

# Modèles Linéaires Appliqués / Régression

## Modèle de Poisson : Exemple

Arthur Charpentier

UQAM

Hiver 2020 - COVID-19 # 11



## Fréquence d'Accident

```
1 > loc = "http://freakonometrics.free.fr/SAuto.RData"
2 > download.file(loc_fichier, "SinistresAuto.RData")
3 > load("SinistresAuto.RData")
4 > str(base)
5 'data.frame': 50000 obs. of 11 variables:
6 $ nocontrat      : int  27 115 121 142 155 186 217 ...
7 $ exposition     : num  0.87 0.72 0.05 0.9 0.12 ...
8 $ zone           : Factor w/ 6 levels "A","B","C", ...
9 $ puissance      : int  7 5 6 10 7 5 5 5 4 4 ...
10 $ agevehicule    : int  0 0 0 10 0 0 4 0 0 0 ...
11 $ ageconducteur  : int  56 45 37 42 59 75 31 41 42 ...
12 $ bonus          : int  50 50 55 50 50 50 64 90 50 ...
13 $ marque         : int  12 12 12 12 12 12 3 12 12 ...
14 $ carburant       : Factor w/ 2 levels "D","E": ...
15 $ region         : int  93 54 11 93 73 42 21 11 ...
16 $ nb             : num  0 0 0 0 0 0 0 0 0 0 ...
```

On veut modéliser la fréquence annuelle de sinistre  
à partir de nb survenus pendant la durée exposition

# Fréquence d'Accident : 1

Modèle :  $Y_i \sim \mathcal{P}(\lambda e_i)$  où  $\lambda_i = \exp(\beta_0)$

```
1 > reg = glm(nb~1+offset(log(exposition)),
2           data=base,family=poisson)
3 > summary(reg)
4
5 Coefficients:
6             Estimate Std. Error z value Pr(>|z|)
7 (Intercept)  -2.6201      0.0228  -114.9   <2e-16 ***
8
9 (Dispersion parameter taken to be 1)
10
11 Null deviance: 12680  on 49999  degrees of freedom
12 Residual deviance: 12680  on 49999  degrees of freedom
13 AIC: 16353
14
15 Number of Fisher Scoring iterations: 6
```

## Fréquence d'Accident : 1

Modèle :  $Y_i \sim \mathcal{P}(\lambda e_i)$  où  $\lambda_i = \exp(\beta_0)$

```
1 > with(base, sum(nb)/sum(exposition))
2 [1] 0.07279295
3 > coefficients(reg)
4 (Intercept)
5 -2.620136
6 > exp(coefficients(reg))
7 (Intercept)
8 0.07279295
9 > predict(reg, newdata=data.frame(exposition=1))
10 1
11 -2.620136
12 > predict(reg, type="response", newdata=data.frame(
13     exposition=1))
14 1
15 0.07279295
```

Aussi,  $\hat{\eta}_i = \hat{\beta}_0 = -2.6201$  et  $\hat{\mu}_i = 0.07276$ .

## Fréquence d'Accident : $x_1$

Modèle :  $Y_i \sim \mathcal{P}(\lambda_i e_i)$  où  $\lambda_i = \exp(\beta_0 + \beta_1 \mathbf{1}(x_i) = E)$ ,  
 $x_i$  est le carburant (diesel, essence)

```
1 > with(base, table(carburant, nb))
2     nb
3 carburant    0     1     2     3     4     6
4         D 23475   886    45     6     1     0
5         E 24722   821    41     1     1     1
6 > reg1 = glm(nb~carburant+offset(log(exposition)),
7 +           data=base, family=poisson)
8 > summary(reg1)
9
10 Coefficients:
11             Estimate Std. Error z value Pr(>|z|)
12 (Intercept) -2.52929    0.03165 -79.904  < 2e-16 ***
13 carburantE   -0.18031    0.04563  -3.952 7.76e-05 ***
```

Aussi  $\hat{\beta}_0 = -2.5293$  et  $\hat{\beta}_1 = -0.1803$ .

## Fréquence d'Accident : $x_1$

$\hat{\beta}_0 = -2.5293$  et  $\hat{\beta}_1 = -0.1803$ .

$$\hat{\lambda}_i = \begin{cases} e^{-2.5293} = 0.079715 & \text{si } x_i = D \\ e^{-2.5293-0.1803} = 0.066563 & \text{si } x_i = E \end{cases}$$

```
1 > predict(reg1,newdata=data.frame(carburant=c("D","E")  
2     ,exposition=1),type="response")  
3     1           2  
0.07971533 0.06656323
```

## Fréquence d'Accident : $x_2$

Modèle :  $Y_i \sim \mathcal{P}(\lambda_i e_i)$  où  $\lambda_i = \exp(\beta_0 + \beta_j \mathbf{1}(x_i = j))$ ,  
 $x_i$  est la zone  $\{A, B, C, D, E, F\}$

```
1 > with(base, table(zone, nb))
2     nb
3 zone    0     1     2     3     4     6
4   A  7411  231     9     0     1     0
5   B  5442  170     5     0     0     0
6   C 13619  462    22     2     0     0
7   D 10749  385    25     2     0     1
8   E  9704  401    23     3     1     0
9   F  1272   58     2     0     0     0
```

La zone A est la zone de référence ici.

## Fréquence d'Accident : $x_2$

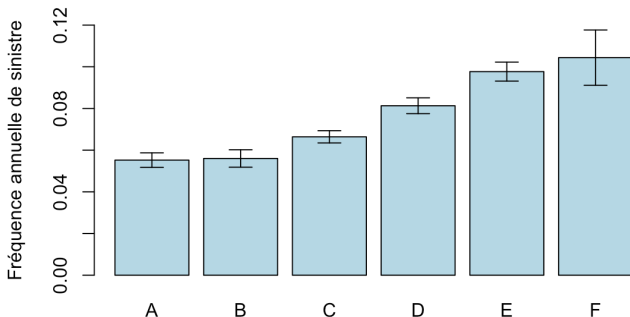
Modèle :  $Y_i \sim \mathcal{P}(\lambda_i e_i)$  où  $\lambda_i = \exp(\beta_0 + \beta_j \mathbf{1}(x_i = j))$ ,  
 $x_i$  est la zone  $\{A, B, C, D, E, F\}$

```
1 > reg2 = glm(nb~zone+offset(log(exposition)),
2 +           data=base,family=poisson)
3 > summary(reg2)
4
5 Coefficients:
6             Estimate Std. Error z value Pr(>|z|)
7 (Intercept) -2.89626    0.06287  -46.069   < 2e-16 ***
8 zoneB        0.01396    0.09751   0.143    0.8862
9 zoneC        0.18407    0.07685   2.395    0.0166 *
10 zoneD        0.38662    0.07836   4.934 8.07e-07 ***
11 zoneE        0.57025    0.07827   7.286 3.20e-13 ***
12 zoneF        0.63659    0.14171   4.492 7.05e-06 ***
```



## Fréquence d'Accident : $x_2$

```
1 y=predict(reg2,newdata=data.frame(zone=LETTERS[1:6],  
    exposition=1),type="response", se.fit =TRUE)  
2 error.bar = function(x,y,upper, lower,length=0.1,...){  
3     arrows(x,y+upper, x, y-lower, angle=90, code=3,  
        length=length, ...)}  
4 barres = barplot(y$fit, beside=TRUE, col="light blue")  
5 error.bar(barres,y$fit,y$se.fit)
```



## Fréquence d'Accident : $x_2$

```
1 > M = matrix(NA,6,6)
2 > for(i in 1:6){
3   base$zone = relevel(base$zone,LETTERS[i])
4   reg2i = glm(nb~zone+offset(log(exposition)),
5               data=base,family=poisson)
6   p=summary(reg2i)$coefficients[2:6,4]
7   names(p)=substr(names(p),5,5)
8   M[i,] = p[LETTERS[1:6]]
9 }
10 > round(M,3)
11      A      B      C      D      E      F
12 A      NA 0.886 0.017 0.000 0.000 0.000
13 B 0.886      NA 0.050 0.000 0.000 0.000
14 C 0.017 0.050      NA 0.002 0.000 0.001
15 D 0.000 0.000 0.002      NA 0.005 0.065
16 E 0.000 0.000 0.000 0.005      NA 0.624
17 F 0.000 0.000 0.001 0.065 0.624      NA
```

## Fréquence d'Accident : $x_2$

```
1 > base$zone = factor(base$zone, levels = LETTERS[1:6])
2 > levels(base$zone) = c("AB","AB","C","D","EF","EF")
3 > reg2b = glm(nb~zone+offset(log(exposition)),data=
4     base,family=poisson)
5
6 Coefficients:
7             Estimate Std. Error z value Pr(>|z|)
8 (Intercept) -2.89048      0.04806 -60.148  < 2e-16 ***
9 zoneC        0.17829      0.06529   2.731  0.00632 **
10 zoneD        0.38084      0.06706   5.679 1.36e-08 ***
11 zoneEF       0.57212      0.06500   8.802  < 2e-16 ***
```

## Fréquence d'Accident : $x_1 + x_2$

On considère ici la régression sur deux variables catégorielles,

```
1 > (E=xtabs(exposition~carburant+zone,data=base))
2           zone
3 carburant    AB      C      D      EF
4           D 4307.494 3745.807 2471.297 1994.951
5           E 3487.379 3966.076 3150.068 3308.060
6 > (N=xtabs(nb~carburant+zone,data=base))
7           zone
8 carburant    AB      C      D      EF
9           D  261  295  210  232
10          E  172  217  247  290
11 > reg12 = glm(nb~carburant+zone+offset(log(exposition)
12             ),data=base,family=poisson)
13 > yp = predict(reg12,type="response")
14 > xtabs(yp~base$carburant+base$zone)
15           base$zone
16 base$carburant    AB      C      D      EF
17           D 264.4026 279.1828 228.0513 226.3634
18           E 168.5974 232.8172 228.9487 295.6366
```

## Fréquence d'Accident : $x_1 + x_2$

Méthode des marges,  $N_{i,j} \sim \mathcal{P}(E_{i,j} \cdot \lambda_{i,j})$  où  $\lambda_{i,j} = A_i \cdot B_j$  où

$$\sum_i E_{i,j} A_i B_j = \sum_i N_{i,j} \text{ et } \sum_j E_{i,j} A_i B_j = \sum_j N_{i,j}$$

soit

$$A_i = \frac{\sum_j N_{i,j}}{\sum_j E_{i,j} B_j} \text{ et } B_j = \frac{\sum_i N_{i,j}}{\sum_i E_{i,j} A_i}$$

```
1 A=rep(1,length(levels(base$carburant)))
2 B=rep(1,length(levels(base$zone)))*sum(N)/sum(E)
3 for(i in 1:1000){
4   A=apply(N,1,sum)/apply(t(B*t(E)),1,sum)
5   B=apply(N,2,sum)/apply(A*E,2,sum)
6 }
```

## Fréquence d'Accident : $x_1 + x_2$

```
1 > B
2           AB           C           D           EF
3 0.05447450 0.06614477 0.08189545 0.10069925
4 > E * A%%t(B)
5           zone
6 carburant           AB           C           D           EF
7           D 264.4026 279.1828 228.0513 226.3634
8           E 168.5974 232.8172 228.9487 295.6366
```

$$\sum_i E_{i,j} A_i B_j = \sum_i N_{i,j} \text{ et } \sum_j E_{i,j} A_i B_j = \sum_j N_{i,j}$$

```
1 > apply(N,2,sum)
2 AB    C    D    EF
3 433 512 457 522
4 > apply(E * A%%t(B),2,sum)
5 AB    C    D    EF
6 433 512 457 522
```

## Fréquence d'Accident : $x_3$

Modèle :  $Y_i \sim \mathcal{P}(\lambda_i e_i)$  où  $\lambda_i = \exp(\beta_0 + \beta_1 x_i)$ ,  $x_i$  est l'âge

```
1 > reg = glm(nb~ageconducateur+offset(log(exposition)),
2 +          data=base,family=poisson)
3 > summary(reg)
4
5 Coefficients:
6             Estimate Std Error z value Pr(>|z|)
7 (Intercept)  -2.163061   0.077380  -27.95  < 2e-16 ***
8 ageconducateur -0.009891   0.001635   -6.05  1.45e-09 ***
9
10 (Dispersion parameter taken to be 1)
11
12     Null deviance: 12680   on 49999   degrees of freedom
13 Residual deviance: 12642   on 49998   degrees of freedom
14 AIC: 16317
15
16 Number of Fisher Scoring iterations: 6
```

Aussi  $\hat{\beta}_0 = -2.1631$  et  $\hat{\beta}_1 = -0.0099$ .

## Fréquence d'Accident : $x_3$

Modèle :  $Y_i \sim \mathcal{P}(\lambda_i e_i)$  où  $\lambda_i = \exp(\beta_0 + \beta_1 x_i)$ ,  $x_i$  est l'âge.

Ici  $\hat{\beta}_0 = -2.1631$  et  $\hat{\beta}_1 = -0.0099$ .

$\hat{\eta}_i = -2.1631 - 0.0099x_i$  et  $\hat{\lambda}_i = e^{\hat{\eta}_i} = e^{-2.1631 - 0.0099x_i}$

On peut aussi utiliser des splines,  $\log \lambda_i = \beta_0 + s(x_i)$

```
1 > library(splines)
2 > regs = glm(nb~bs(ageconducuteur)+
3   offset(log(exposition)), data=base, family=poisson)
```

$$s(x) = \beta_1 x + \beta_2 x^2 + \beta_3 x^3 + \beta_4 (x - s_1)_+^3 + \beta_5 (x - s_2)_+^3$$

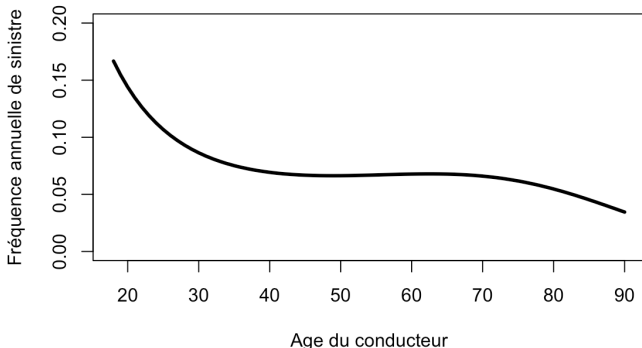
(splines cubiques)



## Fréquence d'Accident : $x_3$

Modèle :  $Y_i \sim \mathcal{P}(\lambda_i e_i)$  où  $\lambda_i = \exp(\beta_0 + \beta_1 x_i)$ ,  $x_i$  est l'âge

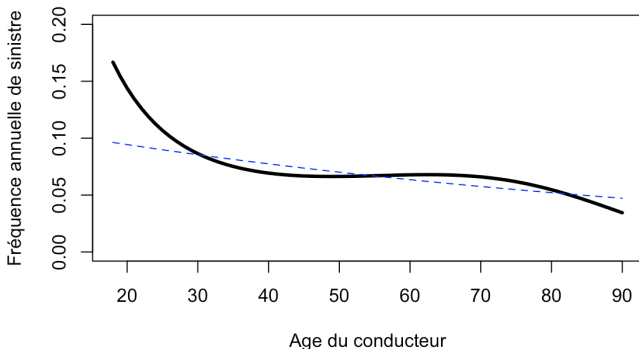
```
1 a = seq(18,90)
2 y = predict(regs,newdata=data.frame(ageconducteur=a,
    exposition=1),type="response", se.fit = TRUE)
3 plot(a,y$fit,type="l")
```



## Fréquence d'Accident

Modèle :  $Y_i \sim \mathcal{P}(\lambda_i e_i)$  où  $\lambda_i = \exp(\beta_0 + \beta_1 x_i)$ ,  $x_i$  est l'âge

```
1 y0=predict(reg,newdata=data.frame(ageconducteur=a,  
    exposition=1),type="response")  
2 lines(a,y0,lty=2)
```



Comparaison avec le modèle linéaire

## Fréquence d'Accident : $x_3$

Modèle :  $Y_i \sim \mathcal{P}(\lambda_i e_i)$  où  $\lambda_i = \exp(\beta_0 + \beta_1 x_i)$ ,  $x_i$  est l'âge

```
1 points(a, y$fit, pch=19)
2 segments(a, y$fit-2*y$se.fit, a, y$fit+2*y$se.fit)
```

