

Corrigé examen 14 décembre 2016

M2 "Modèles de régression"

Exercice 1

1) On modélise une variable qualitative fonction de variables numériques. On a donc un modèle logistique :  $Y = \text{rand}$ , 2 valeurs  $\begin{cases} R : \text{absence} \\ N : \text{présence} \end{cases}$   
 $\pi(x) = \mathbb{P}[Y = "R" | X = x]$  avec  $x = (km91, km118, \dots, km338)$   
 $= (x_1, x_2, \dots, x_7)$

Alors, le modèle logistique est :

(1)  $\underbrace{\mathbb{P}[Y = "R" | X = x]}_{\text{probabilité d'absence de nids}} = \frac{\exp(b_0 + \sum_{j=1}^7 b_j x_j)}{1 + \exp(b_0 + \sum_{j=1}^7 b_j x_j)}$

2) Le nombre de sites est 60. Le nombre de sites où il y a des nids de chouettes = 30.

3) Pour étudier la présence de nids de chouettes on regarde les sorties SAS :

$H_0$ : les 7 variables n'ont aucune influence sur l'absence des nids de chouettes

$$b_1 = b_2 = \dots = b_7 = 0$$

Modèle : (2)  $\pi(x) = \mathbb{P}[Y = "R" | X = x] = \frac{\exp(b_0)}{1 + \exp(b_0)}$

en fait,  $\pi(x)$  est constante, ne dépend pas de  $x$ .

$H_1$ :  $\exists j \in \{1, \dots, 7\}$  t. q.  $b_j \neq 0 \Leftrightarrow \pi(x)$  influence par au moins une des sept variables.

Modèle : (1)

(2)

Par les trois statistiques de test : rapport de vrais score, Wald, le hypothèse  $H_0$  est rejetée ( $p\text{-val} < 0.05$ ) et donc le modèle (1) significatif.

Indicateurs quantitatifs :

Pour (1) AIC = 68.134 BIC = 84.888

(2) AIC = 85.178 BIC = 87.272

les deux critères diminuent pour le modèle (1).

4) On étudie les sorties R. les sorties SAS donnent les mêmes résultats.

$H_0$  :  $X_1$  n'influe pas  $P[Y="R"]$  si  $X_2, \dots, X_7$  sont dans le modèle

$H_0$  :  $b_1 = 0$  |  $X_2, \dots, X_7$  dans le modèle

Modèle :  
(3)  $P[Y="R" | X=x] = \frac{\exp(b_0 + \sum_{j=2}^7 b_j x_j)}{1 + \exp(b_0 + \sum_{j=2}^7 b_j x_j)}$

$H_1$  :  $b_1 \neq 0$  |  $X_2, \dots, X_7$  dans le modèle.

Modèle : (1).

Stat de test :  $Z_1 = \frac{B_1}{\sqrt{\Delta B_1}} \xrightarrow{L} N(0, 1)$

La réalisation :  $z_1 = -1.535 \Rightarrow p\text{-value} = 0.12$

$\Rightarrow H_0$  acceptée.

Conclusion : les var. qui influent  $P[Y="R"]$

sont : km118  
km140

Rq :  $b_0 \neq 0$ . On garde aussi km91 et km241

5) Le Modèle (M2) est :

(3)

$$(4) P[Y = "R" | X = x] = \frac{\exp(b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_6 x_6)}{1 + \exp(b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + b_6 x_6)}$$

les estimations des paramètres :

$$\begin{aligned} \hat{b}_0 &= 10,2758 & \hat{b}_2 &= 0,1199 & \hat{b}_6 &= -0,08 \\ \hat{b}_1 &= -0,0643 & \hat{b}_3 &= -0,12 & & \end{aligned}$$

La var  $x_2 = \text{km}118$  a une influence positive sur la probabilité d'absence de nids pendant que  $\text{km}91$ ,  $\text{km}140$  et  $\text{km}241$  ont une influence négative.

Prévisions par (M2) :

- parmi les 30 sites avec nids, 24 sont bien prévus
- " sans " " " " "

Taux d'erreur:  $\frac{1}{2} \left( \frac{6}{30} + \frac{6}{30} \right) = \frac{1}{5} = 20\%$

Exercice 2

1) Les variables explicatives "site" et "faam" sont de type qualitatives, pendant que la variable expliquée "biomass" est numérique. Donc, on a un modèle d'analyse de variance à 2 facteurs

2) Le nb d'obs totales: 612  
 Nb d'obs utilisées: 450

Le modèle (M3) est :

$$Y_{ijk} = \mu + (\text{site})_i + (\text{faam})_j + (\text{site} \times \text{faam})_{ij} + \epsilon_{ijk}$$

avec:  $\epsilon_{ijk} \sim N(0, \sigma^2)$   $i = 1, \dots, 6$   $j = 1, 2, 3, 4$   
 $k = 1, \dots, n_{ij}$

$\epsilon_{ijk} \perp \epsilon_{i'j'k'}$   
 $(\text{site})_i$ : l'effet du site  $i$  sur la biomasse  
 $y$ : biomasse

(faam)<sub>j</sub> : l'effet de la famille j sur la biomasse  
 (site \* faam)<sub>ij</sub> : l'effet de l'interaction entre le site i et la famille j sur Y.

Contraintes :  $\sum_{i=1}^6 (\text{site})_i = 0$  ,  $\sum_{j=1}^4 (\text{faam})_j = 0$   
 $\forall i = 1, \dots, 6$  ,  $\sum_{j=1}^4 (\text{site} * \text{faam})_{ij} = 0$   
 $\forall j = 1, 2, 3, 4$  ,  $\sum_{i=1}^6 (\text{site} * \text{faam})_{ij} = 0$

Les valeurs des variables :

"site" : 1, 2, ..., 6  
 "faam" : 1, 2, 3, 4

3) H<sub>0</sub> : le modèle (M3) n'est pas significatif  
 (⇒) le site et la famille <sup>et leur interaction</sup> des pins n'ont aucune influence sur la biomasse  
 (⇒)  $\forall \text{site}_i = 0$  ,  $\forall \text{faam}_j = 0$  ,  $\forall (\text{site} * \text{faam})_{ij} = 0$   
 (⇒) Modèle :  $y_{ijk} = \mu + \epsilon_{ijk}$

H<sub>1</sub> : (M3) significatif  
 (⇒) Y influencé par le site ou par la famille ou par leur interaction.  
 (⇒)  $\exists i, \text{site}_i \neq 0$  ou  $\exists \text{faam}_j \neq 0$   
 ou  $\exists i, j (\text{site} * \text{faam})_{ij} \neq 0$

(⇒) Modèle (M3)

Statistique de test :

$$Z = \frac{SM/17}{SR/432} \underset{H_0}{\sim} F(17, 432)$$

Valeur stat de test :  $Z = 16,78$  et  $p\text{-value} \leq 10^{-4}$   
 ⇒ H<sub>0</sub> rejetée ⇒ (M3) significatif.

(5)

Tests des 2 facteurs et de leur interaction :

$H_0$ : le site n'influe pas la biomasse | la famille et l'interaction site \* famille sont dans le modèle

$\Leftrightarrow \forall i = 1, \dots, 6 \text{ site}_i = 0$

Modèle:  $y_{ijk} = \mu + (\text{faam})_j + (\text{site} * \text{faam})_{ij} + \epsilon_{ijk}$

$H_1$ : le site influe | fam et site \* faam dans le modèle

$\Leftrightarrow \exists i, \text{site}_i \neq 0$

Modèle: (M3)

Stat de test:  $z = \frac{S_{\text{site}} / 5}{SR / 432} \underset{H_0}{\sim} F(5, 432)$

Valeur stat de test  $z = 54,97$ , pvalue  $< 10^{-4}$   
 $\Rightarrow$  le site signif.

Parceil, on obtient que "faam" et "site \* faam" n'influent pas  $y$ , si on garde les autres var dans le modèle.

4) (M4):  $y_{ij} = \mu + \text{site}_i + \epsilon_{ij} \quad i = 1, \dots, 6$   
 $\sum_{i=1}^6 \text{site}_i = 0 \quad \epsilon_{ij} \sim N(0, \sigma^2) \quad j = 1, \dots, n_i$   
 $\epsilon_{ij} \perp \epsilon_{i'j'}$

5)  $H_0$ : pour (M4) le site n'influe pas  $y$   
 $\Leftrightarrow \text{site}_i = 0 \quad \forall i = 1, \dots, 6$

$\Leftrightarrow$  Modèle:  $y_{ij} = \mu + \epsilon_{ij}$

$H_1$ : le site influe  $y$  dans (M4)

$\Leftrightarrow$  Modèle (M4)

Stat de test:  $Z = \frac{SM/5}{SR/444} \stackrel{H_0}{\sim} F(5, 444)$

La valeur de la stat de test  $Z = 55.79$

p-value  $< 10^{-4} \Rightarrow$  le site influence  $Y \Rightarrow (M4)_{\text{sigif}}$

6)  $\hat{\sigma}^2 = 129807,23$

$\hat{\mu} = 773,83$

$\hat{site}_1 = 301,5$

$\hat{site}_5 = 301,57$

$\hat{site}_6 = -\sum_{j=1}^5 \hat{site}_j$

} sorties R.

Par rapport à la moyenne, un pin du site = 1 a une biomasse plus grande avec 301,5, le site = 2 a une biomasse inférieure de 479,7

7) On utilise les sorties SAS: on a que l'hypothèse  $site_4 = 0$  est acceptée. Mais la contrainte est  $site_6 = 0$ . Donc les sites 4 et 6 ont la même biomasse.

8) Pour (M3)  $R^2_{\text{adj}} = 0,374$

(M4):  $R^2_{\text{adj}} = 0,379$

donc on a les mêmes  $R^2_{\text{ajustés}}$

parce que (M4) est obtenu à partir de (M3) en supprimant les variables non significatives.

Donc <sup>part</sup> (M4) et (M3) on obtient la même qualité <sup>de prévision</sup>: assez bonne.

