# Modèles Linéaires Appliqués

Arthur Charpentier

Automne 2020

OLS #9 (boostrap & incertitude)

# Régression Linéaire & Bootstrap (1)

Dataset $\{\mathbf{z}_i = (y_i, \mathbf{x}_i)\}$, $i = 1, \cdots, n$.
Use paired sampling by (repeatedly)
resampling $\{\mathbf{z}_1^\star, \cdots, \mathbf{z}_n^\star\}$.
**Idea** :
$\{(y_i, \mathbf{x}_i)\}$ is obtained from (unknown) $\mathbb{P}$
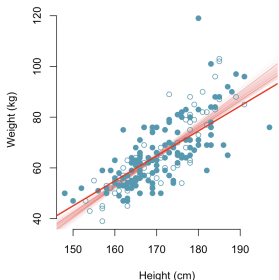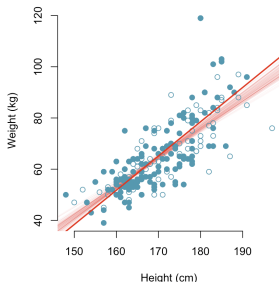Based on $n$ observations, we observe $\mathbb{P}_n$
We generate other samples by resampling
from $\mathbb{P}_n$

1. sample $\{i_1^{(b)}, \cdots, i_n^{(b)}\}$ randomly with replacement in $\{1, 2, \cdots, n\}$
2. consider dataset $(\mathbf{x}_i^{(b)}, y_i^{(b)}) = (\mathbf{x}_{i^{(b)}}, y_{i^{(b)}})$'s, and fit a model
3. let $\widehat{\boldsymbol{\beta}}^{(b)}$ denote the estimated values, or $\widehat{y}_{n+1}^{(b)}$ some prediction

# Régression Linéaire & Bootstrap (1)



```
1 > BETA = matrix(NA,1000,2)
2 > for(s in 1:1000){
3   idx = sample(1:nrow(Davis),nrow(
      Davis),replace=TRUE)
4   reg_sim = lm(weight~height, data=
      Davis[idx,])
5   BETA[s,] = reg_sim$coefficients
6 }
7 > hist(BETA[,1])
8 > hist(BETA[,2])
```
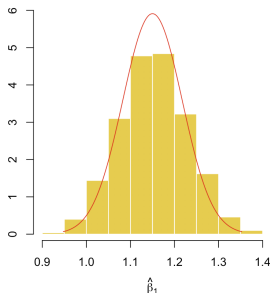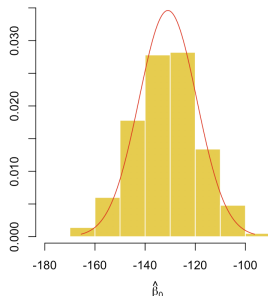
# Régression Linéaire & Bootstrap (1)



```r
1 > BETA = matrix(NA,1000,2)
2 > for(s in 1:1000){
3   idx = sample(1:nrow(Davis),nrow(
      Davis),replace=TRUE)
4   reg_sim = lm(weight~height, data=
      Davis[idx,])
5   BETA[s,] = reg_sim$coefficients
6 }
7 > hist(BETA[,1])
8 > hist(BETA[,2])
```

# Régression Linéaire & Bootstrap (2)

As an alternative model-based resampling

1. sample $\widehat{\varepsilon}_1^{(b)}, \cdots, \widehat{\varepsilon}_n^{(b)}$ resample from $\{\widehat{\varepsilon}_1, \widehat{\varepsilon}_2, \cdots, \widehat{\varepsilon}_n\}$
2. set $y_i^{(b)} = \widehat{\beta}_0 + \widehat{\beta}_1 x_i + \widehat{\varepsilon}_i^{(b)} = \widehat{y}_i + \widehat{\varepsilon}_i^{(b)}$
3. consider dataset $(\mathbf{x}, y^{(b)}) = (\mathbf{x}_i, y_i^{(b)})$'s and fit a model
4. let $\widehat{\boldsymbol{\beta}}^{(b)}$ denote estimated values
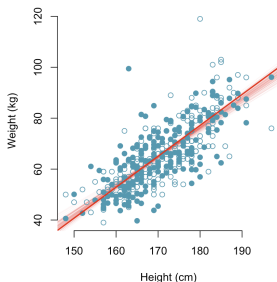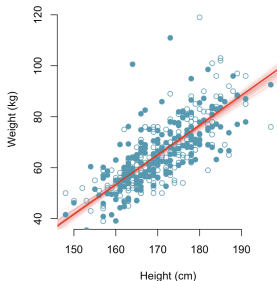
**Note** in a simple regression

$$\widehat{\beta}_1^{(b)} = \frac{\sum [x_i - \overline{x}] \cdot y_i^{(b)}}{\sum [x_i - \overline{x}]^2} = \widehat{\beta}_1 + \frac{\sum [x_i - \overline{x}] \cdot \widehat{\varepsilon}_i^{(b)}}{\sum [x_i - \overline{x}]^2}$$

hence $\mathbb{E}[\widehat{\beta}_1^{(b)}] = \widehat{\beta}_1$, while

$$\mathsf{Var}[\widehat{\beta}_1^{(b)}] = \frac{\sum [x_i - \overline{x}]^2 \cdot \mathsf{Var}[\widehat{\varepsilon}_i^{(b)}]}{\left( \sum [x_i - \overline{x}]^2 \right)^2} \sim \frac{\sigma^2}{\sum [x_i - \overline{x}]^2}$$

# Régression Linéaire & Bootstrap (2)

```
1 > BETA = matrix(NA,1000,2)
2 > reg = lm(weight~height, data=Davis)
3 > epsilon = residuals(reg)
4 > for(s in 1:1000){
5     eps = sample(epsilon,nrow(Davis),
        replace=TRUE)
6     Davis_s = data.frame(height =
        Davis$height, weight =predict(reg)
        +eps)
7     reg_sim = lm(weight~height, data=
        Davis_s)
8     BETA[s,] = reg_sim$coefficients
9 }
10 > hist(BETA[,1])
11 > hist(BETA[,2])
```
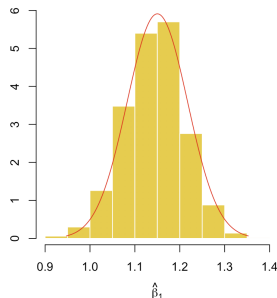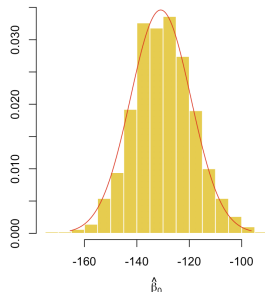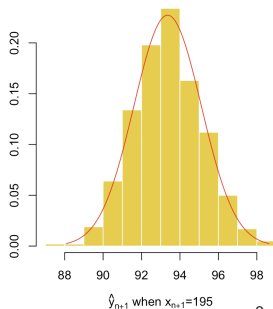
# Régression Linéaire & Bootstrap (2)

```
1  > BETA = matrix(NA,1000,2)
2  > reg = lm(weight~height, data=Davis)
3  > epsilon = residuals(reg)
4  > for(s in 1:1000){
5    eps = sample(epsilon,nrow(Davis),
       replace=TRUE)
6    Davis_s = data.frame(height =
       Davis$height, weight =predict(reg)
       +eps)
7    reg_sim = lm(weight~height, data=
       Davis_s)
8    BETA[s,] = reg_sim$coefficients
9  }
10 > hist(BETA[,1])
11 > hist(BETA[,2])
```
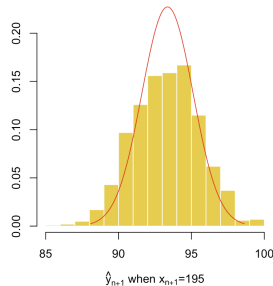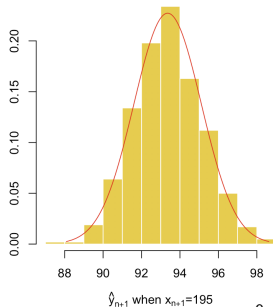
# Régression Linéaire & Bootstrap (1/2)

```r
> PRED = matrix(NA,1000,2)
> nwDavis = data.frame(height = 195)
> for(s in 1:1000){
+    idx = sample(1:n,n,replace=TRUE)
+    reg_sim = lm(weight~height, data=
     Davis[idx,])
+    PRED[s,1] = predict(reg_sim,
     newdata=nwDavis)
+    eps = sample(epsilon,nrow(Davis),
     replace=TRUE)
+    Davis_s = data.frame(height =
     Davis$height, weight =predict(reg)
     +eps)
+    reg_sim = lm(weight~height, data=
     Davis_s)
+    PRED[s,2] = predict(reg_sim,
     newdata=nwDavis)
+ }
```



$\hat{y}_{n+1}$ when $x_{n+1}=195$



$\hat{y}_{n+1}$ when $x_{n+1}=195$

# Régression Linéaire & Bootstrap (1/2)

```
1 > apply(PRED,2,function(x) quantile(x
      ,.025))
2 [1]  89.04203 89.63030
3 > apply(PRED,2,function(x) quantile(x
      ,.975))
4 [1]  97.60345 97.02423
5 > predict(lm(weight~height, data=Davis
      ), newdata=nwDavis,interval="
      confidence",se.fit = TRUE)
6 $fit
7        fit        lwr        upr
8 1 93.35749 89.89571 96.81927
9
10 $se.fit
11 [1] 1.755451
```



$\hat{y}_{n+1}$ when $x_{n+1}=195$



$\hat{y}_{n+1}$ when $x_{n+1}=195$

# Bootstrap heuristics

Here $\widehat{\boldsymbol{\beta}} = \boldsymbol{\beta} + (\mathbf{X}^\mathsf{T}\mathbf{X})^{-1}\mathbf{X}^\mathsf{T}\boldsymbol{\varepsilon} = T_\beta(\boldsymbol{\varepsilon})$ or $\widehat{y}_{n+1} = \mathbf{x}'_{n+1}\widehat{\boldsymbol{\beta}} = T_y(\boldsymbol{\varepsilon})$

Use simulations, we draw $n$ values $\{\epsilon_1, \cdots, \epsilon_n\}$ and

▶ $\mathbb{E}\left[\dfrac{1}{n}\displaystyle\sum_{i=1}^{n} T(\epsilon_i)\right] = \mathbb{E}[T(\boldsymbol{\varepsilon})]$ (unbiased)

▶ $\dfrac{1}{n}\displaystyle\sum_{i=1}^{n} T(\epsilon_i) \xrightarrow{\mathcal{L}} \mathbb{E}[T(\boldsymbol{\varepsilon})]$ as $n \to \infty$ (consistent)

# Bootstrap & Tests

Consider the test of $H_0 : \beta_j = 0$,

1. compute $t_n = \dfrac{(\widehat{\beta_j} - \beta_j)^2}{\widehat{\sigma}_j^2}$

2. generate $B$ boostrap samples, under the null assumption $H_0$

3. for each boostrap sample, compute $t_n^{(b)} = \dfrac{(\widehat{\beta}_j^{(b)} - \widehat{\beta}_j)^2}{\widehat{\sigma}_j^{2(b)}}$

4. reject $H_0$ if $\dfrac{1}{B} \displaystyle\sum_{i=1}^{B} \mathbf{1}(t_n > t_n^{(b)}) < \alpha$.

# Bootstrap & Tests

What does "generate $B$ boostrap samples, under the null assumption $H_0$" mean ?

**Example** : $y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$ with $H_0 : \beta_1 = 0$.

2.1. Estimate the model under $H_0$, i.e. $y_i = \beta_0 + \eta_i$, and save $\{\widehat{\eta_1}, \cdots, \widehat{\eta_n}\}$

2.2. Define $\widetilde{\boldsymbol{\eta}} = \{\widetilde{\eta}_1, \cdots, \widetilde{\eta}_n\}$ with $\widetilde{\eta} = \sqrt{\dfrac{n}{n-1}}\widehat{\eta}$

2.3. Draw (with replacement) residuals $\widetilde{\boldsymbol{\eta}}^{(b)} = \{\widetilde{\eta}_1^{(b)}, \cdots, \widetilde{\eta}_n^{(b)}\}$

2.4. Set $y_i^{(b)} = \widehat{\beta_0} + \widetilde{\eta}_i^{(b)}$

2.5. Estimate the regression model $y_i^{(b)} = \beta_0^{(b)} + \beta_1^{(b)} x_i + \varepsilon_i^{(b)}$

3. for each boostrap sample, compute $t_n^{(b)} = \dfrac{(\widehat{\beta_j}^{(b)} - \widehat{\beta_j})^2}{\widehat{\sigma}_j^{2(b)}}$

4. reject $H_0$ if $\dfrac{1}{B}\sum_{i=1}^{B}\mathbf{1}(t_n > t_n^{(b)}) < \alpha$.