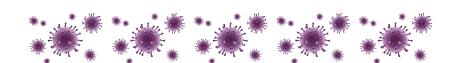
# Modèles Linéaires Appliqués / Régression Sélection de Variables & Stepwise

Arthur Charpentier

**UQAM** 

Hiver 2020 - COVID-19 # 20





## Stepwise - Procédure Pas à Pas

### Retour sur la régression linéaire

```
> data(swiss)
2 > str(swiss)
 'data frame': 47 obs. of 6 variables:
 $ Fertility
                           80.2 83.1 92.5 85.8 76.9
                    : num
  $ Agriculture
                    : num
                           17 45.1 39.7 36.5 43.5
 $ Examination
                    : int
                           15 6 5 12 17 9 16 14 12
 $ Education
                    : int
                           12 9 5 7 15 7 7 8 7 13
                                                     . . .
 $ Catholic
                    : nim
                           9.96 84.84 93.4 33.77 5
                                                     . . .
  $ Infant
                    : num
                           22.2 22.2 20.2 20.3 20.6
                                                     . . .
```

$$y_i = \beta_0 + \beta_1 x_{1,i} + \dots + \beta_{k,i} x_{k,i} + \varepsilon_i = \mathbf{x}_i^{\top} \boldsymbol{\beta} + \varepsilon_i$$



### Test de Student

Idée classique : partir du modèle complet, enlever *la* variable la moins significative, et itérer

```
> summary(lm(Fertility~Examination+Agriculture+InfantM
     +Catholic+Education, data=swiss))
2
 Coefficients:
             Estimate Std. Error t value Pr(>t)
4
  (Intercept) 66.91518 10.70604 6.250 1.91e-07 ***
6 Examination -0.25801 0.25388 -1.016 0.31546
7 Agriculture -0.17211
                        0.07030 - 2.448 0.01873 *
8 InfantM 1.07705
                        0.38172 2.822 0.00734 **
Catholic 0.10412
                        0.03526 2.953 0.00519 **
10 Education -0.87094
                        0.18303 - 4.758 2.43e - 05 ***
```

(jusqu'à ce que toutes les variables soient significatives)

## Stepwise - Procédure Pas à Pas

1.  $y_i = \beta_0 + \beta_1 x_{1,i} + \cdots + \beta_{k,i} x_{k,i} + \varepsilon_i$  et enlever une variable

```
1 > drop1(lm(Fertility~Examination+Agriculture+InfantM+
     Catholic+Education, data=swiss), test="F")
2 Single term deletions
3
4 Model:
5 Fertility ~ Examination + Agriculture + InfantM +
     Catholic + Education
           Df Sum of Sq RSS AIC F value Pr(>F)
6
 <none>
                      2105.0 190.69
8 Examinati 1 53.03 2158.1 189.86 1.0328 0.3154
9 Agricultu 1 307.72 2412.8 195.10 5.9934 0.0187 *
10 Infant.M 1 408.75 2513.8 197.03 7.9612 0.0073 **
11 Catholic 1 447.71 2552.8 197.75 8.7200 0.0051 **
12 Education 1 1162.56 3267.6 209.36 22.6432 2.4e-05 ***
```

## Stepwise - Procédure Pas à Pas

2. Partir de  $y_i = \beta_0 + \varepsilon_i$  et rajouter une variable...

```
1 > add1(lm(Fertility~1, data=swiss), Fertility~
     Agriculture+Infant.Mortality+Catholic+Examination+
     Education , test="F")
2 Single term additions
4 Model:
5 Fertility ~ 1
           Df Sum of Sq RSS AIC F value Pr(>F)
6
 <none>
                      7178.0 238.34
8 Agricultu 1 894.8 6283.1 234.09 6.4089 0.0149 *
9 Infant.M 1 1245.5 5932.4 231.39 9.4477 0.0035 **
10 Catholic 1 1543.3 5634.7 228.97 12.3251 0.0010 **
11 Examinati 1 2994.4 4183.6 214.97 32.2087 9.4e-07 ***
            1 3162.7 4015.2 213.04 35.4456 3.6e-07 ***
12 Education
```

### AIC & BIC

```
Akaike, AIC = -2 \log \mathcal{L}(\widehat{\beta}) + 2k
Bayésien/Schwarz, BIC = -2 \log \mathcal{L}(\widehat{\beta}) + k \log(n)
```

```
1 > n=nrow(swiss)
2 > reg = lm(Fertility ~ Agriculture + Catholic +
      Education, data=swiss)
3 > s2 = sum(resid(reg)^2)
4 > n + n*log(2*pi) + n*log(s2/n) + 2*(1+4)
5 [1] 331.4126
6 > -2*logLik(reg) + 2*(1+4)
7 'log Lik.' 331.4126 (df=5)
8 > AIC(reg, k=2)
9 [1] 331.4126
10 > n + n*log(2*pi) + n*log(s2/n) + log(n)*(1+4)
11 [1] 340.6634
12 > -2*logLik(reg)+log(n)*(1+4)
13 'log Lik.' 340.6634 (df=5)
14 > AIC(reg, k=log(n))
15 [1] 340.6634
```

```
Akaike, AIC = -2 \log \mathcal{L}(\widehat{\beta}) + 2k
```

```
> step(lm(Fertility~1, data=swiss), Fertility~
     Agriculture+Infant.Mortality+Catholic+Examination+
     Education,direction = "forward")
2 Start: AIC=238.35
3 Fertility ~ 1
4
               Df Sum of Sq RSS
                                     AIC
5
6 + Education 1 3162.7 4015.2 213.04
7 + Examination 1 2994.4 4183.6 214.97
8 + Catholic 1 1543.3 5634.7 228.97
9 + Infant.M 1 1245.5 5932.4 231.39
10 + Agriculture 1 894.8 6283.1 234.09
11 <none>
                           7178.0 238.34
```

On enlève une variable + Education, et on continue

```
1 Step: AIC=213.04
2 Fertility ~ Education
3
             Df Sum of Sq RSS AIC
4
5 + Catholic 1 961.07 3054.2 202.18
6 + Infant.M 1 891.25 3124.0 203.25
7 + Examination 1 465.63 3549.6 209.25
                           4015.2 213.04
8 <none>
9 + Agriculture 1 61.97 3953.3 214.31
10
11 Step: AIC=202.18
12 Fertility ~ Education + Catholic
13
               Df Sum of Sq RSS
                                    AIC
14
+ Infant.M 1 631.92 2422.2 193.29
16 + Agriculture 1 486.28 2567.9 196.03
17 <none>
                           3054.2 202.18
18 + Examination 1 2.46 3051.7 204.15
```

```
1 Step: AIC=193.29
2 Fertility ~ Education + Catholic + InfantM
3
               Df Sum of Sq RSS AIC
4
5 + Agriculture 1 264.176 2158.1 189.86
6 <none>
                        2422.2 193.29
7 + Examination 1 9.486 2412.8 195.10
9 Step: AIC=189.86
10 Fertility ~ Education + Catholic + InfantM + Agricult
11
               Df Sum of Sq RSS AIC
12
                            2158.1 189.86
13 <none>
14 + Examination 1 53.027 2105.0 190.69
```

#### et le modèle final est

```
Akaike, AIC = -2 \log \mathcal{L}(\widehat{\beta}) + 2k
```

```
> step(lm(Fertility~Agriculture+InfantM+Catholic+
     Examination+Education, data=swiss), direction = "
     backward")
2 Start: ATC=190.69
3 Fertility ~ Agriculture + InfantM + Catholic +
     Examination +
   Education
4
5
               Df Sum of Sq RSS AIC
6
  - Examination 1
                       53.03 2158.1 189.86
                            2105.0 190.69
8 <none>
9 - Agriculture 1 307.72 2412.8 195.10
10 - Infant.M 1 408.75 2513.8 197.03
11 - Catholic 1 447.71 2552.8 197.75
12 - Education 1 1162.56 3267.6 209.36
```

On enlève une variable - Examination, et on continue

### et le modèle final est (encore)

```
lm(formula = Fertility ~ Education + Catholic + Infant
.Mortality + Agriculture, data = swiss)
```

### AIC & BIC

```
Akaike, AIC = -2 \log \mathcal{L}(\widehat{\beta}) + 2k
Bayésien/Schwarz, BIC = -2 \log \mathcal{L}(\widehat{\beta}) + k \log(n)
... qui peuvent être aussi utilisés pour les GLM
```

```
1 > data("birthwt", package = "MASS")
2 > str(birthwt)
3 'data.frame': 189 obs. of 10 variables:
  $ low : int 0 0 0 0 0 0 0 0 0 ...
 $ age : int 19 33 20 21 18 21 22 17 29 26 ...
6 $ lwt : int 182 155 105 108 107 124 118 103 123 ...
7 $ race : int 2 3 1 1 1 3 1 3 1 1 ...
8 $ smoke: int 0 0 1 1 1 0 0 0 1 1 ...
9 $ ptl : int 0 0 0 0 0 0 0 0 0 ...
10 $ ht : int
              0 0 0 0 0 0 0 0 0 0 ...
11 $ ui : int 1 0 0 1 1 0 0 0 0 ...
12 $ ftv : int 0 3 1 2 0 0 1 1 1 0 ...
  $ bwt : int 2523 2551 2557 2594 2600 2622 2637 ...
13
```

### AIC & BIC

```
1 > bwt <- with(birthwt, {</pre>
     race <- factor(race, labels = c("white", "black
     ". "other"))
3 + ptd \leftarrow factor(ptl > 0)
4 + ftv <- factor(ftv)</pre>
5 + levels(ftv)[-(1:2)] <- "2+"
6 + data.frame(low = factor(low), age, lwt, race,
     smoke = (smoke > 0), ptd, ht = (ht > 0), ui = (ui
     > 0), ftv) })
7 > str(bwt)
8 'data.frame': 189 obs. of 9 variables:
9 $ low : Factor w/ 2 levels "0", "1": 1 1 1 1 1 1 ...
10 $ age : int 19 33 20 21 18 21 22 17 29 26 ...
11 $ lwt : int 182 155 105 108 107 124 118 103 123 ...
12 $ race : Factor w/ 3 levels "white", "black", ...: ...
$ smoke: logi FALSE FALSE TRUE TRUE TRUE FALSE ...
$ ptd : Factor w/ 2 levels "FALSE", "TRUE": 1 1 ...
$ ht : logi FALSE FALSE FALSE FALSE FALSE ...
```

#### On considère les deux modèles extrêmes

```
1 > fullmod = glm(low ~ ., data=bwt,family=binomial)
2 > zeromod = glm(low ~ 1, data=bwt,family=binomial)
```

### Puis on peut avancer, pas à pas, à partir du modèle complet

```
1 > backwards = step(fullmod, direction = "backward")
2 Start: AIC=217.48
3 low ~ age + lwt + race + smoke + ptd + ht + ui + ftv
4
   Df Deviance AIC
5
6 - ftv 2 196.83 214.83
7 - age 1 196.42 216.42
8 <none> 195.48 217.48
9 - ui 1 197.59 217.59
10 - smoke 1 198.67 218.67
11 - race 2 201.23 219.23
12 - lwt 1 200.95 220.95
13 - ht 1 202.93 222.93
14 - ptd 1 203.58 223.58
```

```
1 Step: AIC=214.83
2 low ~ age + lwt + race + smoke + ptd + ht + ui
3
       Df Deviance AIC
5 - age 1 197.85 213.85
6 <none> 196.83 214.83
7 - ui 1 199.15 215.15
8 - race 2 203.24 217.24
9 - smoke 1 201.25 217.25
10 - lwt 1 201.83 217.83
11 - ptd 1 203.95 219.95
12 - ht 1 204.01 220.01
13
14 Step: AIC=213.85
15 low ~ lwt + race + smoke + ptd + ht + ui
16
17 Df Deviance AIC
18 <none > 197.85 213.85
```

On peut aussi avancer, pas à pas, à partir du modèle simple

```
1 > forwards = step(zeromod, scope=list(lower=formula(
     zeromod), upper=formula(fullmod)), direction="
     forward")
2 Start: AIC=236.67
3 low ~ 1
4
         Df Deviance
                      ATC
5
6 + ptd 1 221.90 225.90
7 + lwt 1 228.69 232.69
8 + ui 1 229.60 233.60
9 + smoke
          1 229.81 233.81
          1 230.65 234.65
10 + ht
          2 229.66 235.66
11 + race
          1 231.91 235.91
12 + age
13 <none>
            234.67 236.67
14 + ftv
          2 232.09 238.09
```

```
1 Step: AIC=225.9
2 low ~ ptd
3
        Df Deviance AIC
4
5 + age 1 217.30 223.30
6 + lwt 1 217.50 223.50
7 + ht 1 217.66 223.66
8 + race 2 217.02 225.02
9 + ui 1 219.12 225.12
10 + smoke 1 219.33 225.33
          2 217.88 225.88
11 + ftv
12 <none> 221.90 225.90
13
14 Step: AIC=223.3
15 low ~ ptd + age
16
        Df Deviance AIC
17
18 + ht 1 213.12 221.12
          1 214.31 222.31
19 + lwt
```

```
1 Step: AIC=221.12
2 low ~ ptd + age + ht
3
       Df Deviance AIC
4
5 + lwt 1 207.43 217.43
6 + ui 1 210.13 220.13
7 + smoke 1 210.89 220.89
8 <none> 213.12 221.12
9 + race 2 210.06 222.06
10 + ftv 2 210.38 222.38
11
12 Step: AIC=217.43
13 low ~ ptd + age + ht + lwt
14
Df Deviance AIC
16 + ui 1 205.15 217.15
17 + smoke 1 205.39 217.39
18 <none> 207.43 217.43
```

#### On a ici deux modèles différents

### On peut aussi utiliser une méthode mixte

```
1 > bothways = step(zeromod, list(lower=formula(zeromod)
     ,upper=formula(fullmod)), direction="both")
2 Start: AIC=236.67
3 low ~ 1
4
        Df Deviance
                        AIC
5
6 + ptd 1 221.90 225.90
7 + lwt 1 228.69 232.69
8 + ui 1 229.60 233.60
9 + smoke 1 229.81 233.81
          1 230.65 234.65
10 + ht
          2 229.66 235.66
11 + race
          1 231.91 235.91
12 + age
            234.67 236.67
13 <none>
          2 232.09 238.09
14 + ftv
```

```
(...)
1 Step: AIC=217.15
2 low ~ ptd + age + ht + lwt + ui
3
        Df Deviance AIC
5 <none> 205.15 217.15
6 + smoke 1 203.24 217.24
7 + race 2 201.25 217.25
8 - ui 1 207.43 217.43
9 - age 1 207.51 217.51
10 + ftv 2 202.41 218.41
11 - lwt 1 210.13 220.13
12 - ht 1 212.70 222.70
13 - ptd 1 215.48 225.48
```

#### On retient ici le modèle

### AIC & All Models

On peut aussi utiliser regsubsets,

```
1 > library(leaps)
2 > allreg = regsubsets(low ~ ., data = bwt, nbest = 1,
     nvmax = NULL, force.in = NULL, force.out = NULL,
     method = "exhaustive")
 > as.data.frame(summary(allreg)$outmat)
           age lwt rak rar smo ptd htT uiT ft1 f2+
4
5 1 (1)
62 (1)
9 5 (1)
11 7 (1)
12 8 (1)
13 9 (1)
14 10 (1)
```

### AIC & All Models

que l'on peut aussi visualiser dans un graphique

```
1 > plot(allreg, scale = "bic")
```

