^T = \begin{pmatrix} 1 & 3 & 5 \\ 2 & 4 & 6 \end{pmatrix}$$

$$\|\mathbf{x}\|^2 = \mathbf{x}^T \mathbf{x} = \sum_{i=1}^n x_i^2$$

$$\mathbf{X}^T \mathbf{X} = \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T$$

```
1 > t(1:4) %*% rep(1,4)
2      [,1]
3 [1,]    10
4 > (1:4) %*% t(rep(1,4))
5      [,1] [,2] [,3] [,4]
6 [1,]     1     1     1     1
7 [2,]     2     2     2     2
8 [3,]     3     3     3     3
9 [4,]     4     4     4     4
10 > t(1:4) %*% (1:4)
11      [,1]
12 [1,]    30
13 > (1:4) %*% t(1:4)
14      [,1] [,2] [,3] [,4]
15 [1,]     1     2     3     4
16 [2,]     2     4     6     8
17 [3,]     3     6     9    12
18 [4,]     4     8    12    16
```

# Transposée

Pour  $\mathbf{a}$  et  $\mathbf{b}$ , de dimension  $n$ ,

$$\mathbf{a}^\top \mathbf{b} = \mathbf{b}^\top \mathbf{a} = \sum_{i=1}^n a_i b_i$$

```
1 > a = c(1,2,3,4)
2 > b = c(7,2,4,1)
3 > t(a) %*% b
4      [,1]
5 [1,]    27
```

L'espérance peut s'écrire sous cette forme

$$\mathbb{E}[X] = \sum_{i=0}^n x_i p_i = \mathbf{x}^\top \mathbf{p}$$

```
1 > x = 0:6
2 > p = dbinom(0:6, 6, 1/3)
3 > t(x) %*% p
4      [,1]
5 [1,]    2
```

$$(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$$

$$(\mathbf{A} + \mathbf{B})^\top = \mathbf{A}^\top + \mathbf{B}^\top$$

Une matrice carrée  $\mathbf{A}$  de taille  $(n, n)$  est dite symétrique si  $\mathbf{A} = \mathbf{A}^\top$ .



## Produit ★★★

Si  $\mathbf{A}$  et  $\mathbf{B}$  sont (respectivement) des matrices  $k \times m$  et  $m \times n$ ,

$$C_{ij} = \mathbf{A}_{i\cdot}^\top \mathbf{B}_{\cdot j} = A_{i1}B_{1j} + \cdots + A_{im}B_{mj} = \sum_{k=1}^m A_{ik}B_{kj},$$

$$\begin{pmatrix} A_{11} & A_{12} & \cdots & A_{1m} \\ A_{21} & A_{22} & \cdots & A_{2m} \\ \vdots & \vdots & \ddots & \vdots \\ A_{n1} & A_{n2} & \cdots & A_{nm} \end{pmatrix} \cdot \begin{pmatrix} B_{11} & B_{12} & \cdots & B_{1p} \\ B_{21} & B_{22} & \cdots & B_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ B_{m1} & B_{m2} & \cdots & B_{mp} \end{pmatrix} = \begin{pmatrix} C_{11} & C_{12} & \cdots & C_{1p} \\ C_{21} & C_{22} & \cdots & C_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ C_{n1} & C_{n2} & \cdots & C_{np} \end{pmatrix}$$

```
1 > A = matrix(1:6,2,3)
2 > B = matrix(1:12,3,4)
3 > A %*% B
4      [,1] [,2] [,3] [,4]
5 [1,]    22    49    76   103
6 [2,]    28    64   100   136
```

Le produit matriciel n'est pas commutatif pour deux matrices quelconque de même taille:  $\mathbf{AB} \neq \mathbf{BA}$

## Produit et inverse ★★★

Soit  $\mathbb{I}_n$  la matrice de taille  $(n, n)$  composée de 1 sur la diagonale et de 0 ailleurs. Alors, pour  $\mathbf{A}$  de taille  $(n, n)$ ,  $\mathbb{I}_n$  est l'élément neutre tel que  $\mathbf{A}\mathbb{I}_n = \mathbb{I}_n\mathbf{A} = \mathbf{A}$ .

Soient  $\mathbf{A}$ ,  $\mathbf{B}$  et  $\mathbf{C}$  trois matrices réelles de dimension concordante, alors

- ▶  $(\mathbf{AB})\mathbf{C} = \mathbf{A}(\mathbf{BC})$  (associativité du produit)
- ▶  $\mathbf{A}(\mathbf{B} + \mathbf{C}) = \mathbf{AB} + \mathbf{AC}$  (distributivité du produit)
- ▶  $(\mathbf{AB})^\top = \mathbf{B}^\top \mathbf{A}^\top$ .

### Inverse matricielle

Soit  $\mathbf{A}$  une matrice carrée de taille  $(n, n)$  dont le déterminant est non nul, alors  $\mathbf{A}$  est dite non singulière et il existe une matrice inverse (de même taille) notée  $\mathbf{A}^{-1}$  vérifiant  $\mathbf{AA}^{-1} = \mathbf{A}^{-1}\mathbf{A} = \mathbb{I}_n$

# Produit et inverse ★★★

**Note** le déterminant ne sera pas redéfini ici

Soient **A** et **B** deux matrices inversibles de taille  $(n, n)$  alors

►  $(\mathbf{A}^{-1})^{\top} = (\mathbf{A}^{\top})^{-1}$

$(\mathbf{A}^{-1})$  est symétrique ssi **A** est symétrique)

►  $(\mathbf{A}^{-1})^{-1} = \mathbf{A}.$

►  $(\mathbf{AB})^{-1} = \mathbf{B}^{-1}\mathbf{A}^{-1}$

►  $\det(\mathbf{A}^{-1}) = \frac{1}{\det(\mathbf{A})}$

```
1 > A = matrix(  
2   c(3,2,4,3),2,2)  
3 > A  
4           [,1] [,2]  
5 [1,]         3    4  
6 [2,]         2    3  
7 > solve(A)  
8           [,1] [,2]  
9 [1,]         3   -4  
10 [2,]        -2    3  
11 > A %*% solve(A)  
12           [,1] [,2]  
13 [1,]         1    0  
14 [2,]         0    1
```

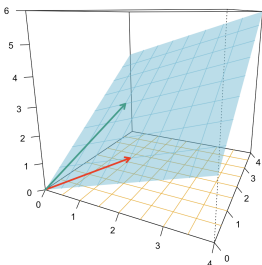
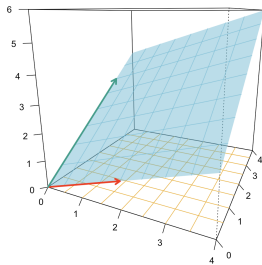
# Espace vectoriel engendré ★★★

Soient  $\mathbf{x}_1, \dots, \mathbf{x}_p \in \mathbb{R}^n$ , on définit  $\mathcal{V}(\mathbf{x}_1, \dots, \mathbf{x}_p)$  comme

$$\left\{ \mathbf{y} \in \mathbb{R}^n : \mathbf{y} = \sum_{i=1}^p a_i \mathbf{x}_i = \mathbf{X} \mathbf{a}, \mathbf{a} \in \mathbb{R}^p \right\}$$

où  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_p]$  est une matrice  $n \times p$ .

La dimension de  $\mathcal{V}(\mathbf{x}_1, \dots, \mathbf{x}_p)$  est le rang de  $\mathbf{X}$ .



# Régression Linéaire

Nous supposons que les données collectées suivent le modèle suivant

$$y_i = \beta_0 + \beta_1 x_{i1} + \cdots + \beta_k x_{ik} + \varepsilon_i,$$

où

- ▶  $x_{ij}$  sont des nombres déterministes (connus).  $\beta_0$  représente la constante (intercept dans les logiciels). On notera souvent  $x_{i0} = 1$ .
- ▶  $\beta_j, j = 0, 1, \dots, k$  paramètres réels à estimer. On pose  $p = k + 1$
- ▶ les variables  $\varepsilon_i$  sont des fluctuations aléatoires (erreur de mesures, mauvaise spécification du modèle, ...).

On peut reformuler le modèle en:

$$\underbrace{\mathbf{y}}_{(n \times 1)} = \underbrace{\mathbf{X}}_{(n \times p)} \underbrace{\boldsymbol{\beta}}_{(p \times 1)} + \underbrace{\boldsymbol{\varepsilon}}_{(n \times 1)}$$

# Régression Linéaire

On peut reformuler le modèle en:  $\underbrace{\mathbf{y}}_{(n \times 1)} = \underbrace{\mathbf{X}}_{(n \times p)} \underbrace{\boldsymbol{\beta}}_{(p \times 1)} + \underbrace{\boldsymbol{\varepsilon}}_{(n \times 1)}$  où

$1 \leq p = 1 + k \leq n$  et

$$\mathbf{y} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \mathbf{X} = \begin{pmatrix} x_{10} & x_{11} & \dots & x_{1k} \\ \vdots & \vdots & \vdots & \vdots \\ x_{n0} & x_{n1} & \dots & x_{nk} \end{pmatrix}, \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_k \end{pmatrix}, \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$\boldsymbol{\beta} \in \mathbb{R}^p$  et  $\mathbf{x}_j = (x_{1j}, \dots, x_{nj})^\top$ . Lorsque  $p = 2$ , ce modèle correspond au modèle de régression linéaire **simple**. On notera

- ▶  $\mathbf{x}_j$  le vecteur de taille  $(n \times 1)$  des  $n$  observations de la  $j$ ème covariable.
- ▶  $\mathbf{x}_i^\top$  le vecteur de taille  $(1 \times p)$  des valeurs des  $p$  covariables pour l'individu  $i$ .
- ▶  $\mathbf{y}$  le vecteur réponse;  $\boldsymbol{\varepsilon}$  vecteur aléatoire (centré sans perte de généralité).

# Moindres carrés ★★★

Formellement,

$\mathcal{H}_1$ : La matrice de design  $\mathbf{X}$  est de plein rang.

$p \leq n$ ,  $\mathcal{H}_1 \Rightarrow \text{rang}(\mathbf{X}) = p$ ,  $\mathbf{X}^\top \mathbf{X}$  de taille  $(p, p)$  est symétrique, définie positive et donc **inversible**.

$\mathcal{H}_2$ : Les erreurs sont centrées, de même variance et non corrélées  
 $\Leftrightarrow \mathbb{E}(\varepsilon_i) = 0$  et  $\text{Var}(\varepsilon_i) = \sigma^2$ .

$\mathcal{H}_2^{\mathcal{N}}$ : Les erreurs sont indépendantes et de même loi  $\mathcal{N}(0, \sigma^2)$

$\mathcal{H}_3$ : La matrice de design  $\mathbf{X}$  est telle que lorsque  $n \rightarrow \infty$ ,  
 $\frac{1}{n}(\mathbf{X}^\top \mathbf{X}) \rightarrow \mathbf{Q}$  où  $\mathbf{Q}$  est une matrice définie positive

- ▶  $\mathcal{H}_1$ : permet de démontrer l'existence de  $\hat{\beta}$ .
- ▶  $\mathcal{H}_2$ : permet de démontrer des propriétés pour  $\hat{\beta}$  (sans biais, calcul de variance).
- ▶  $\mathcal{H}_2^{\mathcal{N}}$ : permet de faire des tests.
- ▶  $\mathcal{H}_3$ : permet de démontrer la convergence de  $\hat{\beta}$ .

## Plan de régression, moindres carrées (OLS)

Soit  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  un échantillon. Le plan de régression qui minimise la somme des carrés des erreurs est  $y = \mathbf{x}^\top \hat{\boldsymbol{\beta}}$  où

$$\hat{\boldsymbol{\beta}} = \underset{\beta_0, \beta_1, \dots, \beta_k}{\operatorname{argmin}} \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \boldsymbol{\beta})^2$$

Sous l'hypothèse  $\mathcal{H}_1$ , l'estimateur des moindres carrées existe et vaut

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

**Note** il existe une unique solution à ce programme d'optimisation



# Coefficient de détermination $R^2$

$$R^2$$

Soit  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  un échantillon, et  $\hat{y}_i$  la prévision par régression linéaire. Alors

$$R^2 = 1 - \frac{\text{SCR}}{\text{SCT}} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

$$R_a^2 \text{ (} R^2 \text{ ajusté)}$$

Soit  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  un échantillon, et  $\hat{y}_i$  la prévision par régression linéaire. Alors

$$R_a^2 = 1 - \frac{\text{SCR}/(n - k - 1)}{\text{SCT}/(n - 1)} = 1 - \frac{n - 1}{n - k - 1} \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

# Coefficient de détermination $R^2$

$R_a^2$  et  $R^2$

$R_a^2 \leq R^2$  et

$$R_a^2 = 1 - (1 - R^2) \cdot \frac{n - 1}{n - k - 1}$$

On pénalise ici les modèles trop complexes, avec trop de variables explicatives.

# Prévision et résidus

## Prévision et résidus

Soit  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  un échantillon. La plande régression qui minimise la somme des carrés des erreurs est  $y = \mathbf{x}^\top \hat{\beta}$ . La différence entre la valeur observée  $y_i$  et la valeur prédite  $\hat{y}_i = \mathbf{x}_i^\top \hat{\beta}$  s'appelle le résidu  $\hat{\varepsilon}_i = y_i - \hat{y}_i$ .

## Résidus

Soient  $\hat{\varepsilon}_i = y_i - \hat{y}_i$  les résidus estimés. Les résidus sont centrés et leur variance  $\sigma^2$  est estimée par  $s^2$  où

$$s^2 = \frac{1}{n-p} \sum_{i=1}^n \hat{\varepsilon}_i^2 \text{ et } \sum_{i=1}^n \hat{\varepsilon}_i = 0.$$

# Regression

```
1 > import numpy as np
2 > import statsmodels.api as sm
3 > x = np.array([[5,1], [15,4], [25,-5], [35,4],
4               [45,-2], [55,2]])
5 > x = sm.add_constant(x)
6 > y = np.array([5, 20, 14, 32, 22, 38])
7 > model = sm.OLS(y, x)
8 > results = model.fit()
9 > print(results.summary())
```

```
=====
10                coef.  std err          t      P>|t|      [0.025      0.975]
11  -----
12 const          4.0581    3.370      1.204    0.315     -6.668     14.785
13 x1              0.5578    0.097      5.770    0.010      0.250      0.865
14 x2              1.5604    0.508      3.071    0.055     -0.057      3.178
15  -----
16 Dep. Variable:    y              R-squared:            0.931
17 Model:            OLS            Adj. R-squared:      0.886
18                               F-statistic:            20.37
```

# Regression

```
1 > df = data.frame(x1 = c(5, 15, 25, 35, 45, 55),
2                   x2 = c(1, 4, -5, 4, -2, 2),
3                   y = c(5, 20, 14, 32, 22, 38))
4 > model = lm(y~x1+x2, data=df)
5 > summary(model)
6
7 Coefficients:
8             Estimate Std. Error t value Pr(>|t|)
9 (Intercept)  4.05810    3.37049   1.204   0.3149
10 x1           0.55783    0.09667   5.770   0.0103 *
11 x2           1.56037    0.50818   3.071   0.0545 .
12 ---
13 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.'
14
15 Residual standard error: 4.037 on 3 degrees of freedom
16 Multiple R-squared:  0.9314, Adjusted R-squared:  0.8857
17 F-statistic: 20.37 on 2 and 3 DF,  p-value: 0.01796
```

## Propriétés de $\hat{\beta}$

Soit  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  un échantillon. Sous les hypothèses  $\mathcal{H}_1$  et  $\mathcal{H}_2$ ,

$$\mathbb{E}[\hat{\beta}] = \beta \text{ et } \text{Var}[\hat{\beta}] = \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1}$$

(admis)

On a vu (section 14) que si  $y = \alpha + \beta x$  (droite de régression)

$$s_{\hat{\beta}}^2 = \frac{\frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = s^2 \left( \sum_{i=1}^n (x_i - \bar{x})^2 \right)^{-1}$$

et ici  $s_{\hat{\beta}_j}^2 = s^2 (\mathbf{X}^\top \mathbf{X})_{jj}^{-1}$ .

# Propriétés de $\hat{\beta}$ ★★★

```
1 > library(DALEX)
2 > reg = lm(m2.price~construction.year+surface+no.rooms
  , data=apartments)
3 > vcov(reg)
4           (Intercept) const.year surface no.rooms
5 (Intercept) 3550207.63   -1807.629 152.913 -3277.951
6 const.year   -1807.629      0.921  -0.080    1.171
7 surface      152.913      -0.080    2.562   -64.046
8 no.rooms     -3277.951      1.171  -64.046  1922.300
9 > summary(reg)
10
11 Coefficients:
12             Estimate Std. Error t value Pr(>t)
13 (Intercept) 6295.7095  1884.1995   3.341 0.000865 ***
14 const.year   -0.8829    0.9599  -0.920 0.357920
15 surface     -9.3827    1.6007  -5.862 6.22e-09 ***
16 no.rooms    -80.6139   43.8440  -1.839 0.066264 .
17
18 Residual standard error: 781.8, 996 degrees of freedom
19 Multiple R-squared: 0.2588, Adjusted R-squared: 0.2566
20 F-statistic: 115.9 on 3 and 996 DF, p-value: < 2.2e-16
```

# Test (possiblement multiples) I

Test simple  $H_0 : \beta_j = 0$  contre  $H_1 : \beta_j \neq 0$

Soit  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  un échantillon. Si le plan de régression est  $y = \mathbf{x}^\top \boldsymbol{\beta}$ , pour tester  $H_0 : \beta_j = 0$  contre  $H_1 : \beta_j \neq 0$ , la statistique de test est

$$T = \frac{\hat{\beta}_j}{s_{\hat{\beta}_j}} \text{ où } s_{\hat{\beta}_j} = \sqrt{s^2 (\mathbf{X}^\top \mathbf{X})_{jj}^{-1}}.$$

Si  $H_0 : \beta = 0$  est vraie, et si  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ ,  $T \sim \text{Std}(n - p)$ .  
Et donc

► on rejette  $H_0$  si  $|t| > T_{n-p}^{-1}(1 - \alpha/2)$

où  $T_\nu$  est la fonction de répartition de la loi de Student  $\text{Std}(\nu)$



# Régression Linéaire

Sous les hypothèses  $\mathcal{H}_1$  et  $\mathcal{H}_2^{\mathcal{N}}$ , on a, pour  $j = 1, \dots, p$

$$T_j = \frac{\hat{\beta}_j - \beta_j}{s_{\hat{\beta}_j}} \sim Std_{n-p} \quad \text{où } s_{\hat{\beta}_j} = s \sqrt{(\mathbf{X}^T \mathbf{X})_{jj}^{-1}}.$$

```
1 > library(DALEX)
2 > reg = lm(m2.price~construction.year+surface+no.rooms
3   , data=apartments)
4 > summary(reg)
5 Coefficients:
6             Estimate Std. Error t value Pr(>t)
7 (Intercept) 6295.7095  1884.1995   3.341 0.000865 ***
8 const.year   -0.8829    0.9599  -0.920 0.357920
9 surface      -9.3827    1.6007  -5.862 6.22e-09 ***
10 no.rooms     -80.6139   43.8440  -1.839 0.066264 .
11
12 Residual standard error: 781.8, 996 degrees of freedom
13 Multiple R-squared: 0.2588, Adjusted R-squared: 0.2566
14 F-statistic: 115.9 on 3 and 996 DF, p-value: < 2.2e-16
```

# Régression Linéaire et test multiple

- ▶  $\xi$  un sous-ensemble d'indices  $\xi \subseteq \{1, \dots, p\}$  de cardinal  $|\xi|$ .
- ▶  $\bar{\xi}$  les indices du complémentaire de  $\xi$  dans  $\{1, \dots, p\}$ ,  
Rappel:  $\xi \cap \bar{\xi} = \emptyset$  et  $\xi \cup \bar{\xi} = \{1, \dots, p\}$
- ▶  $\mathbf{X}_\xi$  sous-matrice des covariables  $\mathbf{x}_j, j \in \xi$ .
- ▶  $\mathbf{X}_{\bar{\xi}}$  sous-matrice des covariables  $\mathbf{x}_j, j \in \bar{\xi}$  (ou  $j \notin \xi$ ).
- ▶  $\beta_\xi$  les paramètres dans le modèle  $(\xi)$  où seules les variables  $\xi$  sont conservées.
- ▶  $[\hat{\beta}]_\xi$ : coordonnées  $\xi$  du vecteur  $\hat{\beta}$ ,  
Note:  $[\hat{\beta}]_\xi \neq \hat{\beta}_\xi$  en général (sauf si  $\mathbf{X}_\xi \perp \mathbf{X}_{\bar{\xi}}$ ).

**Note:**  $\mathbf{u} = (u_1, \dots, u_k) = 0$  signifie  $\forall j, u_j = 0$ .

$\mathbf{u} = (u_1, \dots, u_k) \neq 0$  signifie  $\exists j$  tel que  $u_j \neq 0$ .

# Régression Linéaire et test multiple

Test  $H_0 : \beta_\xi = 0$  contre  $H_1 : \beta_\xi \neq 0$

Soit  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  un échantillon. Pour tester  $H_0 : \beta_\xi = 0$  contre  $H_1 : \beta_\xi \neq 0$  on estime

$$\begin{cases} \mathbf{y} = \mathbf{X}_{\bar{\xi}} \beta_{\bar{\xi}} + \varepsilon_{\bar{\xi}} & (0) \text{ régression contrainte} \\ \mathbf{y} = \mathbf{X} \beta + \varepsilon & (1) \text{ régression non-contrainte} \end{cases}$$

La statistique de test est

$$F = \frac{\text{SCR}(\beta_{\bar{\xi}}) - \text{SCR}(\beta)}{\text{SCR}(\beta)} \frac{n-p}{|\xi|} = \frac{n-p}{|\xi|} \frac{R^2 - R_0^2}{1 - R^2}$$

Si  $H_0 : \beta = 0$  est vraie, et si  $\varepsilon \sim \mathcal{N}(0, \sigma^2)$ ,  $F \sim \mathcal{F}(n-p, |\xi|)$ .  
Et donc

► on rejette  $H_0$  si  $f > F_{n-p, |\xi|}^{-1}(1 - \alpha)$ .

# Régression Linéaire et test multiple

Sur nos données sur le prix des logements en Pologne

$$\begin{cases} y_i = \beta_0 + \beta_2 x_{2,i} + \eta_i & (0) \\ y_i = \beta_0 + \beta_1 x_{1,i} + \beta_2 x_{2,i} + \beta_3 x_{3,i} + \varepsilon_i & (1) \end{cases}$$

```
1 > reg1 = lm(m2.price~construction.year+surface+no.  
  rooms, data=apartments)  
2 > reg0 = lm(m2.price~surface, data=apartments)  
3 > anova(reg0, reg1)  
4 Analysis of Variance Table  
5  
6 Model 1: m2.price ~ surface  
7 Model 2: m2.price ~ construction.year + surface + no.  
  rooms  
8 Res.Df      RSS Df Sum of Sq      F Pr(>F)  
9 1      998 611258600  
10 2      996 608730962  2    2527638 2.0678  0.127
```

# Régression Linéaire et test multiple

```
1 > linearHypothesis(reg1,  
2   c("construction.year = 0", "no.rooms = 0"))  
3 Linear hypothesis test  
4  
5 Hypothesis:  
6 construction.year = 0  
7 no.rooms = 0  
8  
9 Model 1: restricted model  
10 Model 2: m2.price ~ construction.year + surface + no.  
    rooms
```

	Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1	998	611258600				
2	996	608730962	2	2527638	2.0678	0.127

# Régression Linéaire et test multiple

Quand il y a beaucoup de variables, il est possible d'utiliser des méthodes de sélection de variables, les méthodes pas à pas, ou [step-wise](#).

Les algorithmes les plus simples consistent à faire rentrer les variables une à une (méthode ascendante, [forward](#)), ou les à les faire sortir une à une (méthode descendante, [backward](#)).

# Régression Linéaire et test multiple

## Sélection de variable pas à pas

Soit  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  un échantillon.

- ▶  $\xi, \xi \subseteq \{1, \dots, p\}$  de cardinal  $|\xi|$ .
- ▶  $\xi_{+1}, \xi_{+1} \subseteq \{1, \dots, p\}$  de cardinal  $|\xi| + 1$ .
- ▶  $\xi \subset \xi_{+1}$ , autrement dit  $\xi = \xi_{+1} \cup \{j\}, j \in \{1, \dots, p\}$ .

$$\begin{cases} \mathbf{y} = \mathbf{X}_{\xi_{+1}} \boldsymbol{\beta}_{\xi_{+1}} + \boldsymbol{\varepsilon}_{\xi_{+1}} & (\xi) \\ \mathbf{y} = \mathbf{X}_{\xi} \boldsymbol{\beta}_{\xi} + \boldsymbol{\varepsilon}_{\xi} & (\xi_{+1}) \end{cases}$$

On préfère  $(\xi)$  à  $(\xi_{+1})$  si  $R_a^2(\xi) > R_a^2(\xi_{+1})$ .

**Note** notion de [parcimonie](#) et [rasoir d'Ockham](#).

# Omission d'une variable explicative I

Oublier une variable importante peut avoir des conséquences importantes

- ▶  $y_i = \beta_0 + \mathbf{x}_1^\top \beta_1 + \mathbf{x}_2^\top \beta_2 + \varepsilon_i$ : le vrai modèle
- ▶  $y_i = b_0 + \mathbf{x}_1^\top \mathbf{b}_1 + \eta_i$ : le modèle que l'on considère

L'estimateur de  $\mathbf{b}_1$  est

$$\begin{aligned}\hat{\mathbf{b}}_1 &= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{y} \\&= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top [\mathbf{X}_1 \beta_1 + \mathbf{X}_2 \beta_2 + \varepsilon] \\&= (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_1 \beta_1 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \beta_2 + (\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \varepsilon \\&= \beta_1 + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \mathbf{X}_2 \beta_2}_{\beta_{12}} + \underbrace{(\mathbf{X}_1^\top \mathbf{X}_1)^{-1} \mathbf{X}_1^\top \varepsilon}_{\nu_i}\end{aligned}$$

de tel sorte que  $\mathbb{E}[\hat{\mathbf{b}}_1] = \beta_1 + \beta_{12} \neq \beta_1$ , en général.



# Omission d'une variable explicative II

Comme le montrait [Bickel, Hammel & O'Connell \(1975\)](#) (avec un modèle plus complexe car les variables ne sont pas ici continues)

- ▶  $y$  est l'admission aux études graduées
- ▶  $x_1$  est le genre (homme ou femme)
- ▶  $x_2$  est le programme où l'étudiant(e) a postulé

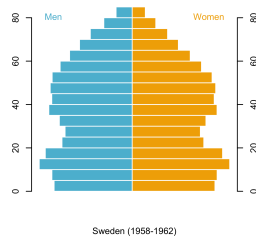
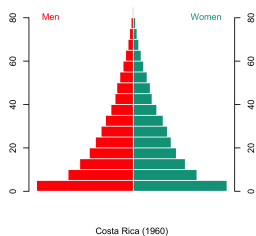
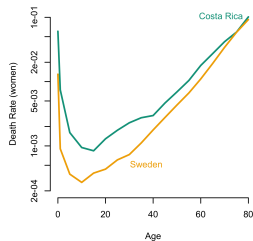
	Total	Men	Women	Proportions
Total	5233/12763 ~ 41%	3714/8442 ~ 44%	1512/4321 ~ 35%	66%-34%
Top 6	1745/4526 ~ 39%	1198/2691 ~ 45%	557/1835 ~ 30%	59%-41%
A	597/933 ~ 64%	512/825 ~ 62%	89/108 ~ 82%	88%-12%
B	369/585 ~ 63%	353/560 ~ 63%	17/ 25 ~ 68%	96%- 4%
C	321/918 ~ 35%	120/325 ~ 37%	202/593 ~ 34%	35%-65%
D	269/792 ~ 34%	138/417 ~ 33%	131/375 ~ 35%	53%-47%
E	146/584 ~ 25%	53/191 ~ 28%	94/393 ~ 24%	33%-67%
F	43/714 ~ 6%	22/373 ~ 6%	24/341 ~ 7%	52%-48%

# Omission d'une variable explicative III

- ▶  $y$  est la durée de vie résiduelle (en années)
- ▶  $x_1$  est le pays (Costa Rica ou Suède)
- ▶  $x_2$  l'âge de la personne

$$\mathbb{P}[Y \leq 1 | \mathbf{X} = \text{Costa Rica}] < \mathbb{P}[Y \leq 1 | \mathbf{X} = \text{Suède}]$$

$$\mathbb{P}[Y \leq 1 | \mathbf{X} = (\text{Costa Rica}, x)] > \mathbb{P}[Y \leq 1 | \mathbf{X} = (\text{Suède}, x)], \quad \forall x$$



## Intervalle de prédiction

Soit  $(\mathbf{X}, \mathbf{y}) = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$  un échantillon. Si le plan de régression est  $y = \mathbf{x}^\top \boldsymbol{\beta}$ . On dispose d'une nouvelle observation  $\mathbf{x}_{n+1}$ . L'intervalle de confiance de la valeur moyenne prédite est:

$$\left[ \mathbf{x}_{n+1}^\top \hat{\boldsymbol{\beta}} \pm t_{1-\alpha/2, n-p} s \sqrt{\mathbf{x}_{n+1}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{n+1}} \right]$$

L'intervalle de confiance pour une valeur particulière est:

$$\left[ \mathbf{x}_{n+1}^\top \hat{\boldsymbol{\beta}} \pm t_{1-\alpha/2, n-p} s \sqrt{1 + \mathbf{x}_{n+1}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{n+1}} \right]$$

# Confidence et prédiction

$$\left[ \mathbf{x}_{n+1}^T \hat{\beta} \pm t_{1-\alpha/2, n-p} s \sqrt{1 + \mathbf{x}_{n+1}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{n+1}} \right]$$

```
1 > predict(reg, newdata = data.frame(construction.year
   surface=80, no.rooms=3), interval = "
   confidence")
2         fit         lwr         upr
3 1 3544.494 3472.005 3616.983
4 > x=c(1,1992,80,3)
5 > t(x)%*%reg$coefficients
6         [,1]
7 [1,] 3544.494
8 > residus = reg$residuals
9 > t(x)%*%reg$coefficients+qt(c(.025,.975),n-4)*sqrt(
   sum(residus^2)/(n-4))*sqrt(t(x)%*%solve(t(X)%*%X)
   %*%x)
10 [1] 3472.005 3616.983
```

# Confidence et prédiction

$$\left[ \mathbf{x}_{n+1}^T \hat{\beta} \pm t_{1-\alpha/2, n-p} s \sqrt{1 + \mathbf{x}_{n+1}^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_{n+1}} \right]$$

```
1 > predict(reg, newdata = data.frame(construction.year
   =1992, surface=80, no.rooms=3), interval = "
   prediction")
2           fit           lwr           upr
3 1 3544.494 2008.663 5080.326
4 > x=c(1,1992,80,3)
5 > t(x)%*%reg$coefficients
6           [,1]
7 [1,] 3544.494
8 > residus = reg$residuals
9 > t(x)%*%reg$coefficients+qt(c(.025,.975),n-4)*sqrt(
   sum(residus^2)/(n-4))*sqrt(1+t(x)%*%solve(t(X)%*%X
   )%*%x)
10 [1] 2008.663 5080.326
```

# Sélection de variables

```
1 > library(olsrr)
2 > model = lm(mpg ~ disp + hp + wt + qsec, data =
  mtcars)
3 > ols_step_all_possible(model)
```

Index	N	Predictors	R-Square	Adj. R-Square
1	1	wt	0.7528328	0.7445939
2	1	disp	0.7183433	0.7089548
3	1	hp	0.6024373	0.5891853
4	1	qsec	0.1752963	0.1478062
5	2	hp wt	0.8267855	0.8148396
6	2	wt qsec	0.8264161	0.8144448
7	2	disp wt	0.7809306	0.7658223
8	2	disp hp	0.7482402	0.7308774
9	2	disp qsec	0.7215598	0.7023571
10	2	hp qsec	0.6368769	0.6118339
11	3	hp wt qsec	0.8347678	0.8170643
12	3	disp hp wt	0.8268361	0.8082829
13	3	disp wt qsec	0.8264170	0.8078189
14	3	disp hp qsec	0.7541953	0.7278591
15	4	disp hp wt qsec	0.8351443	0.8107212

# Sélection de variables

```
1 > model = lm(y ~ ., data = surgical)
2 > ols_step_forward_p(model)
```

## Selection Summary

-----				
Step	Variable Entered	R-Square	Adj. R-Square	RMSE
-----				
1	liver_test	0.4545	0.4440	296.2992
2	alc_heavy	0.5667	0.5498	266.6484
3	enzyme_test	0.6590	0.6385	238.9145
4	pindex	0.7501	0.7297	206.5835
5	bcs	0.7809	0.7581	195.4544
-----				