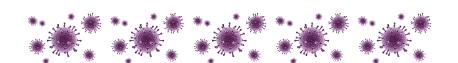
Modèles Linéaires Appliqués / Régression Régression de Poisson : Inférence

Arthur Charpentier

UQAM

Hiver 2020 - COVID-19 # 8



Données homogènes $\{y_1, \dots, y_n\}$ de loi de Poisson, $Y_i \sim \mathcal{P}(\lambda)$

$$\mathbb{P}(Y = y) = e^{-\lambda} \frac{\lambda^{y}}{y!}, \forall y \in \mathbb{N}.$$

Données hétérogènes $\{(y_1, x_1), \dots, (y_n, x_n)\}\$ de loi de Poisson

$$Y_i \sim \mathcal{P}(\lambda_i)$$
 avec $\lambda_i = \exp[\mathbf{x}_i^{\top} \boldsymbol{\beta}].$

Dans ce modèle, $\mathbb{E}(Y_i|\mathbf{x}_i) = \text{Var}(Y_i|\mathbf{x}_i) = \lambda_i = \exp[\mathbf{x}_i^{\top}\boldsymbol{\beta}].$

Modèle dit log-linéaire.

Remarque: on posera parfois $\theta_i = \eta_i = \mathbf{x}_i^{\top} \boldsymbol{\beta}$ et $\mu_i = \lambda_i = \exp[\eta_i]$.



La log-vraisemblance est ici

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

i.e.

$$\log \mathcal{L}(\beta; \mathbf{y}) = \sum_{i=1}^{n} y_i \cdot [\mathbf{x}_i^{\top} \beta] - \exp[\mathbf{x}_i^{\top} \beta] - \log(y_i!)$$

Le gradient est ici

$$\nabla \log \mathcal{L}(\beta; \mathbf{y}) = \frac{\partial \log \mathcal{L}(\beta; \mathbf{y})}{\partial \beta} = \sum_{i=1}^{n} \mathbf{x}_{i}^{\top} (y_{i} - \exp[\mathbf{x}_{i}^{\top} \beta])$$

alors que la matrice Hessienne s'écrit

$$H(\beta) = \frac{\partial^2 \log \mathcal{L}(\beta; \mathbf{y})}{\partial \beta \partial \beta^{\top}} = -\sum_{i=1}^n \exp[\mathbf{x}_i^{\top} \beta] \mathbf{x}_i \mathbf{x}_i^{\top}$$



La condition du premier ordre est ici

$$m{X}^{ op}[m{y} - \exp(m{X}\widehat{m{eta}})] = m{0}, \text{ i.e. } \sum_{i=1}^n m{x}_i^{ op}[m{y} - \exp(m{x}_i^{ op}\widehat{m{eta}})] = m{0}$$

ce qui signifie $\mathbf{X}^{\top}\mathbf{y} = \mathbf{X}^{\top}\widehat{\boldsymbol{\lambda}}$ où $\widehat{\boldsymbol{\lambda}} = \exp(\mathbf{X}\widehat{\boldsymbol{\beta}})$.

La matrice Hessienne est $H(\beta) = -\boldsymbol{X}^{\top} \boldsymbol{\Omega} \boldsymbol{X}, \ \boldsymbol{\Omega} = -\mathrm{diag}(\exp(\boldsymbol{X} \widehat{\beta}))$

Cf slides #2, en posant $\Omega = \operatorname{diag}(\boldsymbol{p}(1-\boldsymbol{p}))$,

$$\nabla \log \mathcal{L}(\beta) = \frac{\partial \log \mathcal{L}(\beta)}{\partial \beta} = \mathbf{X}^{\top}(\mathbf{y} - \mathbf{p}) = \mathbf{0} \quad (\text{avec } \mathbf{p} = \mathbf{p}(\beta))$$

$$H(\beta) = \frac{\partial^2 \log \mathcal{L}(\beta)}{\partial \beta \partial \beta^{\top}} = -\mathbf{X}^{\top} \mathbf{\Omega} \mathbf{X} \quad (\text{avec } \mathbf{\Omega} = \mathbf{\Omega}(\beta))$$



Parfois, les données ne sont pas observées pendant la même durée

Données hétérogènes $\{(y_1, x_1, e_1), \cdots, (y_n, x_n, e_n)\}$ de loi de Poisson (cf #7 processus de Poisson)

$$Y_i \sim \mathcal{P}(e_i \lambda_i) \text{ avec } \tilde{\lambda}_i = e_i \lambda_i = \exp[\mathbf{x}_i^{\top} \boldsymbol{\beta} + \log(e_i)].$$

 $log(e_i)$ est appelée variable offset.

Parfois les données sont regroupées (cf #9 régression binomiale $\mathcal{B}(n,p)$) et on suppose les groupes homogènes

$$Y_i = \sum_{i=1}^{e_i} Y_{i,j}$$
 avec $Y_{i,j} \sim \mathcal{P}(\lambda_i)$, i.i.d., i.e. $Y_i \sim \mathcal{P}(e_i \lambda_i)$.



Modèles Log-Linéaires

Ici $Y_i \sim \mathcal{P}(\exp[\mathbf{x}_i^{\top} \boldsymbol{\beta}])$.

On pourra considérer (#12) $Y_i \sim \mathcal{N}(\exp[\mathbf{x}_i^{\top} \boldsymbol{\beta}], \sigma^2)$, mais

$$\underbrace{Y_i \sim \mathcal{N}(\exp[\boldsymbol{x}_i^\top \boldsymbol{\beta}], \sigma^2)}_{\text{modèle GLM sur } y_i} \quad \neq \quad \underbrace{\log Y_i \sim \mathcal{N}(\boldsymbol{x}_i^\top \boldsymbol{\beta}, \sigma^2)}_{\text{modèle LM sur } \log y_i}$$

Dans le premier cas,
$$\mathbb{E}(Y_i) = \exp[\mathbf{x}_i^{\top} \boldsymbol{\beta}]$$

Dans le second cas,
$$\mathbb{E}(Y_i) = \exp\left[\mathbf{x}_i^{ op} \boldsymbol{\beta} + \frac{\sigma^2}{2} \right]$$

Cf loi lognormale, si
$$Y \sim \mathcal{LN}(\mu, \sigma^2)$$
, $\mathbb{E}(Y) = \exp\left[\mu + \frac{\sigma^2}{2}\right]$

Cf inégalité de Jensen, $\mathbb{E}(\log X) \leq \log \mathbb{E}(\log X)$ car log concave.

Modèles Log-Linéaires

Cas de données $N_{i,j}$ en 'triangle'...

```
1 > library(boot)
2 > data(aids)
3 > head(aids)
   year quarter delay dud time y
5 1 1983
6 2 1983
7 3 1983
```

lci year:quarter = date (i) et delay = développement (j)

```
[33,]
     53
           175
                35
                        13
                    17
                            11
 [34,] 63 135 24 23 12 1
                                NΑ
 [35,] 71 161 48 25 5
                            ΝA
                                NA
36 [36,] 95 178 39 6 NA
                            ΝA
                                NA
 [37,] 76 181 16 NA
                        NΑ
                            NΑ
                                NA
 [38,]
     67
            66
                ΝA
                    ΝA
                        ΝA
                            ΝA
                                NA
```

 \rightarrow utiliser un facteur ligne $\mathbf{L} = (L_1, \cdots, L_l)$ et un facteur colonne $\boldsymbol{C} = (C_1, \cdots, C_J)$

Modèles Log-Linéaires

```
Ici N_{i,j} \sim \mathcal{P}(L_i \cdot C_j) ou N_{i,j} \sim \mathcal{P}(\exp(\ell_i + \gamma_j))
```

On peut ensuite prévoir, $\widehat{\lambda}_{i,j} = \exp(\widehat{\ell}_i + \widehat{\gamma}_j)$

```
4 > p = predict(reg,newdata=Aids,type="response")
5 > matrix(p,length(unique(Aids$time)),byrow = TRUE)
6 [33,] 53.8 150.9 43.3 26.5 16.0 13.5 10.1 7.7
7 [34,] 47.6 133.5 38.3 23.4 14.2 11.9 8.9 6.8
8 [35,] 59.8 167.7 48.1 29.4 17.8 15.0 11.2 8.6
9 [36,] 67.7 189.9 54.5 33.3 20.2 16.9 12.7 9.7
10 [37,] 67.5 189.5 54.4 33.2 20.1 16.9 12.7 9.7
11 [38,] 67.0 188.0 53.9 33.0 19.9 16.8 12.6 9.6
```