

Examen de contrôle continu

Durée 1h30 - Sans document - 14 octobre 2020

Exercice 1 (8 pts)

On considère n variables aléatoires indépendantes y_1, \dots, y_n telles que

$$f(y_i; \pi_i) = \pi_i(1 - \pi_i)^{y_i-1} \mathbf{1}_{\mathbb{N}^*}(y_i).$$

Pour tout $i = 1, \dots, n$, on suppose que $\log(1 - \pi_i) = \alpha + \beta x_i$ où $x_i \in \mathbb{R}$ est supposé connu, α et β sont des paramètres inconnus.

1 (2 pts) Montrer qu'il s'agit d'un modèle linéaire généralisée. La fonction de lien canonique a-t-elle été utilisée ?

2 (3 pts) Donner la log-vraisemblance et les équations de vraisemblance. Pouvons-nous calculer les expressions analytique des estimateurs du maximum de vraisemblance de α et β ?

3 (3 pts) Calculer la matrice d'information de Fisher apportée par (y_1, \dots, y_n) sur les paramètres α et β notée $I_n(\alpha, \beta)$.

Exercice 2 (4 pts)

Nous nous intéressons ici à l'étude de la pièce ayant causé l'explosion de la navette spatiale Challenger en 1986. L'étanchéité du moteur de la navette spatiale est assurée par six pièces identiques appelées "O-ring". L'explosion de la navette Challenger est due à la défaillance d'au moins l'une de ces pièces.

Au cours des 24 vols précédents d'une navette spatiale, nous disposons des données suivantes : la variable `temp` qui correspond à la température au moment du lancement et la variable `defa` qui vaut 0 si aucun des "O-ring" n'a été endommagé au cours du lancement et 1 si au moins l'un d'entre eux a été endommagé.

Proposer une méthode permettant d'estimer la probabilité de défaillance d'au moins un "O-ring" pour une température de 31 degrés Fahrenheit (température au moment du lancement de la navette challenger). Nous supposons que les données sont stockées dans un `data.frame` R nommé `challenger`, donner le code R associé.

Exercice 3 (2 pts)

Expliquer en quoi consiste une régression ridge et donner l'expression explicite de l'estimateur ridge.

Exercice 4 (6 pts)

Expliquer en détails les résultats produits par le code R ci-dessous

```
library(glmnet)
library(caret)

X <- matrix(rnorm(20*20),20)
beta <- rep(0.5,20)
y <- X%*%beta+rnorm(20,sd=2)

app <- data.frame(y=y,X)
names(app)[-1] <- paste("x",1:20,sep="")

model <- lm(y~-1+.,data=app)
beton1 <- as.numeric(model$coefficients)

model <- glmnet(X,y,family="gaussian",nlambda=50,alpha=1,intercept=FALSE)
model <- train(y~.,data=app,method="glmnet",intercept=FALSE,metric="RMSE",
              trControl=trainControl(method="repeatedcv",number=20),
              tuneGrid=data.frame(alpha=1,lambda=model$lambda))

opti <- as.numeric(model$bestTune$lambda)
model$results[model$results[,2]==opti,1:3]
#   alpha   lambda   RMSE
# 29     1 0.2516718 1.873211
beton2 <- coef(model$finalModel,opti)[-1]

sum(abs(beton1))
# [1] 26.90182
sum(beton2==0)
# [1] 9
sum(abs(beton2))
# [1] 4.785677

library(plotmo)
plot_glmnet(glmnet(X,y,family="gaussian",nlambda=50,alpha=1,intercept=FALSE),
            s=as.numeric(model$bestTune$lambda),label=2)

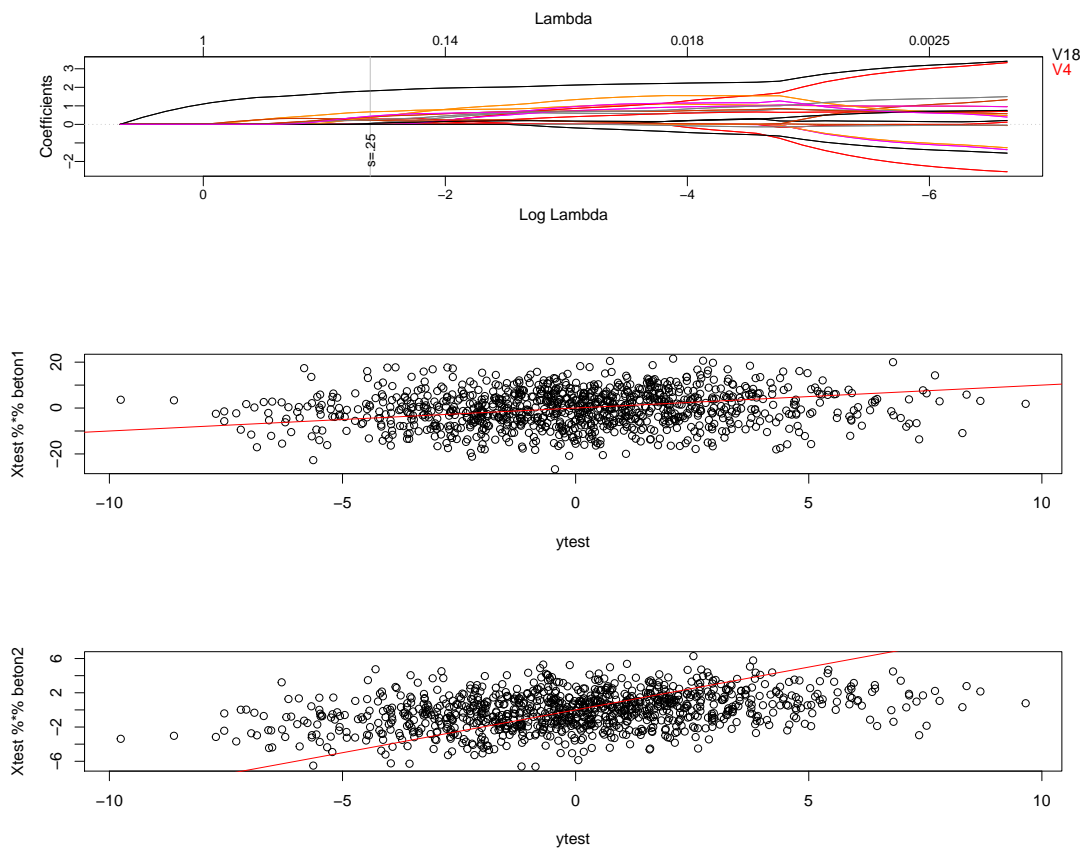
Xtest <- matrix(rnorm(1000*20),1000)
test <- data.frame(Xtest)
```

```
names(test) <- paste("x",1:20,sep="")
ytest <- Xtest%*%beta+rnorm(1000,sd=2)
```

```
plot(ytest, Xtest%*%beton1)
abline(a=0,b=1,col="red")
```

```
plot(ytest, Xtest%*%beton2)
abline(a=0,b=1,col="red")
```

```
sqrt(mean((Xtest%*%beta-ytest)^2))
# [1] 1.975905
sqrt(mean((Xtest%*%beton1-ytest)^2))
# [1] 8.144402
sqrt(mean((Xtest%*%beton2-ytest)^2))
# [1] 2.988543
```



①

Correction exam
de contrôle
- continu # ППА304
14/10/2020

Exercice 1

$$\underline{1} \quad f(y; \pi) = \pi (1-\pi)^{y-1} \frac{\pi(y)}{N^*}$$

$$f(y; \pi) = \exp \left[y \log \frac{\pi(y)}{N^*} + \log \left(\frac{\pi}{1-\pi} \right) \right]$$

$$\theta = \log(1-\pi)$$

$$\Rightarrow \pi = 1 - e^\theta$$

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

Famille exponentielle
 Oui - c'est le lien
 canonique.

(2)

$$E(Y) = \frac{1}{\pi} = \frac{1}{1 - e^{\alpha + \beta \mu}}$$

$$2] \quad L(\alpha, \beta) = \sum_{i=1}^M \binom{M}{\mu_i} [\alpha + \beta \mu_i]$$

$$+ \sum_{i=1}^M \log(1 - e^{\alpha + \beta \mu_i})$$

$$\frac{dL}{d\alpha}(\alpha, \beta) = \sum_{i=1}^M \mu_i - M - \sum_{i=1}^M \left[\frac{e^{\alpha + \beta \mu_i}}{1 - e^{\alpha + \beta \mu_i}} \right]$$

$$\frac{dL}{d\beta}(\alpha, \beta) = \sum_{i=1}^M \mu_i \mu_i - \sum_{i=1}^M \mu_i - \sum_{i=1}^M \left[\frac{\mu_i e^{\alpha + \beta \mu_i}}{1 - e^{\alpha + \beta \mu_i}} \right]$$

Les 2 expressions analytiques
 de l'ENT.

$$\frac{d^2 L V}{(d\alpha)^2}(\alpha, \beta) = - \sum_{i=1}^M \frac{e^{\alpha + \beta N_i} (1 - e^{\alpha + \beta N_i}) + (e^{\alpha + \beta N_i})^2}{(1 - e^{\alpha + \beta N_i})^2} \quad (3)$$

$$\frac{d^2 L V}{(d\alpha)^2}(\alpha, \beta) = - \sum_{i=1}^M \frac{e^{\alpha + \beta N_i}}{(1 - e^{\alpha + \beta N_i})^2}$$

$$\frac{d^2 L V}{d\alpha d\beta}(\alpha, \beta) = - \sum_{i=1}^M \frac{N_i e^{\alpha + \beta N_i}}{(1 - e^{\alpha + \beta N_i})^2}$$

$$\frac{d^2 L V}{(d\beta)^2}(\alpha, \beta) = - \sum_{i=1}^M \frac{N_i^2 e^{\alpha + \beta N_i}}{(1 - e^{\alpha + \beta N_i})^2}$$

These - as terms sont déterministes
 On m'explique visiblement l'
 information de Fisher

Exercice 2

4

Théorie de régression logistique

→ Variable à expliquer
deja qui est binaire

→ Variable explicative
temp

$$P(\text{deja} = 1 | \text{temp}) = \frac{e^{\alpha + \beta \text{temp}}}{1 + e^{\alpha + \beta \text{temp}}}$$

model <- glm(deja ~ 1 + temp,

family = "binomial")

predict(model, data.frame(temp = 31))

5

Exercício 3

Problema de regressão regularizada
Penalidade - norma semi - euclidiana

$$\hat{\beta}^R \in \arg \min_{\beta} \sum_{i=1}^M (y_i - \mathbf{1}_i^T \beta)$$

norma - euclidiana $\sum_{j=1}^p \beta_j^2 \leq \lambda$

Problema equivalente

$$\min_{\beta} \sum_{i=1}^M (y_i - \mathbf{1}_i^T \beta)^2 + \lambda \sum_{j=1}^p \beta_j^2$$

função objetivo

$$\hat{\beta}^R = (X^T X + \lambda I_p)^{-1} X^T y$$

matr. $X = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$

matriza dos normalizados
esquadrados

Exercice 4

(6)

On dispose de 20 données suivant
le modèle :

$$y = \sum_{j=1}^{20} \beta_j x_j + \mu$$

avec $\mu \sim N(0, 4)$

On estime le vecteur β par la
méthode des moindres-carrés
 \Rightarrow objet R beton 1

On estime le vecteur β par une
technique LASSO dont le paramètre
de régularisation est calibré par
validation croisée \Rightarrow 20 ensembles
(leave-one-out cross validation)

\Rightarrow objet R beton 2

On trouve que g coefficients
sont estimés à 0 et que la
somme des valeurs absolues de beton 1

est largement supérieure à
celle des valeurs objectives
en item 2. C'est l'effet
de la régularisation LASSO.

(7)

Graph 1: évolution des valeurs
des estimations des paramètres LASSO
en fonction de la contrainte

Graph 2: estimation PLO sans
mais trois variables

Graph 3: estimation LASSO divisée
mais 2 variables

Sur les données étudiées, très bons
résultats pour le LASSO, très
proche de l'Oracle.