

Examen de contrôle continu

Durée 1h - Sans document - 15 octobre 2018

Exercice 1 (5 pts)

1 (2 pts) Quelle est la différence entre les modèles logit et probit ?

2 (3 pts) Quelle est l'objectif des méthodes de régularisation en regression ?

Exercice 2 (3 pts)

On considère le modèle de régression suivant, pour tout $i = 1, \dots, n$

$$y_i = \beta_1 + \beta_2 x_i + u_i$$

avec $\mathbb{E}(u_i) = 0$, $\mathbb{V}(u_i) = 1$ si i est pair et $\mathbb{V}(u_i) = 2$ si i est impair, $\mathbb{C}(u_i, u_j) = 0$ si $i \neq j$.
Donner l'expression de l'estimateur des moindres carrés généralisés de β_1 , commenter.

Exercice 3 (2 pts)

On considère n variables aléatoires indépendantes y_1, \dots, y_n telles que y_i est distribuée suivant une loi de Poisson de paramètre λ_i .

1 (0.5 pt) Montrer que la loi de Poisson appartient à la famille exponentielle.

2 (0.5 pt) Calculer $\mathbb{E}(y_i)$ et $\mathbb{V}(y_i)$.

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

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

Exercice 4 (10 pts)

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

```
library(glmnet)
library(caret)

set.seed(875)
n <- 1000
p <- 500
pm <- 50
x <- matrix(rnorm(n*p), nrow=n, ncol=p)
colnames(x) <- paste("v", 1:p, sep="")
y <- apply(x[,1:pm], 1, sum) + rnorm(n)

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

y.train <- y[train_rows]
y.test <- y[-train_rows]

model1 <- glmnet(x.train, y.train, family="gaussian", alpha=1, nlambda=50)
model2 <- train(x.train, y.train ,method="glmnet", metric="RMSE", trControl=
trainControl(method="repeatedcv", number=5, repeats=10),
tuneGrid=data.frame(alpha=1, lambda=model1$lambda))
model3 <- glmnet(x.train, y.train, family="gaussian", alpha=0)
model4 <- train(x.train, y.train, method="glmnet", metric="RMSE", trControl=
trainControl(method="repeatedcv", number=5, repeats=10),
tuneGrid=data.frame(alpha=0, lambda=model3$lambda))

dtrain <- data.frame(y=y.train, x.train)
model5 <- lm(y~., data=dtrain)
dtest <- data.frame(x.test)

mean((predict(model2, x.test) - y.test)^2)
mean((predict(model4, x.test) - y.test)^2)
mean((predict(model5, dtest) - y.test)^2)
mean((apply(x.test[,1:pm], 1, sum)-y.test)^2)

## [1] 1.475314
## [1] 3.141372
## [1] 3.729522
## [1] 1.032703
```

Correction examen
contrôle continu HPGA 304
15/10/2018

1

Exercice 1

Il s'agit de deux modèles
linéaires généralisés associés
à la loi de Bernoulli,
la variable à expliquer est
binomiale. Ces deux modèles
ont des fonctions de lien
différentes. Soit $p = \mathbb{E}(y)$
 $= \mathbb{P}(y=1)$
Pour le modèle logit, on suppose que
$$p = \frac{e^{\kappa\beta}}{1 + e^{\kappa\beta}}$$

Com le modèle probit, on suppose que $P = \Phi(X/\sigma)$.

(2)

2] Les méthodes de régularisation sont utiles lorsque le nombre de régressions (variables explicatives) est proche du nombre d'observations.
Objectif: minimiser la variance des estimateurs des paramètres.

Exercice 2

$$\hat{\beta}_1^{GLS} = \underset{\beta_1 \in \mathbb{R}}{\text{arg min}} \sum_{i=1}^n \left[\frac{(y_i - \beta_1 x_i)^2}{\alpha_i} \right]$$

avec $\alpha_i = 1$ si i est pair et
 $\alpha_i = 2$ si i est impair.

(3)

Somit $h(\beta_1, \beta_2) = \sum_{i=1}^m \left[\frac{(y_i - \beta_1 - \beta_2 x_i)^2}{d_i} \right]$

Neues Problem

$$\frac{dh}{d\beta_1}(\beta_1, \beta_2) = -2 \left[\sum_{i=1}^m \left(\frac{y_i}{d_i} \right) - \beta_1 \sum_{i=1}^m \left(\frac{1}{d_i} \right) - \beta_2 \sum_{i=1}^m \frac{x_i}{d_i} \right]$$

$$\frac{dh}{d\beta_2}(\beta_1, \beta_2) = -2 \left[\sum_{i=1}^m \frac{y_i x_i}{d_i} - \beta_1 \sum_{i=1}^m \frac{x_i}{d_i} - \beta_2 \sum_{i=1}^m \frac{x_i^2}{d_i} \right]$$

Aimmi, $\frac{dh}{d\beta_1}(\beta_1^*, \beta_2^*) = 0$ & $\frac{dh}{d\beta_2}(\beta_1^*, \beta_2^*) = 0$

$$\Leftrightarrow \left\{ \beta_1^* \sum \left(\frac{1}{d_i} \right) + \beta_2^* \sum \frac{x_i}{d_i} = \sum \frac{y_i}{d_i} \right.$$

$$\left. \beta_1^* \sum \frac{x_i}{d_i} + \beta_2^* \sum \frac{x_i^2}{d_i} = \sum \frac{y_i x_i}{d_i} \right\}$$

$$\Leftrightarrow \left\{ \begin{aligned} \beta_1^* &= \frac{\sum \frac{y_i}{d_i} - \beta_2^* \sum \frac{x_i}{d_i}}{\sum \frac{1}{d_i}} \\ \beta_2^* &= \frac{\sum \frac{y_i x_i}{d_i} - \beta_1^* \sum \frac{x_i}{d_i}}{\sum \frac{x_i^2}{d_i}} \end{aligned} \right. \quad (4)$$

$$\Rightarrow \sum \frac{1}{d_i} \beta_1^* = \sum \frac{y_i}{d_i} - \frac{\left[\sum \frac{y_i x_i}{d_i} - \beta_1^* \sum \frac{x_i}{d_i} \right] \sum \frac{x_i}{d_i}}{\sum \frac{x_i^2}{d_i}}$$

$$\Rightarrow \beta_1^* \left[\sum \frac{1}{d_i} \sum \frac{x_i^2}{d_i} - \left(\sum \frac{x_i}{d_i} \right)^2 \right] = \sum \frac{y_i}{d_i} \sum \frac{x_i^2}{d_i} - \left(\sum \frac{x_i}{d_i} \right)^2$$

$$\Rightarrow \beta_1^* = \frac{\sum \frac{y_i}{d_i} \sum \frac{x_i^2}{d_i} - \left(\sum \frac{x_i}{d_i} \right)^2}{\sum \frac{1}{d_i} \sum \frac{x_i^2}{d_i} - \left(\sum \frac{x_i}{d_i} \right)^2}$$

En effet,

$$\frac{d^2 h}{(d\beta_1)^2}(\beta_1, \beta_2) = 2 \sum \left(\frac{1}{\alpha_i} \right)$$

$$\frac{d^2 h}{(d\beta_2)^2}(\beta_1, \beta_2) = 2 \sum \frac{\mu_i^2}{\alpha_i}$$

$$\frac{d^2 h}{d\beta_1 d\beta_2}(\beta_1, \beta_2) = 2 \sum \frac{\mu_i}{\alpha_i}$$

Comme $2 \sum \left(\frac{1}{\alpha_i} \right) > 0$ et

$$4 \sum \left(\frac{1}{\alpha_i} \right) \sum \frac{\mu_i^2}{\alpha_i} - 4 \left(\sum \frac{\mu_i}{\alpha_i} \right)^2 > 0$$

(en effet, soit $M_i = \frac{1}{\sqrt{\alpha_i}}$ et $N_i = \frac{\mu_i}{\sqrt{\alpha_i}}$
d'après Cauchy - Schwarz

$$\left(\sum M_i N_i \right)^2 \leq \sum M_i^2 \sum N_i^2$$

alors la matrice hessienne de h est
définie positive et donc h strictement
convexe.

$$\hat{\beta}_1^{OLS} = \beta_1^*$$

(6)

Exercice 3

$$1] f(y; \lambda) = e^{-\lambda} \frac{\lambda^y}{y!} \prod_{i=1}^n \pi_i(y)$$

$$= \exp\{-\lambda + y \log(\lambda)\} \frac{1}{y!} \prod_{i=1}^n \pi_i(y)$$

$$\theta = \log(\lambda) \Rightarrow \lambda = e^\theta$$

$$h(\theta) = e^\theta$$

$$v(\omega y) = \frac{1}{y!} \prod_{i=1}^n \pi_i(y) S(\omega y)$$

où S est la mesure de comptage

$$2] E(y) = h'(\theta) = e^\theta = \lambda$$

$$V(y) = h''(\theta) = e^\theta = \lambda$$

$$3] E(y) = \lambda = e^{\beta_1 + \beta_2 x}$$

(8)

$$\Rightarrow E(y) = f(\beta_1 + \beta_2 x)$$

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

son lien, $f(\mu) = e^\mu = b'(\mu)$

Ainsi, c'est bien la fonction de lien canonique qui a été utilisée.

Exercice 4

Le vrai modèle contient 50 régresseurs générés à partir d'une loi gaussienne centrée réduite.

Nous avons

$$y = x_1 + \dots + x_{50} + \mu$$

$$\text{où } \mu \sim N(0, 1)$$

(7)

La première partie du code consiste à générer des données suivant ce modèle.

En plus, 450 variables aléatoires sont ajoutées au modèle. On compare alors les résultats produits par des régressions ridge, lasso et classique sur l'ensemble variables prédictives. Les paramètres de régularisation des régressions lasso et ridge est fixé à l'aide d'une stratégie de validation croisée à 5 ensembles répétée 10 fois.

(8)