

Examen partiel - 4 novembre 2015
Durée 1h30 - Documents interdits

Exercice 1 (7 pts) Soient X le vecteur aléatoire de dimension d contenant des variables explicatives quantitatives et Y le facteur binaire à valeurs dans $\{0, 1\}$ à expliquer.

1) (1 pt) Donner la probabilité que $Y = 1$ sachant $X = x$ en fonction de la loi marginale de Y et de la loi conditionnelle de X .

2) (3 pts) Pour la fonction de coût élémentaire

$$h(y, d) = \begin{cases} 0 & \text{si } y = d \\ a & \text{si } y = 1 \text{ et } d = 0 \\ 10a & \text{si } y = 0 \text{ et } d = 1 \end{cases}$$

donner le prédicteur $g(X)$ minimisant le coût moyen.

3) (3 pts) Expliquer en détails le résultat produit par le code R donné ci-dessous.

```
pi0=0.3;pi1=0.7
mu0=c(3,9);mu1=c(5,5)
Sigma0=matrix(c(3,3*0.5,3*0.5,3),2,2)
Sigma1=matrix(c(5,-5*0.9,-5*0.9,5),2,2)
library(mvtnorm)
print(c(pi0*dmvnorm(c(4,4),mu0,Sigma0),pi1*dmvnorm(c(4,4),mu1,Sigma1)))
```

Exercice 2 (6 pts) Expliquer en détails le résultat produit par le code R donné ci-dessous.

```
n=200;pi0=0.3;pi1=0.7;mu0=c(3,9);mu1=c(5,5)
Sigma0=matrix(c(3,3*0.5,3*0.5,3),2,2)
Sigma1=matrix(c(5,-5*0.9,-5*0.9,5),2,2)
y=sample(c(0,1),size=n,prob=c(pi0,pi1),replace=TRUE);y=sort(y)
```

```
n0=sum(y==0);n1=sum(y==1);library(mvtnorm)
x=rbind(rmvnorm(n0,mu0,Sigma0),rmvnorm(n1,mu1,Sigma1))
plot(x,pch="*",cex=4,col=y+1)
```

```
donnees=data.frame(y=y,x1=x[,1],x2=x[,2])
```

```
nrep=100;npart=n/5
errorreg=rep(0,nrep)
for (l in 1:nrep)
{
  indi=sample(1:n,n)
  for (i in 1:5)
  {
    out=indi[(npart*(i-1)+1):(i*npart)]
    model.lm=lm(y~x1+x2,data=donnees[-out,])
    classif[out]=as.numeric(predict(model.lm,donnees[out,])>1/2)
  }
  errorreg[l]=mean(classif!=y)
}
```

Exercice 3 (7 pts) On se place dans le contexte de la classification supervisée binaire.

1) (2 pts) Expliquer en quoi consiste la méthode des k -plus-proches-voisins ?

2) (2 pts) Quel est le modèle de l'analyse discriminante linéaire ?

3) (4 pts) Nous supposons que nous observons 100 réalisations indépendantes $(x_1, y_1), \dots, (x_n, y_n)$ du couple (X, Y) stockées dans les objets R : \mathbf{x} et \mathbf{y} , \mathbf{y} est le vecteur binaire à expliquer. Donner le code R mettant en oeuvre la technique de la validation croisée à 2 ensembles pour calculer l'erreur associée à l'analyse discriminante quadratique.

Correction examen
partiel 04/11/2015
HPPA 303

①

Exercice 1

$$1] \mathbb{P}(Y=1|X=\kappa) = \frac{f_1(\kappa)\pi_1}{f_0(\kappa)\pi_0 + f_1(\kappa)\pi_1}$$

avec $\mathbb{P}(Y=i) = \pi_i$, $\forall i=0,1$
et $f_i(\kappa)$ la densité de la loi
conditionnelle de X sachant $Y=i$,
 $\forall i=0,1$.

$$2] C(g) = \mathbb{E}[h(Y, g(X))] (*)$$
$$= \mathbb{E}\left[\mathbb{E}[h(Y, g(X)) | X=\kappa]\right]$$

(2)

$$\begin{aligned} & \mathbb{E}[h(Y, g(x)) | X=x] \\ &= \alpha \mathbb{P}(Y=1 | X=x) \mathbb{1}_{\{g(x)=0\}} \\ &+ (1-\alpha) \mathbb{P}(Y=0 | X=x) \mathbb{1}_{\{g(x)=1\}} \quad (**) \end{aligned}$$

Minimiser (*) en $g(\cdot)$ est équivalent à minimiser (**) en $g(\cdot)$ uniformément en x .

$$g^*(x) \in \arg \min_{g(\cdot)} \mathbb{E}[h(Y, g(x)) | X=x]$$
$$g^*(x) = \begin{cases} 0 & \text{si } \alpha \mathbb{P}(Y=1 | X=x) < \\ & (1-\alpha) \mathbb{P}(Y=0 | X=x) \\ 1 & \text{sinon} \end{cases}$$

$$g^*(k) = \begin{cases} 0 & \text{si } \mathbb{P}(Y=0|X=k) > \frac{1}{11} \\ 1 & \text{si } \mathbb{P}(Y=0|X=k) \leq \frac{1}{11} \end{cases} \quad (3)$$

On veut surtout pas se tromper
lorsque $Y=0$.

3] Nous avons

$$\pi_0 = \mathbb{P}(Y=0) = 0.3$$

$$\pi_1 = \mathbb{P}(Y=1) = 0.7$$

$$X|Y=0 \sim \mathcal{N}_2 \left(\begin{pmatrix} 3 \\ 9 \end{pmatrix}, \begin{bmatrix} 3 & 3/2 \\ 3/2 & 3 \end{bmatrix} \right)$$

$$X|Y=1 \sim \mathcal{N}_2 \left(\begin{pmatrix} 5 \\ 5 \end{pmatrix}, \begin{bmatrix} 5 & -\frac{45}{10} \\ -\frac{45}{10} & 5 \end{bmatrix} \right)$$

Le cas où on calcule le vecteur

$$(\pi_0 f_0((4,4)), \pi_1 f_1((4,4)))$$

Exercice 2

Simulation de $n=200$ données
suivant le modèle présenté exercice 1
question 3.

Calcul de l'erreur associée au
classification par régression linéaire
par une méthode de validation
croisée à 5 ensembles répétée
100 fois.

Exercice 3

- 1) `train`
- 2) `train`

3) `model = lda ← train(yN, XN,`
`data = data.frame(y = y, X = X),`
`method = "lda", metric = "Accuracy",`
`trainControl = trainControl(method = "repeated", number = 5,`
`repeats = 100))`