

Corrigé examen 29 janvier 2020

Exercice 1

$Y$ : Conc et trois facteurs:  $F_1$ : LAB à 6 niveaux  
 $F_2$ : SPC à 7 niveaux  
 $F_3$ : BAT à 3 niveaux

1) On a un modèle d'analyse de variance à 3 facteurs avec une interaction:

$$(M1): Y_{ijkl} = \mu + (F_1)_i + (F_2)_j + (F_3)_k + (F_1 * F_3)_{ik} + \epsilon_{ijkl}$$

$n_{ijk} = 2$       $i = 1, \dots, 6$       $j = 1, \dots, 7$       $k = 1, 2, 3$       $l = 1, \dots, n_{ijk}$

$F_1_i$  c'est l'effet du niveau  $i$  de  $F_1$  sur  $Y$   
 $F_2_j$       $j$  de  $F_2$  sur  $Y$   
 $F_3_k$       $k$  de  $F_3$  sur  $Y$   
 $(F_1 * F_3)_{ik}$  du niveau  $i$  de  $F_1$  et de  $k$  de  $F_3$  sur  $Y$

Contraintes :

$$\sum_{i=1}^6 F_1_i = 0, \quad \sum_{j=1}^7 F_2_j = 0, \quad \sum_{k=1}^3 F_3_k = 0$$
$$\forall i = 1, \dots, 6, \quad \sum_{k=1}^3 (F_1 * F_3)_{ik} = 0$$
$$\forall k = 1, 2, 3, \quad \sum_{i=1}^6 (F_1 * F_3)_{ik} = 0$$

$\epsilon_{ijkl} \sim N(0, \sigma^2)$  indép.

2)  $H_0$ : (M1) non significatif  $\Leftrightarrow$  aucun des 3 facteurs et interaction n'influent  $Y$ .

$$\left\{ \begin{array}{l} F_1_i = 0, F_2_j = 0, F_3_k = 0, (F_1 * F_3)_{ik} = 0 \\ \forall i = 1, \dots, 6 \text{ et } \forall j = 1, \dots, 7 \text{ et } \forall k = 1, 2, 3 \end{array} \right.$$

Modèle réduit:  $Y_{ijkl} = \mu + \epsilon_{ijkl}$

(2)

$H_1: (M1) \text{ signif} \Leftrightarrow Y \text{ influencé par au moins un des facteurs ou par l'interaction entre } F1 \text{ et } F3$

$$\exists F1_i \neq 0 \text{ ou } \exists F2_j \neq 0 \text{ ou } \exists F3_k \neq 0 \text{ ou } \exists (F1 * F3)_{ik} \neq 0$$

Modèle complet: (M1)

Statistique de test:  $Z = \frac{SM/23}{SR/228} \underset{H_0}{\sim} F(23, 228)$

Valeur de la stat. de test  $Z = 541.4$  avec la p-value  $< 10^{-16} \Rightarrow H_0 \text{ rejetée} \Rightarrow (M1) \text{ significatif}$

3) On teste  $F1 * F3$ :

$H_0$ : l'interaction  $F1 * F3$  n'influe pas  $Y$  si  $F1, F2, F3$  sont dans le modèle.

Modèle réduit:  $\Leftrightarrow (F1 * F3)_{ik} = 0 \quad \forall i = 1, \dots, 6 \text{ et } \forall k = 1, 2, 3 \mid F1, F2, F3 \text{ dans le modèle}$   
 $Y_{ijk} = \mu + (F1)_i + (F2)_j + (F3)_k + \epsilon_{ijk}$

$H_1: (F1 * F3)$  influe  $Y$  si  $F1, F2, F3$  sont dans le modèle:

$\Leftrightarrow \exists (F1 * F3)_{ik} \neq 0$  si  $F1, F2, F3$  sont dans le modèle  
Modèle: (M1)

Stat de test:  $Z = \frac{S_{F1 * F3} / 10}{SR / 228} \underset{H_0}{\sim} F(10, 228)$

Valeur stat de test:  $Z = 2,6226$  avec p-value = 0,001  
 $\Rightarrow H_0 \text{ rejetée} \Rightarrow F1 * F3 \text{ influe } Y \text{ si } F1, F2, F3 \text{ sont dans le modèle}$

Facteurs significatifs  $F1$  et  $F2$ , pendant que  $F3$  n'est pas significatif. Donc il faut enlever le facteur BAT du modèle.

4)  $\hat{F}_{1,1} = -0,32$   $\hat{F}_{1,2} = 0,04$  ...  $\hat{F}_{1,5} = 0,01$ ,  $\hat{F}_{1,6} = -\sum_{i=1}^5 \hat{F}_{1,i}$   
 $\hat{F}_{2,1} = -1,41$ , ...  $\hat{F}_{2,6} = -0,13$ ,  $\hat{F}_{2,7} = -\sum_{i=1}^6 \hat{F}_{2,i}$   
 $\hat{F}_{3,1} = -0,03$   $\hat{F}_{3,2} = 0,05$   $\hat{F}_{3,3} = -(0,05 - 0,03) = -0,02$   
 $(F1 * F3)_{11} = -0,02$   $(F1 * F3)_{12} = 0,02$ ,  $(F1 * F3)_{13} = 0$   
 $\hat{\mu} = 1,92$

Par rapport à la moyenne, la concentration de y pour le laboratoire 1 est plus petite de 0.32

$\hat{F}_{3,2} = 0,05 \Rightarrow$  par rapport à la moyenne, l'échantillon 2 a une concentration plus grande de 0.05

$(F1 * F3)_{12} = 0,02$  : par rapport à la moyenne, l'interaction du laboratoire 1 et du lot 3 augmente la valeur de y de 0.02.

5)  $R^2_{adj} = 0,98$  donc proche de 1  $\Rightarrow$  (M1) de très bonne qualité d'ajustement.

Exercice 2

Y: class, var qualitative avec les valeurs 1 et 2  
 $V_1, V_2, V_3$  sont des var numériques

1) On a un modèle de régression logistique, qui a la forme :

(M2)  $P[Y_i = 2 | V_{1,i}, V_{2,i}, V_{3,i}] = \frac{\exp(\mu + b_1 V_{1,i} + b_2 V_{2,i} + b_3 V_{3,i})}{1 + \exp(\mu + \sum_{j=1}^3 b_j V_{j,i})}$   
*i=1, ..., n*  $n=699$

2)  $H_0$  : (M2) non signif  $\Leftrightarrow b_1 = b_2 = b_3 = 0 \Leftrightarrow$  aucune des trois var n'influe la proba  $P[Y=2]$ .

Modèle:  $P[Y_i = 2] = \frac{\exp(\mu)}{1 + \exp(\mu)}$

(4)

$H_1$ : (M2) significatif  $\Leftrightarrow \exists b_j \neq 0 \quad j=1,2,3 \Leftrightarrow$  au moins une des trois var influence  $P[Y=2]$

Modèle: (M2)

Tests utilisés: LR, SCORE, Wald: p-value  $< 10^{-4} \Rightarrow H_0$  rejetée.  
Donc (M2) significatif.

3) On étudie la proba que le cancer soit malin:  $P[Y=2]$

$H_0$ :  $V_1$  n'influe pas  $P[Y=2]$  ni  $V_2$  et  $V_3$  sont dans le modèle

$\Leftrightarrow b_1 = 0 \Rightarrow$  Modèle réduit:  $P[Y_i=2 | V_2, V_3] = \frac{\exp(\mu + b_2 V_2 + b_3 V_3)}{1 + \exp(\mu + b_2 V_2 + b_3 V_3)}$

$H_1$ :  $V_1$  influe  $P[Y=2]$  ni  $V_2$  et  $V_3$  dans le modèle

$\Leftrightarrow b_1 \neq 0 \Rightarrow$  Modèle (M2)

Stat de test:

$$Z = \frac{\hat{b}_1}{\sqrt{\hat{J}_{b_1}}} \xrightarrow{H_0} \chi^2(1) \quad (\text{avec SAS})$$

Valeur stat de test:  $Z = 36.54 \Rightarrow$  p-value =  $10^{-4} \Rightarrow H_0$  rejetée

Remarque: avec R, la stat de test est:

$$Z_2 = \frac{\hat{b}_1}{\sqrt{\hat{J}_{b_1}}} \xrightarrow{H_0} N(0,1)$$

$Z_2 = 6.04$  p-value  $< 10^{-9} \Rightarrow H_0$  rejetée.

Pareil:  $V_2$  et  $V_3$  influent aussi  $P[Y=2]$  donc aussi  $P[Y=1]$  c'est à dire la proba que le cancer soit bénin. Pour rappel:  $P[Y=1] + P[Y=2] = 1$ .

4)  $\hat{\beta}_1 = 0.59$   $\hat{\beta}_2 = 0.55$   $\hat{\beta}_3 = 0.71$   $\hat{\mu} = -7.59$

5)  $P[Y=2 | V_1=3, V_2=1, V_3=1] = \frac{\exp(\hat{\mu} + \hat{\beta}_1 \cdot 3 + \hat{\beta}_2 \cdot 1 + \hat{\beta}_3 \cdot 1)}{1 + \exp(\hat{\mu} + \hat{\beta}_1 \cdot 3 + \hat{\beta}_2 \cdot 1 + \hat{\beta}_3 \cdot 1)}$

6)  $P[Y=2]$  augment avec  $V_3$  car  $\hat{\beta}_3 > 0$ . Donc la première patiente a une proba plus grande pour que le cancer soit malin.

7) On utilise SAS :

Sur 458 patientes qui ont un cancer bénin, le modèle (M2) prévoit bien 446 cas.

Sur 241 patientes avec cancer malin, (M2) prévoit bien 220 cas.

Taux de mauvaise prévision :  $\frac{1}{2} \left( \frac{12}{458} + \frac{21}{241} \right) = 0.05$   
Donc très bonne prévision par le modèle.

Exercice 3

1) C'est un modèle de régression multiple : avec  $Y = V_1$   
 $Y_i = \mu + \sum_{j=2}^9 \beta_j V_{j,i} + \epsilon_i$   $i=1, \dots, n$   $\epsilon_i \sim N(0, \sigma^2)$  indep.

2) Param  $\mu$  et  $\beta_2, \dots, \beta_9$   
Pour (M3) les param  $\beta_2, \dots, \beta_9$  sont estimés par la méthode des moindres carrés pendant que pour (M4) ils sont estimés par LASSO adaptatif.

3) Par (M3) aucune estimation de  $\beta_j$  n'est égale à 0.  
Par (M4)  $\hat{\beta}_2, \hat{\beta}_4, \hat{\beta}_5, \hat{\beta}_6, \hat{\beta}_8, \hat{\beta}_9$  sont 0 (ou proches)  
Par (M3) et (M4) :  $\hat{\beta}_1$  et  $\hat{\beta}_3$  sont  $\neq 0$ .