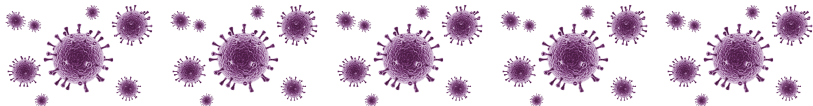


Modèles Linéaires Appliqués / Régression GLM & Poids et Modèles Tweedie

Arthur Charpentier

UQAM

Hiver 2020 - COVID-19 # 16



Les Poids en Régression

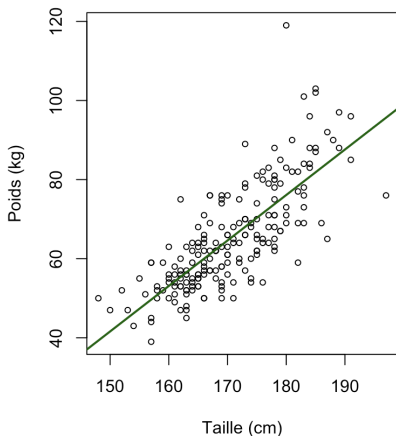
Dans la régression classique (Gaussienne) on cherche à résoudre

$$\hat{\beta} = \operatorname{argmin} \left\{ \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2 \right\}$$

Avec des **poids** $\omega = (\omega_i)$

$$\hat{\beta}_\omega = \operatorname{argmin} \left\{ \sum_{i=1}^n \omega_i (y_i - \mathbf{x}_i^\top \beta)^2 \right\}$$

Problème des données agrégées,
poids = taille de la population
cf exemple taille et poids des élèves
agrégés ici par 'province' ou 'région'



Les Poids en Régression

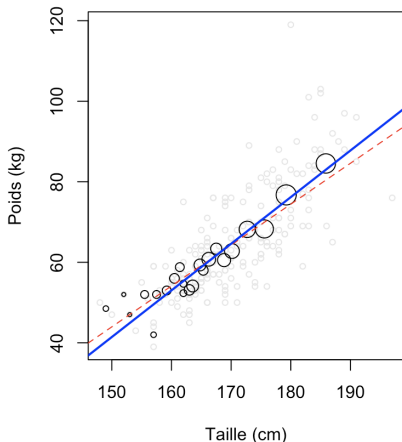
Dans la régression classique (Gaussienne) on cherche à résoudre

$$\hat{\beta} = \operatorname{argmin} \left\{ \sum_{i=1}^n (y_i - \mathbf{x}_i^\top \beta)^2 \right\}$$

Avec des **poids** $\omega = (\omega_i)$

$$\hat{\beta}_\omega = \operatorname{argmin} \left\{ \sum_{i=1}^n \omega_i (y_i - \mathbf{x}_i^\top \beta)^2 \right\}$$

Problème des données agrégées,
poids = taille de la population
cf exemple taille et poids des élèves
agrégés ici par 'province' ou 'région'



Les Poids en Régression

$$\hat{\beta}_{\omega} = \operatorname{argmin} \left\{ \sum_{i=1}^n \omega_i (y_i - \mathbf{x}_i^{\top} \beta)^2 \right\}$$

donne la condition du premier ordre, en dérivant par rapport à β_j

$$\sum_{i=1}^n \omega_i x_{j,i} (y_i - \mathbf{x}_i^{\top} \beta) = 0, \quad \forall j, \quad \text{ou} \quad \mathbf{X}^{\top} \Omega (\mathbf{y} - \mathbf{X} \beta) = \mathbf{0}$$

où $\Omega = \operatorname{diag}(\omega_i)$, soit $\hat{\beta} = (\mathbf{X}^{\top} \Omega \mathbf{X})^{-1} \mathbf{X}^{\top} \Omega \mathbf{y}$

Dans les GLM, la log-vraisemblance est alors

$$\log \mathcal{L}(\theta, \varphi | \mathbf{y}) = \sum_{i=1}^n \omega_i \frac{y_i \theta_i - b(\theta_i)}{\varphi} + \text{terme indépendant de } \theta_i$$

Les Poids en Régression

- **Régression de Poisson** : poids vs. exposition (offset)

Exposition e_i : $\log \mathcal{L}(\beta; \mathbf{y}) = \sum_{i=1}^n [y_i \log(e_i \lambda_i) - e_i \lambda_i - \log(y_i!)]$

```
1 > reg = glm(Y~YEARMARRIAGE+offset(log(EXPOSITION)),  
  data=DF, family=poisson)
```

Poids y_i/e_i : $\log \mathcal{L}^*(\beta; \mathbf{y}) = \sum_{i=1}^n e_i \left[\frac{y_i}{e_i} \log(\lambda_i) - \lambda_i - \log(y_i!) \right]$

```
1 > reg = glm(Y/EXPOSITION~YEARMARRIAGE, data=DF,  
  weights=EXPOSITION, family=poisson)
```

Les deux modèles sont semblables...

Les Poids en Régression

- Régression Binomiale : poids vs. taille - $Y \sim \mathcal{B}(n, p)$

```
1 > reg = glm(cbind(Y,n-Y)~X1+X2, data=DF, binomial)
```

Déboulement $(y_i, x_{1,i}, x_{2,i}, n_i) \rightarrow n_i$ lignes
 $(1, x_{1,i}, x_{2,i}) \times y_i$ et $(0, x_{1,i}, x_{2,i}) \times n_i - y_i$

```
1 > reg = glm(y~X1+X2, data=bigDF, family=binomial)
```

Poids modéliser y_i/n_i avec un poids n_i

```
1 > reg = glm(y/n~X1+X2, data=DF, binomial, weights=n)
```

Les trois modèles sont semblables...

Les Poids en Régression

- Tableau de contigence, y vs x

```
1 > df = read.table("http://freakonometrics.free.fr/
    baseexo.csv", sep=";", header=TRUE)
2 > (M=as.matrix(df[,2:ncol(df)]))
3 M
4      14 19 21 23 25 27 29 31 33 35
5 18 74 13  7  1  0  0  0  0  0
6 23  6  4  2  7  4  0  0  0  0
7 25  2  3  2  2  4  0  0  0  0
8 27  1  1  2  3  3  2  0  0  0
9 29  2  0  1  3  2  0  6  0  0
10 31  2  0  2  1  0  0  1  2  1
11 33  0  0  0  2  0  0  1  1  3
12 35  0  1  0  1  0  0  0  0  2
13 37  0  0  0  0  1  1  0  1  0  1
```

Les Poids en Régression

Poids (x, y) et le poids n



```
1 > W = as.vector(M)
2 > X = rep(df[,1], ncol(M))
3 > y = as.numeric(substr(names(df)[-1], 2, 3))
4 > Y = rep(y, each=nrow(M))
5 > cbind(X, Y, W)
6      X  Y  W
7 1  16 14 74
8 2  23 14  6
9 3  25 14  2
10 4  27 14  1
11 5  29 14  2
12 > summary(lm(Y ~ X, weights = W))
13
14 Coefficients:
15             Estimate Std. Error t value Pr(>|t|)
16 (Intercept)  4.35569    2.03022   2.145    0.038 *
17 X            0.68263    0.09016   7.572 3.04e-09 ***
```


Les Poids en Régression

Duplication (x, y) répété n fois



```
1 > vX = vY = rep(NA, sum(W))
2 > sumW = c(0, cumsum(W))
3 > for(i in 1:length(W)){
4   if(W[i]>0){
5     vX[(1+sumW[i]):sumW[i+1]] = X[i]
6     vY[(1+sumW[i]):sumW[i+1]] = Y[i] } }
7 > base2 = data.frame(x=vX,y=vY)
8 > tail(base2,10)
9     x  y
10 179 35 33
11 180 35 33
12 181 37 35
13 > summary(lm(y ~ x))
14
15 Coefficients:
16             Estimate Std. Error t value Pr(>|t|)
17 (Intercept)  4.35569    0.95972   4.538 1.04e-05 ***
18 X2           0.68263    0.04262  16.017 < 2e-16 ***
```

Modèle Tweedie

Consider a Tweedie distribution, with variance function power $p \in (1, 2)$, mean μ and scale parameter ϕ , then it is a compound Poisson model,

- $N \sim \mathcal{P}(\lambda)$ with $\lambda = \frac{\phi\mu^{2-p}}{2-p}$
- $Y_i \sim \mathcal{G}(\alpha, \beta)$ with $\alpha = -\frac{p-2}{p-1}$ and $\beta = \frac{\phi\mu^{1-p}}{p-1}$

Conversely, consider a compound Poisson model $N \sim \mathcal{P}(\lambda)$ and $Y_i \sim \mathcal{G}(\alpha, \beta)$, then

- variance function power is $p = \frac{\alpha+2}{\alpha+1}$
- mean is $\mu = \frac{\lambda\alpha}{\beta}$
- scale parameter is $\phi = \frac{[\lambda\alpha]^{\frac{\alpha+2}{\alpha+1}-1}\beta^{2-\frac{\alpha+2}{\alpha+1}}}{\alpha+1}$

Modèle Tweedie

In the context of regression

$$N_i \sim \mathcal{P}(\lambda_i) \text{ with } \lambda_i = \exp[\mathbf{x}_i^\top \boldsymbol{\beta}_\lambda]$$

$$Y_{j,i} \sim \mathcal{G}(\mu_i, \phi) \text{ with } \mu_i = \exp[\mathbf{x}_i^\top \boldsymbol{\beta}_\mu]$$

Then $S_i = Y_{1,i} + \dots + Y_{N,i}$ has a Tweedie distribution

- variance function power is $p = \frac{\phi + 2}{\phi + 1}$
- mean is $\lambda_i \mu_i$
- scale parameter is $\frac{\lambda_i^{\frac{1}{\phi+1}-1}}{\mu_i^{\frac{\phi}{\phi+1}}} \left(\frac{\phi}{1 + \phi} \right)$

There are $1 + 2\dim(\mathbf{X})$ degrees of freedom.

Remark Note that the scale parameter should not depend on i .

Modèle Tweedie

A Tweedie regression is

- variance function power is $p \in (1, 2)$
- mean is $\mu_i = \exp[\mathbf{x}_i^T \boldsymbol{\beta}_{\text{Tweedie}}]$
- scale parameter is ϕ

There are $2 + \dim(\mathbf{X})$ degrees of freedom.

