

Examen final - 7 novembre 2016 Durée 2h - Documents interdits

Exercice 1 (8 pts)

1) (2 pts) Expliquer l'objectif d'une modélisation log-linéaire graphique ?

2) (2 pts) Représenter le modèle log-linéaire graphique associé à l'équation suivante

$$x_1 * x_2 * x_3 + x_1 * x_4 + x_4 * x_5 * x_6 .$$

Donner deux relations d'indépendance conditionnelle entre les variables qui le composent.

3) (2 pts) Nous considérons 100 réalisations suivant deux variables qualitatives ayant chacune deux modalités. Donner la vraisemblance associée au modèle log-linéaire graphique poissonien d'indépendance. Montrer qu'il s'agit bien d'un modèle linéaire généralisé.

4) (2 pts) Nous disposons d'un jeu de données R (`data.frame`) nommée `vars` contenant 8 variables qualitatives. Donner le code R permettant de mettre en oeuvre le modèle saturé, le modèle d'indépendance et de déterminer le meilleur modèle log-linéaire graphique au sens du critère AIC par une méthode descendante.

Exercice 2 (7 pts)

On considère un n -échantillon $(Y_1, x_1), \dots, (Y_n, x_n)$ du couple (Y, x) où Y est une variable aléatoire à valeurs dans $\{0, 1\}$ et x est une variable fixée prenant également les valeurs $\{0, 1\}$. On suppose qu'il existe une variable latente Y_i^* telle que $Y_i^* | x_i \sim \mathcal{N}(a + bx_i, 1)$ et

$$Y_i = \begin{cases} 1 & \text{si } Y_i^* \leq 0 \\ 0 & \text{si } Y_i^* > 0 \end{cases} .$$

On note $\Phi(\cdot)$ la fonction de répartition de la loi normale centrée réduite ($\Phi(u) = \mathbb{P}(\mathcal{N}(0, 1) \leq u)$).

1 (2 pts) Montrer qu'il s'agit d'un modèle linéaire généralisé. En suivant les notations du cours quant aux familles exponentielles à un paramètre de nuisance, on explicitera les paramètres ϕ , θ , les fonctions r et b ainsi que les régresseurs à considérer.

2 (1 pt) Montrer que ce n'est pas le lien canonique qui a été choisi.

3 (2 pts) Nous disposons de 100 observations (y_i, x_i) décrites dans la table de contingence suivante

	$y = 0$	$y = 1$	
$x = 0$	26	26	.
$x = 1$	32	16	

On note $\alpha = \Phi(a)$ et $\beta = \Phi(a + b)$ donner la log-vraisemblance en fonction de α et β .

4 (2 pts) Calculer l'estimateur du maximum de vraisemblance du couple (α, β) .

Exercice 3 (7 pts)

Commenter en détails le fichier R Markdown fourni en annexes. Il s'agit de données recueillies par l'Insee de février à avril 2003 dans le cadre de l'étude "Histoire de vie".

Annexes Examen HMMA304

Jean-Michel Marin

7 novembre 2016

```
library(questionr)
data(hdv2003)
d <- hdv2003
d <- d[,-c(1,5,7,9,11,12,13)]
summary(d)
```

```
##      age      sexe
## Min.   :18.00  Homme: 899
## 1st Qu.:35.00  Femme:1101
## Median :48.00
## Mean   :48.16
## 3rd Qu.:60.00
## Max.   :97.00
##
##                                     nivetud
## Enseignement technique ou professionnel court      :463
## Enseignement superieur y compris technique superieur:441
## Derniere annee d'etudes primaires                   :341
## 1er cycle                                           :204
## 2eme cycle                                           :183
## (Other)                                             :256
## NA's                                               :112
##
##      occup      freres.soeurs
## Exerce une profession:1049  Min.   : 0.000
## Chomeur                  : 134  1st Qu.: 1.000
## Etudiant, eleve         : 94  Median : 2.000
## Retraite                 : 392  Mean   : 3.283
## Retire des affaires     : 77  3rd Qu.: 5.000
## Au foyer                 : 171  Max.   :22.000
## Autre inactif           : 83
##
##      relig      lecture.bd peche.chasse cuisine
## Praticquant regulier    :266  Non:1953  Non:1776  Non:1119
## Praticquant occasionnel :442  Oui: 47   Oui: 224   Oui: 881
## Appartenance sans pratique :760
## Ni croyance ni appartenance:399
## Rejet                   : 93
## NSP ou NVPR             : 40
##
##      bricol      cinema      sport      heures.tv
## Non:1147  Non:1174  Non:1277  Min.   : 0.000
## Oui: 853  Oui: 826  Oui: 723  1st Qu.: 1.000
##
##                                     Median : 2.000
##                                     Mean   : 2.247
##                                     3rd Qu.: 3.000
##                                     Max.   :12.000
##                                     NA's   :5
```

```
d$age <- cut(d$age, c(16, 25, 45, 65, 97), right = FALSE, include.lowest = TRUE)
d$etud <- d$nivetud
levels(d$etud)
```

```
## [1] "N'a jamais fait d'etudes"
## [2] "A arrete ses etudes, avant la derniere annee d'etudes primaires"
## [3] "Derniere annee d'etudes primaires"
## [4] "1er cycle"
## [5] "2eme cycle"
## [6] "Enseignement technique ou professionnel court"
## [7] "Enseignement technique ou professionnel long"
## [8] "Enseignement superieur y compris technique superieur"
```

```
levels(d$etud) <- c("Primaire", "Primaire", "Primaire",
  "Secondaire", "Secondaire", "Technique/Professionnel",
  "Technique/Professionnel", "Supérieur")
freq(d$etud)
```

```
##           n    % val%
## Primaire    466 23.3 24.7
## Secondaire  387 19.4 20.5
## Technique/Professionnel 594 29.7 31.5
## Supérieur   441 22.1 23.4
## NA          112  5.6  NA
```

```
levels(d$etud)
```

```
## [1] "Primaire"          "Secondaire"
## [3] "Technique/Professionnel" "Supérieur"
```

```
d$etud <- addNA(d$etud)
levels(d$etud)
```

```
## [1] "Primaire"          "Secondaire"
## [3] "Technique/Professionnel" "Supérieur"
## [5] NA
```

```
d <- na.omit(d[,-3])
summary(d)
```

```
##      age      sexe      occup      freres.soeurs
## [16,25]:168 Homme: 895 Exerce une profession:1047 Min.   : 0.000
## [25,45]:705 Femme:1100 Chomeur           : 131 1st Qu.: 1.000
## [45,65]:742          Etudiant, eleve   :  94 Median : 2.000
## [65,97]:380          Retraite           : 392 Mean   : 3.283
##          Retire des affaires :  77 3rd Qu.: 5.000
##          Au foyer           : 171 Max.   :22.000
##          Autre inactif      :  83
##          relig      lecture.bd peche.chasse cuisine
## Praticquant regulier      :266 Non:1948 Non:1772 Non:1116
## Praticquant occasionnel   :441 Oui:  47  Oui: 223  Oui: 879
## Appartenance sans pratique :759
## Ni croyance ni appartenance:397
## Rejet                      : 92
## NSP ou NVPR                : 40
##
## bricol      cinema      sport      heures.tv
```

```
## Non:1145 Non:1171 Non:1275 Min. : 0.000
## Oui: 850 Oui: 824 Oui: 720 1st Qu.: 1.000
##
## Median : 2.000
## Mean : 2.247
## 3rd Qu.: 3.000
## Max. :12.000
##
##
##          etud
## Primaire      :465
## Secondaire    :387
## Technique/Professionnel:591
## Supérieur     :440
## NA            :112
##
##
```

```
levels(d$sport)
```

```
## [1] "Non" "Oui"
```

```
model0 <- glm(sport ~ sexe,data=d,family=binomial)
summary(model0)
```

```
##
## Call:
## glm(formula = sport ~ sexe, family = binomial, data = d)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.0255 -1.0255 -0.8813  1.3373  1.5058
##
## Coefficients:
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.36836    0.06799  -5.418 6.03e-08 ***
## sexeFemme   -0.37707    0.09374  -4.022 5.76e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
##      Null deviance: 2609.2  on 1994  degrees of freedom
## Residual deviance: 2593.0  on 1993  degrees of freedom
## AIC: 2597
##
## Number of Fisher Scoring iterations: 4
```

```
d$sexe <- relevel(d$sexe, ref="Femme")
model1 <- glm(sport ~ sexe,data=d,family=binomial)
summary(model1)
```

```
##
## Call:
## glm(formula = sport ~ sexe, family = binomial, data = d)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
```

```

## -1.0255 -1.0255 -0.8813 1.3373 1.5058
##
## Coefficients:
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) -0.74543 0.06454 -11.550 < 2e-16 ***
## sexeHomme 0.37707 0.09374 4.022 5.76e-05 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 2609.2 on 1994 degrees of freedom
## Residual deviance: 2593.0 on 1993 degrees of freedom
## AIC: 2597
##
## Number of Fisher Scoring iterations: 4

```

```

levels(d$sport)

```

```

## [1] "Non" "Oui"

```

```

d$sport <- relevel(d$sport, ref="Oui")
table(d$sport,d$sexe)

```

```

##
##      Femme Homme
## Oui   354   366
## Non   746   529

```

```

model2 <- glm(sport ~ sexe,data=d,family=binomial)
summary(model2)

```

```

##
## Call:
## glm(formula = sport ~ sexe, family = binomial, data = d)
##
## Deviance Residuals:
##      Min       1Q   Median       3Q      Max
## -1.5058 -1.3373  0.8813  1.0255  1.0255
##
## Coefficients:
## Estimate Std. Error z value Pr(>|z|)
## (Intercept) 0.74543 0.06454 11.550 < 2e-16 ***
## sexeHomme -0.37707 0.09374 -4.022 5.76e-05 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## (Dispersion parameter for binomial family taken to be 1)
##
## Null deviance: 2609.2 on 1994 degrees of freedom
## Residual deviance: 2593.0 on 1993 degrees of freedom
## AIC: 2597
##
## Number of Fisher Scoring iterations: 4

```

```

model3 <- glm(sport ~ sexe + etud + relig + heures.tv,data=d,family=binomial)
anova(model3,test="Chisq")

```

```

## Analysis of Deviance Table
##
## Model: binomial, link: logit
##
## Response: sport
##
## Terms added sequentially (first to last)
##
##
##           Df Deviance Resid. Df Resid. Dev Pr(>Chi)
## NULL                1994      2609.2
## sexe                1    16.21      1993      2593.0 5.668e-05 ***
## etud                4   319.19      1989      2273.8 < 2.2e-16 ***
## relig              5     1.65      1984      2272.2 0.8953896
## heures.tv          1    13.18      1983      2259.0 0.0002831 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

model4 <- glm(sport~.,data=d,family=binomial)
anova(model4,test="Chisq")

```

```

## Analysis of Deviance Table
##
## Model: binomial, link: logit
##
## Response: sport
##
## Terms added sequentially (first to last)
##
##
##           Df Deviance Resid. Df Resid. Dev Pr(>Chi)
## NULL                1994      2609.2
## age                3   213.226      1991      2396.0 < 2.2e-16 ***
## sexe                1    19.253      1990      2376.7 1.145e-05 ***
## occup              6    30.610      1984      2346.1 3.009e-05 ***
## freres.soeurs      1    24.625      1983      2321.5 6.963e-07 ***
## relig              5     4.401      1978      2317.1 0.4932429
## lecture.bd         1     3.730      1977      2313.4 0.0534341 .
## peche.chasse       1     0.123      1976      2313.2 0.7260114
## cuisine            1    17.894      1975      2295.3 2.335e-05 ***
## bricol             1    13.545      1974      2281.8 0.0002329 ***
## cinema             1    71.717      1973      2210.1 < 2.2e-16 ***
## heures.tv          1    21.002      1972      2189.1 4.588e-06 ***
## etud               4    67.891      1968      2121.2 6.325e-14 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

```

library(caret)

```

```

## Loading required package: lattice
## Loading required package: ggplot2

```

```
res3 <- train(formula(model3),data=d,method="glm",family=binomial,
  metric="Accuracy",trControl=trainControl(method="repeatedcv",number=10,
  repeats=50))
confusionMatrix(res3)
```

```
## Cross-Validated (10 fold, repeated 50 times) Confusion Matrix
##
## (entries are percentual average cell counts across resamples)
##
##           Reference
## Prediction  Oui  Non
##           Oui 16.0  9.7
##           Non 20.1 54.2
##
## Accuracy (average) : 0.7023
```

```
res4 <- train(formula(model4),data=d,method="glm",family=binomial,
  metric="Accuracy",trControl=trainControl(method="repeatedcv",number=10,
  repeats=50))
confusionMatrix(res4)
```

```
## Cross-Validated (10 fold, repeated 50 times) Confusion Matrix
##
## (entries are percentual average cell counts across resamples)
##
##           Reference
## Prediction  Oui  Non
##           Oui 18.8 10.7
##           Non 17.3 53.2
##
## Accuracy (average) : 0.7208
```

Correction Exam final

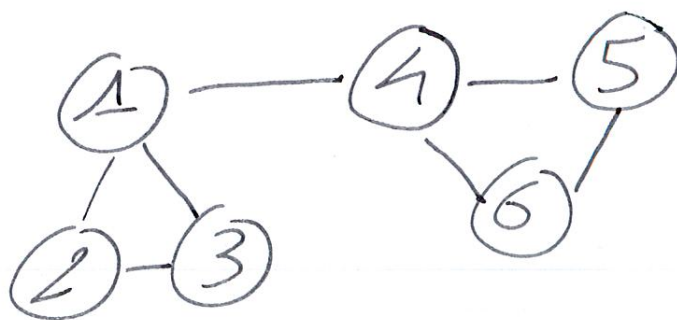
HPPA304 - 07/11/2016

①

Exercice 1

1] Les modèles log-linéaires graphiques sont utilisés pour étudier les relations existant entre différentes variables qualitatives.

2]



$$(K_5, K_6) \perp\!\!\!\perp (K_1, K_2, K_3) \mid K_4$$

$$K_2 \perp\!\!\!\perp (K_4, K_5, K_6) \mid K_1, K_3$$

3] 2 variables qualitatives

(2)

Les données sont synthétisées dans une table de contingence 2×2 .

Le modèle d'indépendance s'écrit

$$\log(M_{ij}) = \mu + \alpha_i + \beta_j \quad \begin{array}{l} i \in \{1, 2\} \\ j \in \{1, 2\} \end{array}$$

avec les contraintes

$$\alpha_1 = -\alpha_2$$

$$\beta_1 = -\beta_2$$

$$Y_{ij} \sim \mathcal{P}(\mu_{ij})$$

$$Y = (Y_{11}, Y_{12}, Y_{21}, Y_{22})$$

$$V(\mu, \alpha_1, \beta_1; Y) = \prod_{i=1}^2 \prod_{j=1}^2 \left[e^{-e^{\mu + \alpha_i + \beta_j}} \right]$$

$$e^{Y_{ij} [\mu + \alpha_i + \beta_j]}$$

$$Y_{ij}!$$

Il s'agit bien d'un modèle
linéaire généralisé car

(3)

i) famille exponentielle (triviale)

ii) $\text{log} [E(\mu_{ij})] = \mu + \alpha_i + \beta_j$
 $= \mathbf{K}_{ij}^T \theta$ avec

$$\theta = (\mu, \alpha_1, \alpha_2)$$

$$\mathbf{K}_{11} = (1, 1, 1)$$

$$\mathbf{K}_{12} = (1, 1, -1)$$

$$\mathbf{K}_{21} = (1, -1, 1)$$

$$\mathbf{K}_{22} = (1, -1, -1)$$

4) library (yRim)

model 1 ← vlmmod (n.r., data = vory)

model 2 ← vlmmod (n.r.1, data = vory)

model 3 ← backward (model 1,
crit = "aic")

4

Exercice 2

1) Soit $\pi_\kappa = \mathbb{P}(Y=1|\kappa)$

$$f(y; \pi_\kappa) = \pi_\kappa^y (1-\pi_\kappa)^{1-y} \quad \Gamma_{\{y \in \{0,1\}\}}$$

$$f(y; \pi_\kappa) = \exp \{ y \log(\pi_\kappa) + (1-y) \log(1-\pi_\kappa) \}$$

$$f(y; \pi_\kappa) = \exp \left\{ \underset{\Gamma(y)}{\overset{\{0,1\}}{y}} \log \left(\frac{\pi_\kappa}{1-\pi_\kappa} \right) + \log(1-\pi_\kappa) \right\}$$

La loi de Bernoulli appartient à la famille exponentielle

$$\theta = \log \left(\frac{\pi_\kappa}{1-\pi_\kappa} \right) \Rightarrow \pi_\kappa = \frac{e^\theta}{1+e^\theta}$$

$$h(\theta) = -\log \left(1 - \frac{e^\theta}{1+e^\theta} \right) = \log(1+e^\theta)$$

$$\phi = 1, \quad C(y, \phi) = 0$$

$$V(y) = \pi(y) \quad \{0,1\}$$

En velleurs,

(5)

$$\begin{aligned} \mathbb{E}(Y) &= \mathbb{P}(Y=1|K) \\ &= \mathbb{P}(Y^* \leq 0 | K) \\ &= \mathbb{P}(Y^* - \alpha - bK \leq -\alpha - bK | K) \\ &= \mathbb{P}(N(0,1) \leq -\alpha - bK | K) \\ &= \Phi(-\alpha - bK) \end{aligned}$$

$$\Rightarrow -\Phi^{-1}(\mathbb{E}(Y)) = \alpha + bK$$

IC s'agit donc bien d'un modèle linéaire généralisé.

$$\eta(q) = \Phi(-q)$$

$$\underline{2]} \text{ Nous avons } h'(q) = \frac{e^q}{1+e^q}$$

$\neq \eta(q)$ - ce n'est pas le lien canonique qui a été utilisé.

$$3] \underline{y} = (y_1, \dots, y_m) \quad \underline{\kappa} = (\kappa_1, \dots, \kappa_m) \quad (6)$$

$$\Delta V(\alpha, \beta; \underline{y}, \underline{\kappa})$$

$$= \sum_{i=1}^m \left[y_i \log(\Phi(-\alpha - \beta \kappa_i)) + (1 - y_i) \log(1 - \Phi(-\alpha - \beta \kappa_i)) \right]$$

$$\Delta V(\alpha, \beta; \underline{y}, \underline{\kappa}) =$$

$$\sum_{i=1}^m \left[y_i \log(\Phi(-\alpha - \beta \kappa_i)) + (1 - y_i) \log(\Phi(\alpha + \beta \kappa_i)) \right]$$

$$\Delta V(\alpha, \beta; \underline{y}, \underline{\kappa}) = 16 \log(\Phi(-\alpha - \beta)) + 26 \log(\Phi(-\alpha)) + 32 \log(\Phi(\alpha + \beta)) + 26 \log(\Phi(\alpha))$$

(7)

$$L V(\alpha, \beta; \underline{y}, \kappa)$$

$$= 16 \log(1-\beta) + 26 \log(1-\alpha) \\ + 32 \log(\beta) + 26 \log(\alpha)$$

$$\Downarrow \frac{dL V}{d\alpha}(\alpha^*, \beta^*; \underline{y}, \kappa) = -\frac{26}{1-\alpha^*} + \frac{26}{\alpha^*} = 0$$

$$\frac{dL V}{d\beta}(\alpha^*, \beta^*; \underline{y}, \kappa) = -\frac{16}{1-\beta^*} + \frac{32}{\beta^*} = 0$$

$$\Leftrightarrow \begin{cases} 26\alpha^* = 26(1-\alpha^*) \\ 32(1-\beta^*) = 16\beta^* \end{cases}$$

$$\Leftrightarrow \begin{cases} \alpha^* = \frac{26}{52} = \frac{1}{2} \\ \beta^* = \frac{32}{8} = \frac{2}{3} \end{cases}$$

(8)

$$\frac{d^2 L V}{d\alpha d\beta} (\alpha, \beta; \underline{y}, \underline{K}) = 0$$

$$\frac{d^2 L V}{(d\alpha)^2} (\alpha, \beta; \underline{y}, \underline{K}) = -\frac{26}{\alpha^2} - \frac{16}{(1-\alpha)^2} < 0$$

$$\frac{d^2 L V}{(d\beta)^2} (\alpha, \beta; \underline{y}, \underline{K}) = -\frac{32}{\beta^2} - \frac{26}{(1-\beta)^2} < 0$$

La fonction $L V$ est strictement
- concave -

Ainsi,

$$\hat{\alpha} = \frac{13}{21} \quad \hat{\beta} = \frac{16}{29}$$

Exercice 3

modèle 0, modèle 1, modèle 2
trois modèles de régression logistique
où la variable spat est expliquée
par la psea

Les trois modèles sont identiques, seules les modalités de référence des facteurs (à expliquer et explicatif) - changent. (9)

Pour le modèle 0,

$$TP(\text{Sport} = \text{"Oui"} \mid \text{Sexe} = \text{"Femme"})$$

$$\# \frac{e^{-0.37 - 0.38}}{1 + e^{-0.37 - 0.38}} \# 0.32$$

$$TP(\text{Sport} = \text{"Oui"} \mid \text{Sexe} = \text{"Homme"})$$

$$\# \frac{e^{-0.37}}{1 + e^{-0.37}} \# 0.40$$

Pour le match 1,

$$P(\text{Sport} = \text{"Oui"} \mid \text{Sexe} = \text{"Homme"})$$

$$\# \frac{e^{-0.74} + 0.37}{1 + e^{-0.74} + 0.37} \# 0.40$$

$$P(\text{Sport} = \text{"Oui"} \mid \text{Sexe} = \text{"Femme"})$$

$$\# \frac{e^{-0.75}}{1 + e^{-0.75}} \# 0.32$$

Pour le match 2,

$$P(\text{Sport} = \text{"Non"} \mid \text{Sexe} = \text{"Homme"})$$

$$\# \frac{e^{0.74} - 0.37}{1 + e^{0.74} - 0.37} \# 0.60$$

$$P(\text{Sport} = \text{"Non"} \mid \text{Sexe} = \text{"Femme"})$$

$$\# \frac{e^{0.75}}{1 + e^{0.75}} \# 0.68$$

Les modèles de régression

logistique modèle 3 et modèle 4
sont comparés à l'aide
d'une technique de validation
croisée à 10 ensembles
répétés 50 fois.

Le modèle 4 a de meilleures
performances en terme de
taux d'erreur de classification.