

Examen de contrôle continu sans document - 12 octobre 2016

Exercice 1 (5 pts)

Nous considérons le modèle de régression linéaire tel que pour tout $i \in \{1, \dots, 10\}$ et $j \in \{1, \dots, 3\}$, nous avons

$$y_{i,j} = \alpha x_i + \epsilon_{i,j}$$

où α est un paramètre inconnu et $\epsilon_{i,j}$ est une variable aléatoire réelle telle que :

$$\mathbb{E}(\epsilon_{i,j}) = 0, \quad \mathbb{V}(\epsilon_{i,j}) = 1,$$

$$\mathbb{C}(\epsilon_{i,j}, \epsilon_{i',j'}) = 0 \text{ lorsque } i \neq i',$$

$$\mathbb{C}(\epsilon_{i,j}, \epsilon_{i',j'}) = 1/2 \text{ lorsque } i = i'.$$

1 (2 pts) Donner l'expression de l'estimateur des Moindres Carrés Ordinaires $\hat{\alpha}$ de α et calculer sa variance.

2 (3 pts) Donner le code R permettant de calculer l'estimateur des Moindres Carrés Généralisés de α .

Exercice 2 (9 pts)

On considère n variables aléatoires indépendantes y_1, \dots, y_n telles que y_i est distribué suivant une loi négative binomiale de paramètres (h, π_i) . Nous supposons que h est fixé. La loi négative binomiale modélise, dans le contexte d'une suite d'épreuves de Bernoulli indépendantes, le nombre d'essais nécessaires pour obtenir h succès, π_i représente la probabilité de succès. Nous avons

$$f(y_i; \pi_i) = C_{y_i-1}^{h-1} \pi_i^h (1 - \pi_i)^{y_i-h} \mathbf{1}_{\{h, h+1, \dots\}}(y_i).$$

Pour tout $i = 1, \dots, n$, on suppose que $\log(\pi_i/(1 - \pi_i)) = \alpha + \beta x_i$ 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 généralisé. La fonction de lien canonique a-t-elle été utilisée ?

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

3 (2 pts) Calculer la matrice d'information de Fisher apportée par (y_1, \dots, y_n) sur les paramètres α et β notée $I_n(\alpha, \beta)$.

4 (3 pts) Modifier le modèle en utilisant le lien canonique et donner les équations de vraisemblance.

Exercice 3 (6 pts)

1 (3 pts) Expliquer les résultats produits par la code R suivant (le jeu de données est décrit en annexe).

```
library(klaR)
data(GermanCredit)
model1 <- glm(credit_risk~duration,family=binomial,data=GermanCredit)
model2 <- glm(credit_risk~.,data=GermanCredit,family=binomial)
model3 <- step(model2,direction="backward")
model4 <- step(model2,direction="backward",k=log(n))
model5 <- glm(credit_risk~1,data=GermanCredit,family=binomial)
model6 <- step(model5,direction="both",scope=formula(model2))
```

2 (3 pts) Donner le code R permettant d'estimer par validation croisée à 2 ensembles répétée 10 fois les taux d'erreurs associés aux modèles considérés dans la question précédente.

Annexes

GermanCredit (klaR) R Documentation
Statlog German Credit

Description

The dataset contains data of past credit applicants. The applicants are rated as good or bad. Models of this data can be used to determine if new applicants present a good or bad credit risk.

A data frame containing 1,000 observations on 21 variables.

status

factor variable indicating the status of the existing checking account, with levels ... < 100 DM, 0 <= ... < 200 DM, ... >= 200 DM/salary for at least 1 year and no checking account

duration

duration in months

credit_history

factor variable indicating credit history, with levels no credits taken/all credits paid back duly, all credits at this bank paid back duly, existing credits paid back duly till now, delay in paying off in the past and critical account/other credits existing

purpose

factor variable indicating the credit's purpose, with levels car (new), car (used), furniture/equipment, radio/television, domestic appliances, repairs, education, retraining, business and others

amount

credit amount.

savings

factor. savings account/bonds, with levels ... < 100 DM, 100 <= ... < 500 DM, 500 <= ... < 1000 DM, ... >= 1000 DM and unknown/no savings account

employment_duration

ordered factor indicating the duration of the current employment, with levels unemployed, ... < 1 year, 1 <= ... < 4 years, 4 <= ... < 7 years and ... >= 7 years

installment_rate

installment rate in percentage of disposable income.

personal_status_sex

factor variable indicating personal status and sex, with levels male:divorced/separated, female:divorced/separated/married, male:single, male:married/widowed and female:single

other_debtors

factor. Other debtors, with levels none, co-applicant and guarantor

present_residence

present residence since

property

factor variable indicating the client's highest valued property, with levels real estate, building society savings agreement/life insurance, car or other and unknown/no property

age

client's age

other_installment_plans

factor variable indicating other installment plans, with levels bank, stores and none

housing

factor variable indicating housing, with levels rent, own and for free

number_credits

number of existing credits at this bank

job

factor indicating employment status, with levels unemployed/unskilled - non-resident, unskilled - resident, skilled employee/official and management/self-employed/highly qualified employee/officer

people_liable

number of people being liable to provide maintenance

telephone

binary variable indicating if the customer has a registered telephone number

foreign_worker

binary variable indicating if the customer is a foreign worker

credit_risk

binary variable indicating credit risk, with levels good and bad

Correction examen
 - contrôle - continu
 12 octobre 2016
 НППА 304

Exercice 1

1) $\hat{\alpha} \in \text{arg min}_{\alpha} \sum_{i=1}^{10} \sum_{j=1}^3 (y_{ij} - \alpha \kappa_i)^2$

$\hat{\alpha} = (X^T X)^{-1} X^T y$

où $X = (\kappa_1, \kappa_1, \kappa_1, \dots, \kappa_{10}, \kappa_{10}, \kappa_{10})$

Ainsi, $X^T X = 3 \sum_{i=1}^{10} \kappa_i^2$

$(X^T X)^{-1} = \frac{1}{3 \sum_{i=1}^{10} \kappa_i^2}$ et $X^T y = \sum_{i=1}^{10} \sum_{j=1}^3 \kappa_i y_{ij}$
 $= \sum_{i=1}^{10} \kappa_i \sum_{j=1}^3 y_{ij}$

Nous obtenons

$$\hat{\alpha} = \frac{\sum_{i=1}^{10} \kappa_i \left(\frac{\sum_{j=1}^3 y_{ij}}{3} \right)}{\sum_{i=1}^{10} \kappa_i^2}$$

$$V(\hat{\alpha}) = \frac{1}{\left(\sum_{i=1}^{10} \kappa_i^2 \right)^2} \sum_{i=1}^{10} \kappa_i^2 \left(\frac{1}{9} \right) V\left(\sum_{j=1}^3 y_{ij} \right)$$

$$\begin{aligned} V\left(\sum_{j=1}^3 y_{ij} \right) &= C(y_{i1} + y_{i2} + y_{i3}, y_{i1} + y_{i2} + y_{i3}) \\ &= V(y_{i1}) + V(y_{i2}) + V(y_{i3}) + 2C(y_{i1}, y_{i2}) \\ &\quad + 2C(y_{i1}, y_{i3}) + 2C(y_{i2}, y_{i3}) \\ &= 3 + 3 \times 2 \times \frac{1}{9} = 6 \text{ puis on termine.} \end{aligned}$$

2] Soit $\tilde{\alpha}$ l'estimateur des moindres carrés généralisés de α

Nous avons

$$\tilde{\alpha} = (X^T \Omega^{-1} X)^{-1} X^T \Omega^{-1} y$$

$$\text{var } \Omega = V(y) = V(\varepsilon) = E(\varepsilon \varepsilon^T)$$

Prise at exemple, nous

(3)

avons

$$\forall (y_{i,1}, y_{i,2}, y_{i,3}) =$$

$$\begin{bmatrix} 1 & 0.5 & 0.5 \\ 0.5 & 1 & 0.5 \\ 0.5 & 0.5 & 1 \end{bmatrix}$$

$$= \frac{1}{2} I_3 + \frac{1}{2} J_3 \quad \text{avec}$$

I_3 matrice identité de dimension 3×3

J_3 matrice constante dont tous les

termes sont égaux à 1 de dimension 3×3

Nous avons donc

$$R = \begin{bmatrix} \frac{1}{2} I_3 + \frac{1}{2} J_3 & 0 & 0 & \dots \\ 0 & \frac{1}{2} I_3 + \frac{1}{2} J_3 & 0 & \\ 0 & 0 & \frac{1}{2} I_3 + \frac{1}{2} J_3 & 0 \\ \vdots & 0 & 0 & \ddots \end{bmatrix}$$

Nous pouvons maintenant
 écrire le code R, mais supposons
 que l'objet R y contient
 le vecteur (y₁₁, y₁₂, y₁₃, ..., y_{10,3})
 et l'objet R x contient
 le vecteur (x₁, x₂, x₃, ..., x₁₀, x₁₀, x₁₀)

Code R

```

H <- matrix(rep(1/2, 9), 3) + diag(rep(1/2), 3)
OMEGA <- bdiag(H, H, H, H, H, H, H, H, H, H)
# ( fonction bdiag de la
# bibliothèque Matrix
IOMEGA <- solve(OMEGA)
alphatilde <- solve(+ (x) %*% IOMEGA %*% y
                    %*% + (x) %*% IOMEGA %*% y
  
```


Code 2 si l'on ne connaît pas (5)
la fonction `diag` de
la bibliothèque `Matrix`

$H \leftarrow \text{matrice}(\text{rep}(1/2, 9), 3) + \text{diag}(\text{rep}(1/2), 3)$

$IH \leftarrow \text{inv}(H)$

$z \leftarrow \text{rep}(0, 30); w \leftarrow \text{rep}(0, 30)$

for (i in 1:10) {

$z[(1+3(i-1)):3i] \leftarrow IH \%* \% x[(1+3(i-1)):3i]$

$w[(1+3(i-1)):3i] \leftarrow IH \%* \% y[(1+3(i-1)):3i]$

}

$\text{alphabet} \leftarrow \text{inv}(+ (1) \%* \% z) \%* \% + (1) \%* \% w$

L'exécution du code ci-dessus vient
du fait que l'on a une matrice
bloc-diagonale et aussi bloc-diagonale
avec les matrices inverses des blocs en

Exercice 2

(6)

1)

$$f(y; \pi) = C_{y-1}^{\pi^{h-1} (1-\pi)^{y-h}} \pi^h (1-\pi)^{y-h} \pi^h (y) \quad \{h, h+1, \dots\}$$

$$f(y; \pi) = \exp \left\{ y \log(1-\pi) + h \log \left(\frac{\pi}{1-\pi} \right) \right\} \\ C_{y-1}^{\pi^{h-1} (1-\pi)^{y-h}} \pi^h (y) \quad \{h, h+1, \dots\}$$

Donc la loi négative binomiale appartient à une famille exponentielle avec paramètre de nuisance :

$$\alpha = \log(1-\pi) \Rightarrow \pi = 1 - e^{-\alpha}$$

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

$$\phi = 1$$

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

soit la mesure de référence

$$V(y) = C_{y-1}^{\pi^{h-1} (1-\pi)^{y-h}} \pi^h (y) \quad \{h, h+1, \dots\}$$

Donc toujours, nous savons que

$$E(y) = f'(a) = h + \frac{h e^a}{1 - e^a}$$

et donc $E(y) = \frac{h}{1 - e^a} = \frac{h}{\pi}$

D'après l'énoncé, $\log\left(\frac{\pi}{1 - \pi}\right) = \alpha + \beta x$

et donc $\pi = \left[\frac{e^{\alpha + \beta x}}{1 + e^{\alpha + \beta x}} \right]$

Ainsi, $E(y) = \frac{[1 + e^{\alpha + \beta x}] h}{e^{\alpha + \beta x}}$

et $E(y) = \pi(\alpha + \beta x)$

avec $\pi(x) = h \left(\frac{1 + e^q}{e^q} \right)$

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

Donc toujours,

(8)

$$\pi(q) = h \left[\frac{1+e^q}{e^q} \right] \neq h'(q) = \frac{h}{1-e^q}$$

ce n'est pas le lim comme d'habitude
qui a été utilisé -

$$\begin{aligned} 2) \quad L(\alpha, \beta) &= \sum_{i=1}^m y_i \log \left(\frac{1}{1+e^{\alpha+\beta x_i}} \right) \\ &+ h \sum_{i=1}^m [\alpha + \beta x_i] \\ &+ c \end{aligned}$$

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

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

On ne peut pas remarquer

$$\frac{dL}{d\alpha}(\alpha^*, \beta^*) = 0 \quad \& \quad \frac{dL}{d\beta}(\alpha^*, \beta^*) = 0$$

3]

(9)

$$\frac{d^2 \mathcal{L}}{d\alpha^2}(\alpha, \beta) = - \sum_{i=1}^M \gamma_i \left[\frac{e^{\alpha + \beta \mu_i} (1 + e^{\alpha + \beta \mu_i}) - e^{\alpha + \beta \mu_i} e^{\alpha + \beta \mu_i}}{(1 + e^{\alpha + \beta \mu_i})^2} \right]$$

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

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

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

$$\mathcal{I}_M(\alpha, \beta) = h \left[\begin{array}{c|c} \sum_{i=1}^M \frac{1}{(1 + e^{\alpha + \beta \mu_i})} & \sum_{i=1}^M \frac{\mu_i}{(1 + e^{\alpha + \beta \mu_i})} \\ \hline \sum_{i=1}^M \frac{\mu_i}{(1 + e^{\alpha + \beta \mu_i})} & \sum_{i=1}^M \frac{\mu_i^2}{(1 + e^{\alpha + \beta \mu_i})} \end{array} \right]$$

4) Lemme liim - arnorniqum

(10)

$$f(y) = \frac{h}{\pi} = \left[\frac{h}{1 - e^{\alpha + \beta x}} \right]$$

Dans ce cas, $\pi = 1 - e^{\alpha + \beta x}$

$$L(\alpha, \beta) = \sum_{i=1}^m y_i [\alpha + \beta x_i]$$

$$+ h \sum_{i=1}^m \log \left[\frac{1 - e^{\alpha + \beta x_i}}{e^{\alpha + \beta x_i}} \right]$$

+ cst

$$= \alpha \sum_{i=1}^m y_i + \beta \sum_{i=1}^m y_i x_i - h m \alpha$$
$$- h \beta \sum_{i=1}^m x_i + h \sum_{i=1}^m \log (1 - e^{\alpha + \beta x_i})$$

+ cst

$$\frac{dL}{d\alpha}(\alpha, \beta) = \sum y_i - m h$$

$$- h \frac{\sum_{i=1}^m e^{\alpha + \beta \pi_i}}{1 - e^{\alpha + \beta \pi_i}}$$

(11)

$$\frac{dL}{d\beta}(\alpha, \beta) = \sum_{i=1}^m y_i \pi_i - h \sum_{i=1}^m \pi_i$$

$$- h \sum_{i=1}^m \pi_i \left(\frac{e^{\alpha + \beta \pi_i}}{1 - e^{\alpha + \beta \pi_i}} \right)$$

Exercice 3

1] Il s'agit de six modèles de régression logistique

modèle 1: le risque de crédit expliqué par le salaire.

modèle 2: le risque de crédit expliqué par toutes les variables

modèle 3: le risque de crédit expliqué par les variables les plus pertinentes ou sans critère AIC et méthode des combinaisons

model 4 : idem que model 3
pour le critère BIC

model 5 : le risque de -crédit
expliqué uniquement par la
-constante (model null)

model 6 : le risque de -crédit
expliqué par les variables
les plus pertinentes en sens
du critère AIC et méthode
progressive.

⊆ Nous donnons le code pour le
model 1, il s'applique à tous
les autres en changeant
model 1 par model 1.

library(-caret)

twinn(formula(model 1), data = GermanCredit,
method = "glm", family = "binomial",
metric = "Accuracy", control = twinnControl(repeats = 10),
method = "repeatedcv", number = 2)