Examen final

Durée 2h - Sans document - 5 novembre 2018

Exercice 1 (10 pts)

 $\bf 1$ (4 pts) On rappelle que la loi Beta de paramètre $\alpha>0$ et $\beta>0$ admet pour densité de probabilité

$$f(x) = \frac{\Gamma(\alpha + \beta)}{\Gamma(\alpha)\Gamma(\beta)} x^{\alpha - 1} (1 - x)^{\beta - 1} \mathbf{1}_{]0,1[}(x).$$

Décrire les procédures mises en oeuvre par l'intermédiaire du code R suivant.

```
f <- function(x)
{
    sqrt((1-x)*x)/(exp(x)+log(1+x))
}

x <- runif(10000)
mean(f(x))
sqrt(var(f(x)))

x <- rbeta(10000,3/2,3/2)
h <- function(x)
{
    1/(exp(x)+log(1+x))
}
gamma(3/2)^2/gamma(3)*mean(h(x))
gamma(3/2)^2/gamma(3)*sqrt(var(h(x)))</pre>
```

- **2 (2 pts)** Donner le code R permettant de générer une réalisation suivant une variable aléatoire X telle que $\mathbb{P}(X=1)=0.6$, $\mathbb{P}(X=2)=0.3$ et $\mathbb{P}(X=3)=0.1$.
- **2 (4 pts)** Pour quel type de données et répondre à quelle question les modèles loglinéaires graphiques sont-ils utilisés? Donner un exemple. Quelle procédure peut-on mettre en oeuvre pour sélectionner le modèle log-linéaire graphique le plus adapté à un ensemble de données?

Exercice 2 (4 pts)

On considère n variables aléatoires indépendantes y_1, \ldots, y_n telles que y_i est distribué suivant une loi binomiale négative de paramètres (h, π_i) . Nous supposons que h est fixé. De manière générique, la loi négative binomiale modélise, dans le contexte d'une suite d'épreuves de Bernoulli indépendantes, le nombre d'échecs nécessaires pour obtenir h succès, π représentant la probabilité de succès :

$$f(y;\pi) = C_{h+y-1}^{y} \pi^{h} (1-\pi)^{y} \mathbf{1}_{\mathbb{N}}(y)$$
.

1 (1 pt) Montrer que la loi binomiale négative appartient à la famille exponentielle.

```
2 (1 pt) Calculer \mathbb{E}(y) et \mathbb{V}(y).
```

Pour tout i = 1, ..., n, on suppose que $\log(1 - \pi_i) = \beta_1 + \beta_2 x_i$.

3 (1 pt) Montrer qu'il s'agit d'un modèle linéaire généralisé. La fonction de lien canonique a-t-elle été utilisée?

4 (1 pt) Donner la log-vraisemblance et les équations de vraisemblance.

Exercice 3 (6 pts)

Expliquer en détails les procédures mises en oeuvre par l'intermédiaire du code R suivant.

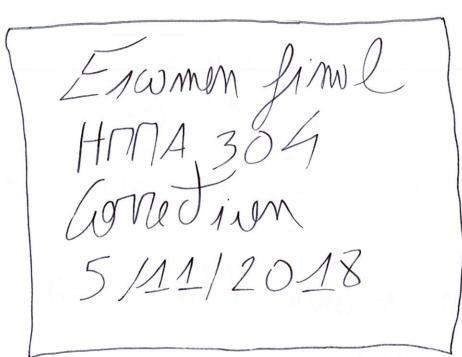
```
library(glmnet)
library(nnet)
library(caret)

set.seed(8750)
n <- 1000
p <- 200
x <- matrix(rnorm(n*p), nrow=n, ncol=p)
colnames(x) <- paste("v", 1:p, sep="")

y <- rep(0,n)
for (i in 1:n)
{
    p <- exp(sum(x[i,]))/(1+exp(sum(x[i,])))
    y[i] <- sample(c(0,1),1,prob=c(1-p,p))
}

train_rows <- sample(1:n,0.66*n)
x.train <- x[train_rows, ]
x.test <- x[-train_rows, ]</pre>
```

```
y.train <- y[train_rows]</pre>
y.test <- y[-train_rows]</pre>
train <- data.frame(y=y.train,x.train)</pre>
test <- data.frame(y=y.test,x.test)</pre>
model1 <- multinom(y~.,family="binomial",data=train)</pre>
# weights: 202 (201 variable)
# initial value 457.477139
# iter 10 value 3.602629
# iter 20 value 1.458515
# iter 30 value 0.921504
# iter 40 value 0.191383
# iter 50 value 0.009916
# iter 60 value 0.001003
# iter 70 value 0.000214
# final value 0.000086
# converged
mean(predict(model1,test)!=as.factor(y.test))
# [1] 0.1647059
model <- glmnet(x.train, as.factor(y.train), family="binomial", alpha=1)</pre>
model2 <- train(x.train, as.factor(y.train), method="glmnet", family="binomial",</pre>
                metric="Accuracy", trControl=
                   trainControl(method="repeatedcv", number=5, repeats=10),
                tuneGrid=data.frame(alpha=1, lambda=model$lambda))
mean(predict(model2,x.test)!=as.factor(y.test))
# [1] 0.1735294
y.star < - rep(0,340)
for (i in 1:340)
  p \leftarrow \exp(sum(x.test[i,]))/(1+exp(sum(x.test[i,])))
  if (p>=1/2) y.star[i] <- 1
mean(y.star!=y.test)
# [1] 0.02647059
```



Escencia 1 Now soulwitons reproches $\int_{0}^{\infty} \frac{\sqrt{\chi(1-\chi)}}{e^{\chi} + \log(11/4)} d\chi = I$ Dum im premier temps, me methods who Thanks and to the stronger of mine in sewl. En effet, $T = E V_{(2)1} \left[\frac{V \times (1-x)}{e^{x} + \log_{x}(1+x)} \right]$ generi 10000/Califotigns! juint

On utime I jon Duns em seame Temps, un T = t $Stu\left(\frac{3}{2},\frac{3}{2}\right) \left[\frac{\left(\Gamma\left(\frac{3}{2}\right)\right)^{2}}{\Gamma\left(3\right)} \frac{1}{e^{x} + loy(1tx)} \right]$ Rais for simply 10000 realisations furious to some loss between $\left(\frac{3}{2},\frac{3}{2}\right)$ at some estimant. 10000 i=2 e/i+lly [1+/i] [7/3]

Leant-type de Te et myolement
et im.

Jumple(C(1,2,3), 1, pull = (0.6, 0.3, 0.1))3 Vioin aurs. Escencia 2 J(J) = eng {h by (T) + y by (1-T)} (y) 0= loy (1-TT) => T=1-e 1 (0) = - h loy (1-e0) 2) E[4] = 1/0/= h = h(1-11) $V[y] = f'(0) = h \frac{e^{0}}{1-e^{0}/2} = h \left(\frac{1-11}{1-2}\right)$

3) logy (1-TT) = B1+B2 K $=> 1-11 = e^{B_1+B_2K}$ $=> 11 = 1 - e^{B_1+B_2K}$ => E[y] = \[\langle \langle \beta \langle \ = J(B1+B2K) J(M) = M/C = 1/M)

(1st m mouth himmin

generaling and atilism

vingi que lim amonique. 1/2/B1/B2 $= \int_{i=1}^{\infty} \frac{1}{2^{i}} \int_{i=1}^{\infty} \frac{1}{2^{i}} \left(\frac{1}{2^{i}} + \frac{1}{2^{i}} \frac{1}{2^{i}} \right) + \sum_{i=1}^{\infty} \frac{1}{2^{i}} \left(\frac{1}{2^{i}} + \frac{1}{2^{i}} \frac{1}{2^{i}} \right)$

 $\frac{dV(\beta_1, \beta_2) = 0}{d\beta_1}$ $= \sum_{i=1}^{m} \mathcal{J}_i - \lambda \sum_{i=1}^{m} \left[\frac{e^{\beta_1 + \beta_2 M_i}}{1 - e^{\beta_1 + \beta_2 M_i}} \right] = 0$ (1/2) $=\sum_{i=1}^{m} M_{i}N_{i} - M_{i} = \sum_{i=1}^{m} \frac{M_{i}e^{B_{1}+B_{2}M_{i}}}{1 - e^{B_{1}+B_{2}M_{i}}} = 0$ $=\sum_{i=1}^{m} M_{i}N_{i} - M_{i} = \sum_{i=1}^{m} \frac{M_{i}e^{B_{1}+B_{2}M_{i}}}{1 - e^{B_{1}+B_{2}M_{i}}} = 0$ $=\sum_{i=1}^{m} M_{i}N_{i} - M_{i} = \sum_{i=1}^{m} M_{i}e^{B_{1}+B_{2}M_{i}}$ $=\sum_{i=1}^{m} M_{i}N_{i} - M_{i} = \sum_{i=1}^{m} M_{i}e^{B_{1}+B_{2}M_{i}}$ $=\sum_{i=1}^{m} M_{i}N_{i} - M_{i}e^{B_{1}+B_{2}M_{i}}$ On gener m jen de donnen Contempt 1000 invhirintes, 200 vouinthe peutitives to me vouinth trimin it expliques smirunt me much in regression hogistique start tous les coefficients in rome blus sont eyeme

Om isola 33% ver dumin (6) On amon astonicalonem regression legistique - chessique it regularism logistique avece regularisotion 11 about la jurumita de regularisot ison est culibré por volintation origine it 5 ensembles résters 10 fois-