

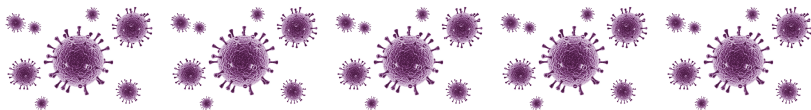
Modèles Linéaires Appliqués / Régression

Modèle de Poisson : Méthode des Marges

Arthur Charpentier

UQAM

Hiver 2020 - COVID-19 # 10



Méthode des Marges

Supposons que l'on prenne en compte ici deux classes de risques.

Tableau de contingence et biais minimal, Bailey (1963) et

Mildenhall (1999)

On suppose que

$$n_{i,j} \approx e_{i,j} L_i C_j, \text{ i.e. } \mathbf{n} \approx \mathbf{e} \cdot \mathbf{L} \mathbf{C}^\top$$

Une idée naturelle pour estimer $\mathbf{L} = (L_i)$ et $\mathbf{C} = (C_j)$ est d'utiliser une régression de Poisson,

$$N_{i,j} \sim \mathcal{P}(e_{i,j} L_i C_j) \text{ ou } N_{i,j} \sim \mathcal{P}(e_{i,j} \exp[\ell_i + \gamma_j])$$

```
1 > loc = "http://freakonometrics.free.fr/baseaffairs.  
    txt"  
2 > base = read.table(loc, header=TRUE)
```

Méthode des Marges

```
1 > (E=xtabs((Y>=0)~as.factor(RELIGIOUS)+as.factor(
    SATISFACTION),data=base))
2           as.factor(SATISFACTION)
3 as.factor(RELIGIOUS)  1  2  3  4  5
4                     1  1  1  4 17 19
5                     2  5 19 24 48 56
6                     3  2 12 20 38 46
7                     4  3 14 30 63 73
8                     5  1  5 11 21 30
9 > (N=xtabs(Y~as.factor(RELIGIOUS)+as.factor(
    SATISFACTION),data=base))
10          as.factor(SATISFACTION)
11 as.factor(RELIGIOUS)  1  2  3  4  5
12                     1  2  3 18  6 12
13                     2  6 17 14 53 13
14                     3  3 23 16 55 11
15                     4  2 17 15 27 12
16                     5  0  4  6 11  9
```

Méthode des Marges

Régression sur données individuelles :

```
1 > reg12 = glm(Y~as.factor(RELIGIOUS)+as.factor(
    SATISFACTION),data=base,family=poisson)
2 > yp = predict(reg12,type="response")
3 > xtabs(yp~as.factor(RELIGIOUS)+as.factor(SATISFACTION
    ),data=base)
```

		as.factor(SATISFACTION)				
(RELIGIOUS)		1	2	3	4	5
1	1	1.800	2.145	5.396	23.423	8.238
2	2	5.371	24.325	19.327	39.482	14.495
3	3	3.022	21.610	22.654	43.965	16.748
4	4	2.028	11.278	15.200	32.605	11.889
5	5	0.779	4.642	6.423	12.525	5.631

$$N_{i,j} = \sum_{k=1}^n \hat{\lambda}_k \mathbf{1}(x_{k,1} = i, x_{k,2} = j)$$

Méthode des Marges

Régression sur données individuelles :

```
1 > A=as.numeric(exp(coefficients(reg12)[1]+c(0,
    coefficients(reg12)[2:5])))
2 > B=as.numeric(exp(c(0,coefficients(reg12)[6:9])))
3 > as.numeric(A)
4 [1] 1.7995049 1.0742896 1.5111070 0.6759457 0.7789961
5 > as.numeric(B)
6 [1] 1.0000000 1.1917506 0.7495891 0.7656542 0.2409377
7 > E* (A%*%t(B))
8
9               as.factor(SATISFACTION)
10 (RELIGIOUS)      1      2      3      4      5
11      1  1.800  2.145  5.396 23.423  8.238
12      2  5.371 24.325 19.327 39.482 14.495
13      3  3.022 21.610 22.654 43.965 16.748
14      4  2.028 11.278 15.200 32.605 11.889
      5  0.779  4.642  6.423 12.525  5.631
```

Méthode des Marges

Régression sur données groupées : $N_{i,j} \sim \mathcal{P}(e_{i,j} L_i C_j)$,

```
1 > B=data.frame(N=as.numeric(N),E=as.numeric(E),  
    RELIGIOUS=rep(1:5,5),SATISFACTION=rep(1:5,each=5))  
2 > reg12g = glm(N~as.factor(RELIGIOUS)+as.factor(  
    SATISFACTION)+offset(log(E)),data=B,family=poisson  
    )  
3 > matrix(predict(reg12g,type="response"),5,5)  
4      [,1] [,2] [,3] [,4] [,5]  
5 [1,] 1.800 2.145 5.396 23.423 8.238  
6 [2,] 5.371 24.325 19.327 39.482 14.495  
7 [3,] 3.022 21.610 22.654 43.965 16.748  
8 [4,] 2.028 11.278 15.200 32.605 11.889  
9 [5,] 0.779 4.642 6.423 12.525 5.631
```

$$n_{i,j} = e_{i,j} \hat{L}_i \hat{C}_j, \text{ i.e. } \mathbf{n} = \mathbf{e} \cdot \mathbf{LC}^\top$$

Méthode des Marges

Oublions un instant la régression de Poisson,

Cherchons $\mathbf{L} = (L_i)$ et de $\mathbf{C} = (C_j)$ tels que $\mathbf{n} \approx \mathbf{e} \cdot \mathbf{L}\mathbf{C}^\top$

L'estimation de $\mathbf{L} = (L_i)$ et de $\mathbf{C} = (C_j)$ se fait généralement de trois manières:

- par moindres carrés
- par minimisation d'une distance (e.g. du chi-deux)
- par un principe de balancement (ou méthode des marges)

Méthode des Marges

Il est possible d'utiliser une méthode par moindres carrés (pondérée). On va chercher à minimiser la somme des carrés des erreurs, i.e.

$$D = \sum_{i,j} e_{i,j} (n_{i,j} - L_i C_j)^2$$

La condition du premier ordre donne ici

$$\frac{\partial D}{\partial L_i} = -2 \sum_j C_j e_{i,j} (n_{i,j} - L_i C_j) = 0$$

$$L_i = \frac{\sum_j C_j e_{i,j} n_{i,j}}{\sum_j e_{i,j} C_j^2} = \frac{\sum_j c_j y_{i,j}}{\sum_j e_{i,j} C_j^2}$$

L'autre condition du premier ordre donne

$$C_j = \frac{\sum_i L_i e_{i,j} n_{i,j}}{\sum_i e_{i,j} L_i^2} = \frac{\sum_i L_i y_{i,j}}{\sum_i e_{i,j} L_i^2}$$

On résoud alors ce petit système de manière itérative (car il n'y a pas de solution analytique simple).

Méthode des Marges

Il est aussi possible d'utiliser une méthode basée sur la distance du chi-deux. On va chercher à minimiser

$$Q = \sum_{i,j} \frac{e_{i,j}(n_{i,j} - L_i C_j)^2}{L_i C_j}$$

Là encore on utilise les conditions du premier ordre, et on obtient

$$L_i = \left(\frac{\sum_j \left(\frac{e_{i,j} y_{i,j}^2}{C_j} \right)}{\sum_j e_{i,j} C_j} \right)^{\frac{1}{2}} \quad \text{et} \quad L_j = \left(\frac{\sum_i \left(\frac{e_{i,j} y_{i,j}^2}{L_i} \right)}{\sum_i e_{i,j} L_i} \right)^{\frac{1}{2}}$$

où $y_{i,j} = e_{i,j} n_{i,j}$ (que l'on résout itérativement).

Méthode des Marges

Dans la méthode des marges – Bailey (1963) – formellement, on veut

$$\sum_j y_{i,j} = \sum_j e_{i,j} n_{i,j} = \sum_j e_{i,j} L_i C_j,$$

en somment sur la ligne i , pour tout i , ou sur la colonne j ,

$$\sum_i y_{i,j} = \sum_i e_{i,j} n_{i,j} = \sum_i e_{i,j} L_i C_j,$$

La première, et la seconde, équation donnent respectivement

$$L_i = \frac{\sum_j y_{i,j}}{\sum_j e_{i,j} C_j} \text{ et } C_j = \frac{\sum_i y_{i,j}}{\sum_i e_{i,j} L_i}.$$

Cette solution... correspond à la régression de Poisson.

Méthode des Marges

La première, et la seconde, équation donnent respectivement

$$L_i = \frac{\sum_j y_{i,j}}{\sum_j e_{i,j} C_j} \text{ et } C_j = \frac{\sum_i y_{i,j}}{\sum_i e_{i,j} L_i}.$$

```
1 > SATISFACTION=as.factor(base$SATISFACTION)
2 > RELIGIOUS=as.factor(base$RELIGIOUS)
3 > A=rep(1,length(levels(SATISFACTION)))
4 > B=rep(1,length(levels(RELIGIOUS)))*sum(N)/sum(E)
5 > for(i in 1:1000){
6     A=apply(N,1,sum)/apply(t(B*t(E)),1,sum)
7     B=apply(N,2,sum)/apply(A*E,2,sum) }
8 > E * A%*%t(B)
9
10      as.factor(SATISFACTION)
11 (RELIGIOUS)      1      2      3      4      5
12      1  1.800  2.145  5.396 23.423  8.238
13      2  5.371 24.325 19.327 39.482 14.495
14      3  3.022 21.610 22.654 43.965 16.748
15      4  2.028 11.278 15.200 32.605 11.889
16      5  0.779  4.642  6.423 12.525  5.631
```

Méthode des Marges

Cette solution vérifie

$$\sum_j y_{i,j} = \sum_j e_{i,j} L_i C_j, \text{ et } \sum_i y_{i,j} = \sum_i e_{i,j} L_i C_j,$$

```
1 > N0 = E * A%%t(B)
2 > apply(N0,1,sum)
3   1   2   3   4   5
4  41 103 108  73  30
5 > apply(N,1,sum)
6   1   2   3   4   5
7  41 103 108  73  30
8 > apply(N0,2,sum)
9   1   2   3   4   5
10 13  64  69 152  57
11 > apply(N,2,sum)
12   1   2   3   4   5
13 13  64  69 152  57
```