

Modèles Linéaires Appliqués

Arthur Charpentier

Automne 2020

OLS #4 (régression sur une variable continue - 3)

Poids & Taille

```
1 > str(Davis)
2 'data.frame': 200 obs. of 5 variables:
3 $ sex      : Factor w/ 2 levels "F","M": 2 1 ...
4 $ weight   : int 77 58 53 68 59 76 76 69 71 ...
5 $ height   : int 182 161 161 177 157 170 167 ...
6 > X = Davis$height
7 > Y = Davis$weight
8 > (B1hat = cor(X,Y) * sd(Y)/sd(X))
9 [1] 1.150092
10 > (B0hat = mean(Y) - B1hat*mean(X))
11 [1] -130.9104
```

$$y_i = \underbrace{-130.91}_{\hat{\beta}_0} + \underbrace{1.15}_{\hat{\beta}_0} x_i + \widehat{\varepsilon}_i$$

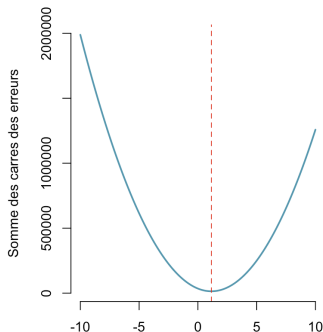
```
1 > plot(X,Y)
2 > abline(B0hat, B1hat)
```

Poids & Taille

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{s_{xy}}{s_x^2} = \text{corr}(\mathbf{x}, \mathbf{y}) \cdot \frac{s_y}{s_x}$$

$$\text{and } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

```
1 > (B1hat = cor(X,Y) * sd(Y)/sd(X))
2 [1] 1.150092
3 > (B0hat = mean(Y) - B1hat*mean(X))
4 [1] -130.9104
5 > SCR = fonction(B){sum(((Y-mean(Y)
6   ) -B*(X-mean(X)))^2)}
7 > optim(0,SCR)$par
8 [1] 1.15
9 > x = seq(-10,10,length=251)
10 > y = Vectorize(SCR)(x)
```



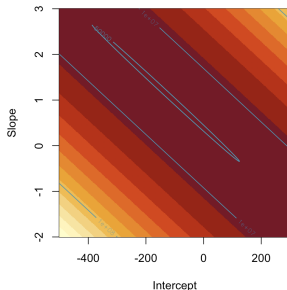
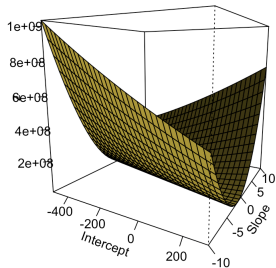
Note: $\hat{\beta}_1 = \operatorname{argmin} \left\{ \sum_{i=1}^n ((y_i - \bar{y}) - \beta_1(x_i - \bar{x}))^2 \right\}$

Poids & Taille

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \text{corr}(\mathbf{x}, \mathbf{y}) \cdot \frac{s_y}{s_x}$$

$$\text{and } \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

```
1 > (B1hat = cor(X,Y) * sd(Y)/sd(X))
2 [1] 1.150092
3 > (B0hat = mean(Y) - B1hat*mean(X))
4 [1] -130.9104
5 > SCR2 = function(B){sum((Y-B[1] -B
6   [2]*X)^2)}
7 > optim(c(0,0),SCR2)
8 $par
9 [1] -130.96193      1.15037
10 $value
11 [1] 14321.11
```



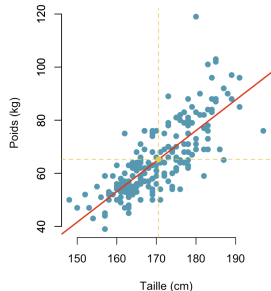
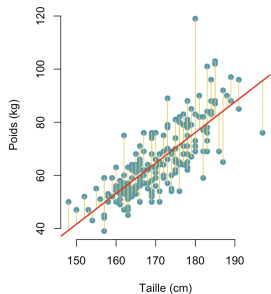
Poids & Taille

$$\widehat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i \text{ and } \widehat{\varepsilon}_i = y_i - \widehat{y}_i$$

```
1 > Y_hat = B0hat + B1hat*X
2 > E_hat = Y-Y_hat
3 > sum(E_hat)
4 [1] -2.046363e-12
5 > sum(E_hat*X)
6 [1] -3.517471e-10
7 > B0hat + B1hat*mean(X)
8 [1] 65.255
9 > mean(Y)
10 [1] 65.255
```

i.e. $\bar{y} = \hat{\beta}_0 + \hat{\beta}_1 \bar{x}$

$$\sum_{i=1}^n \widehat{\varepsilon}_i = 0 \text{ and } \sum_{i=1}^n x_i \widehat{\varepsilon}_i = 0$$



Poids & Taille

$$\underbrace{\sum_{i=1}^n (y_i - \bar{y})^2}_{\substack{\text{total sum of squares} \\ \text{TSS}}} = \underbrace{\sum_{i=1}^n (y_i - \hat{y}_i)^2}_{\substack{\text{residual sum of squares} \\ \text{RSS}}} + \underbrace{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}_{\substack{\text{explained sum of squares} \\ \text{ESS}}}$$

```
1 > (TSS=sum( (Y-mean(Y))^2 ))
2 [1] 35322
3 > RSS=sum( E_hat^2 )
4 > ESS=sum( (Y_hat-mean(Y))^2 )
5 > ESS+RSS
6 [1] 35322
```

$$R^2 = \frac{ESS}{TSS} = 1 - \frac{RSS}{TSS} = \widehat{\beta}_1^2 \frac{s_x}{s_y} = \frac{s_{xy}^2}{s_x s_y} = \text{corr}(\mathbf{x}, \mathbf{y})^2$$

```
1 > (R2 = 1-RSS/TSS)
2 [1] 0.5945555
```

Poids & Taille

$$\hat{\sigma}^2 = \frac{1}{n-2} \sum_{i=1}^n \hat{\varepsilon}_i^2 = \frac{RSS}{n-2}$$

```
1 > (s = sqrt(sum(E_hat^2)/(length(Y)-2)))  
2 [1] 8.504635
```

$$\widehat{\text{Var}}[\widehat{\beta}_0] = \widehat{\sigma}^2 \left(\frac{1}{n} + \frac{\bar{x}^2}{s_x^2} \right) \text{ and } \widehat{\text{Var}}[\widehat{\beta}_1] = \frac{\widehat{\sigma}^2}{s_x^2}$$

```
1 > (s1 = sqrt(s^2/sum((X-mean(X))^2)))  
2 [1] 0.06749465  
3 > (s0 = sqrt(s^2*(1/length(Y)+mean(X)^2/sum((X-mean(X))^2))))  
4 [1] 11.52792
```

Poids & Taille

To test $H_0 : \beta_1 = 0$ against $H_1 : \beta_1 \neq 0$

$$T = \frac{\hat{\beta}_1 - 0}{\sqrt{\widehat{\text{Var}}[\hat{\beta}_1]}} \sim \text{Std}(n-2) \text{ if } H_0 \text{ is true}$$

```
1 > (t0 = B0hat/s0)
2 [1] -11.35594
3 > (t1 = B1hat/s1)
4 [1] 17.03975
```

or

$$F = \frac{(TSS - RSS)/1}{RSS/(n-2)} = \frac{ESS}{RSS/(n-2)} \sim \mathcal{F}_{1,n-2} \text{ if } H_0 \text{ is true}$$

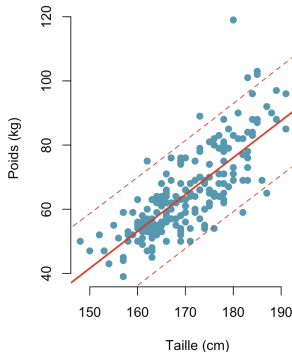
```
1 > (F = ESS*(length(Y)-2)/RSS)
2 [1] 290.353
```


Poids & Taille

For prediction, $\hat{y}_x = \hat{\beta}_0 + \hat{\beta}_1 x$ and

$$\sqrt{\widehat{\text{Var}}[y_x - \hat{y}_x]} = \hat{\sigma} \cdot \sqrt{1 + \frac{1}{n} + \frac{(x - \bar{x})^2}{s_x^2}}$$

```
1 > pred = function(x){  
2 +   yx = B0hat + B1hat*x  
3 +   sd_yx = s*sqrt(1+1/length(Y)+(x-  
4 +     mean(X))^2/sum((X-mean(X))^2))  
5 +   c(yx+qt(.025,length(Y)-2)*sd_yx,  
6 +     yx,  
7 +     yx+qt(.975,length(Y)-2)*sd_yx)  
8 + }  
9 > pred(170)  
      lower      pred      upper  
47.79186  64.60520  81.41853
```



Using R Functions

To fit a model, use

```
1 > model = lm(weight~height, data=Davis)
2 > summary(model)
3
4 Coefficients:
5             Estimate Std. Error t value Pr(>t)
6 (Intercept) -130.91040    11.52792   -11.36  <2e-16 ***
7 height       1.15009     0.06749    17.04  <2e-16 ***
8
9 Residual standard error: 8.505 on 198 DF
10 Multiple R-squared:  0.5946, Adjusted R-squared:  0.5925
11 F-statistic: 290.4 on 1 and 198 DF, p-value: < 2.2e-16
```

and for the prediction

```
1 > predict(model,newdata=data.frame(height=170),
2         interval = "prediction",level = .95)
3         fit      lwr      upr
4 1 64.6052 47.79186 81.41853
```