

Modèles Linéaires Appliqués / Régression

Régression Logistique : x continue

Arthur Charpentier

UQAM

Hiver 2020 - COVID-19 # 5

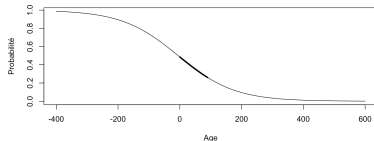
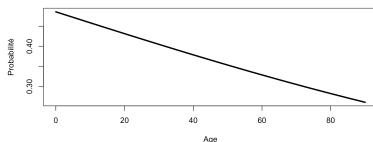


Régression sur x_1

```
1 > loc = "http://freakonometrics.free.fr/titanic.RData"
2 > download.file(loc, "titanic.RData")
3 > load("titanic.RData")
4 > base = base[!is.na(base$Age),1:7]
5 > reg = glm(Survived~Age,family=binomial,data=base)
6 > summary(reg)
```

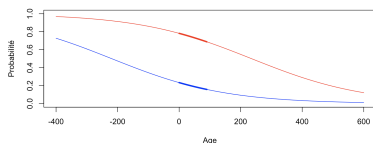
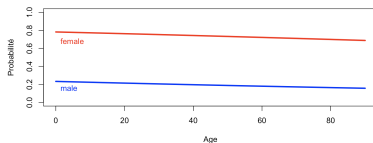
Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-0.05672	0.17358	-0.327	0.7438
Age	-0.01096	0.00533	-2.057	0.0397 *



Régression sur $x_1 + \mathbf{1}_2$

```
1 > reg = glm(Survived~Age+Sex,family=binomial)
2
3 Coefficients:
4             Estimate Std. Error z value Pr(>z)
5 (Intercept)  1.277273   0.230169   5.549 2.87e-08 ***
6 Age         -0.005426   0.006310  -0.860    0.39
7 Sexmale     -2.465920   0.185384 -13.302 < 2e-16 ***
```



Régression sur $x_1 + \mathbf{1}_2$

$$\text{Prédiction : } \hat{p}_{x,s} = \frac{e^{\hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 \mathbf{1}_M}}{1 + e^{\hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 \mathbf{1}_M}} \text{ où } s \in \{F, M\}.$$

Pas de relation simple entre $\hat{p}_{x,M}$ et $\hat{p}_{x,F}$

$$\text{ou } \text{cote}_{x,s} = \frac{\hat{p}_{x,s}}{1 - \hat{p}_{x,s}} = e^{\hat{\beta}_0 + \hat{\beta}_1 x + \hat{\beta}_2 \mathbf{1}_M}, \text{ aussi,}$$

$$\text{cote}_{x,M} = \text{cote}_{x,F} \cdot e^{\hat{\beta}_2}$$

Régression sur x_1

Comme dans les 'modèles linéaires', on peut chercher une transformation non-linéaire de x_1

- régression polynomiale

```
1 > reg3=glm(Survived~poly(Age,3), family=binomial)
2 > summary(reg3)
3
4 Coefficients:
5             Estimate Std. Error z value Pr(>z)
6 (Intercept)  -0.39327    0.07739  -5.082 3.74e-07 ***
7 poly(Age, 3)1 -4.81433    2.15139  -2.238  0.02524 *
8 poly(Age, 3)2  2.40500    2.29333   1.049  0.29432
9 poly(Age, 3)3 -6.83820    2.42311  -2.822  0.00477 **
```

$\text{logit}(p) = P(x_1)$ avec P polynôme de degré 3 (significatif).
avec `poly(,3)` écrit dans une base de polynômes orthogonaux



Régression sur x_1

écriture dans la base usuelle, $\text{logit}(p) = \beta_0 + \beta_1 x_1 + \beta_2 x_1^2 + \beta_3 x_1^3$

```
1 > reg3b=glm(Survived~Age+I(Age^2)+I(Age^3), family=
    binomial, data=base[!is.na(base$Age),])
2 > summary(reg3b)
3
4 Coefficients:
5             Estimate Std. Error z value Pr(>z)
6 (Intercept)  9.258e-01  3.612e-01   2.563  0.010364 *
7 Age         -1.428e-01  4.078e-02  -3.501  0.000463 ***
8 I(Age^2)      4.392e-03  1.418e-03   3.096  0.001960 **
9 I(Age^3)     -4.069e-05  1.442e-05  -2.822  0.00477 **
```

Remarque La p -value du test $H_0 : \beta_3 = 0$ est identique...

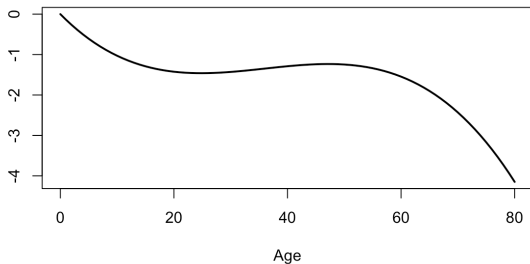
Les deux modèles sont **équivalents**



Régression sur x_1

On peut visualiser $x_1 \mapsto \hat{\beta}_1 x_1 + \hat{\beta}_2 x_1^2 + \hat{\beta}_3 x_1^3$

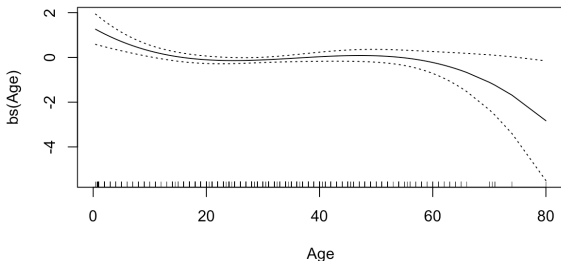
```
1 > beta = coefficients(reg3b)
2 > xage = 0:80
3 > yage = beta[2]*xage+beta[3]*xage^2+beta[4]*xage^3
4 > plot(xage,yage)
```



Régression sur x_1

- régression spline

```
1 > library(splines)
2 > res = glm(Survived~bs(Age),family=binomial)
3 > library(gam)
4 > regam=gam(Survived~bs(Age),family=binomial)
5 > plot(reg, se=TRUE)
```



Cf cas linéaire : splines cubiques $\cdots + (x_1 - s_1)_+^3 + \cdots + (x_1 - s_2)_+^3$

