

Examen de rattrapage - 13 janvier 2016
Durée 1h30 - Documents interdits

Exercice 1 (6 pts) On se place dans le contexte de la classification supervisée binaire.

- 1) (2 pts) Décrire le modèle de la régression logistique.
- 2) (2 pts) Donner trois critères de sélection entre deux modèles de régression logistique.
- 3) (2 pts) Donner deux stratégies de sélection de variables basées sur les critères précédents.

Exercice 2 (6 pts) Nous considérons un jeu de données contenu dans le data.frame spam de la bibliothèque kernlab.

A data frame with 4601 observations and 58 variables.

The first 48 variables contain the frequency of the variable name (e.g., business) in the e-mail. If the variable name starts with num (e.g., num650) the it indicates the frequency of the corresponding number (e.g., 650). The variables 49-54 indicate the frequency of the characters ';', '(', '[' , '!', '\\$', and '\#'. The variables 55-57 contain the average, longest and total run-length of capital letters. Variable 58 indicates the type of the mail and is either ‘nonspam’ or ‘spam’, i.e. unsolicited commercial e-mail.

The data set contains 2788 e-mails classified as ‘nonspam’ and 1813 classified as ‘spam’.

Donner le code R comparant les classifieurs produits par les méthodes suivantes :

- classification par regression,
- analyse discriminante quadratique,
- forêts aléatoires.

Exercice 3 (8 pts) Commenter en détails le fichier R Markdown fourni en annexes.

1

ΗΠΠΑ 303
Discrimination et
scoring - correction
examen notrophe
13/01/2016

Exercice 1

1] $Y \in \{0, 1\}$ et $X \in \mathbb{R}^d$

$$P(Y=1 | X=x) = \frac{e^{x^T \beta}}{1 + e^{x^T \beta}}$$

où $\beta \in \mathbb{R}^d$ est un paramètre
inconnu.

2] Le modèle de régression est \textcircled{L}
utilisé dans le contexte de
la classification supervisée
linéaire. On peut donc
utiliser le critère de
l'erreur de classification,
erreur que l'on pourra
estimer par une technique
de validation croisée.

Deux autres critères sont
basés sur des notions probabilistes
de la fonction de vraisemblance,
il s'agit des critères AIC et
BIC.

3] On peut utiliser une
méthode de sélection
ascendante ou une stratégie
descendante. (3)

Exercice 2

```
library(kernlab); library(MASS)
data(spam); library(randomForest)
model.rf <- randomForest(type="n", data=spam)
```

```
model.rf$var.m.rtree[500,1]
```

```
K <- 10; M <- 4600
```

```
g <- M/K; res.hola <- (0,M)
```

```
res.reg <- rep(0,M); ind <- sample(M)
```

```
for (i in 1:K)
```

```
{
```

```

ent ← imobli [((K-1)*g+1):(K*g)]
model. lolo ← lolo (type n., data = spom
res. lolo [imobli] ← predict (model. lolo,
                             spom [imobli,]) $ posterior [2]
model. rey ← lm (as.numeric (type) - 1 n.,
                data = spom [-imobli,])
res. rey [imobli] ← predict (model. rey,
                             spom [imobli,])
}

```

```

meom (Gamma (res. lolo) ! = (as.numeric (type) - 1))
MEMM (as.numeric (res. rey) >= 1/2) ! = (as.numeric (type) - 1))

```

Dans le code présenté ci-dessus, nous avons calculé une seule fois l'erreur par validation croisée à $K=10$ ensemble.

Exercice 3

5

On a environ 500 observations de
notre échantillon pour les utiliser
comme données de validation.

Avec ces données nous pourrions
évaluer les performances de
différentes méthodes de classification.

On le fait par LDA, les forêts
aléatoires et la classification bayésienne
naïve. Les forêts aléatoires donnent
les meilleurs résultats.

Puis, on évalue par validation
croisée à 10 - ensemble le nombre
de voisins optimal pour la méthode

als k -plus-proches voisins. ⑥
On trouve $k = 1$ -