

Examen de contrôle continu
Durée 1h30 - Sans document - 09 octobre 2020

Exercice 1 (3 pts) On se place dans le contexte de la classification supervisée binaire. Qu'est qu'une courbe ROC? Que permet-elle de faire?

Exercice 2 (3 pts) On se place dans le contexte de la classification supervisée binaire. On suppose que l'on dispose d'une seule variable explicative réelle X . La variance empirique de X est égale 1, la moyenne des X pour $Y = 0$ est égale à 0 et la moyenne des X pour $Y = 1$ est égale à 2. Par ailleurs, pour les données observées, les deux classes sont équilibrées. Donner l'expression explicite du classifieur LDA.

Exercice 3 (8 pts) Soit Y un facteur binaire à expliquer, $Y \in \{0, 1\}$ et X une variable explicative réelle. Nous supposons que $\mathbb{P}(Y = 1) = 1/3$ et

$$X|Y = 1 \sim f_1(x) = 2 \exp(-2x) \mathbb{I}_{\mathbb{R}_+}(x)$$
$$X|Y = 0 \sim f_0(x) = \frac{1}{2} \exp(-x/2) \mathbb{I}_{\mathbb{R}_+}(x).$$

- 1) (0.5 pt) Donner la loi marginale de X .
- 2) (0.5 pt) Donner la loi de Y sachant $X = x$.
- 3) (2 pts) Donner le classifieur optimal g^* lorsque

$$h(y, d) = \begin{cases} 0 & \text{si } y = d \\ 2 & \text{si } y = 1 \text{ et } d = 0 \\ 1 & \text{si } y = 0 \text{ et } d = 1 \end{cases} .$$

- 4) (2 pts) Calculer le taux d'erreur moyen associé à g^* .
- 5) (3 pts) Donner deux classifieurs aléatoires différents et calculer leurs taux d'erreur moyens. Comparer tous les taux d'erreur.

Exercice 3 (6 pts) Expliquer en détails les résultats produits par le code R donné ci-dessous.

```
> cancer.train <- read.csv("~/Desktop/data-cancer-train.csv")
> cancer.train <- cancer.train[,-1]
> cancer.test <- read.csv("~/Desktop/data-cancer-test.csv")
> cancer.test <- cancer.test[,-1]
> cancer.train$diagnosis <- as.factor(cancer.train$diagnosis)
> cancer.test$diagnosis <- as.factor(cancer.test$diagnosis)
```

```

> library(MASS)
> model1 <- qda(diagnosis~.,data=cancer.train)
> library(caret)
> confusionMatrix(predict(model1)$class,cancer.train$diagnosis,positive="M")

```

Confusion Matrix and Statistics

	Reference	
Prediction	B	M
B	318	5
M	3	141

```

Accuracy : 0.9829
95% CI : (0.9665, 0.9926)
No Information Rate : 0.6874
P-Value [Acc > NIR] : <2e-16

```

```

Kappa : 0.96

```

```

Mcnemar's Test P-Value : 0.7237

```

```

Sensitivity : 0.9658
Specificity : 0.9907
Pos Pred Value : 0.9792
Neg Pred Value : 0.9845
Prevalence : 0.3126
Detection Rate : 0.3019
Detection Prevalence : 0.3084
Balanced Accuracy : 0.9782

```

```

'Positive' Class : M

```

```

> model2 <-train(diagnosis~.,data=cancer.train,method="qda",metric="Accuracy",
  trControl=trainControl(method="repeatedcv",number=10,repeats=100))
> confusionMatrix(model2)

```

Cross-Validated (10 fold, repeated 100 times) Confusion Matrix

(entries are percentual average cell counts across resamples)

	Reference	
Prediction	B	M
B	66.6	1.5
M	2.2	29.7

```

Accuracy (average) : 0.963

```

```

> dim(cancer.test)
[1] 102 31

> sum(predict(model1,cancer.test)$class==
      predict(model2,cancer.test,type="raw"))
[1] 102

> confusionMatrix(predict(model1,cancer.test)$class,
                  cancer.test$diagnosis,positive="M")

```

Confusion Matrix and Statistics

```

          Reference
Prediction B  M
   B  35  5
   M   1 61

          Accuracy : 0.9412
          95% CI   : (0.8764, 0.9781)
No Information Rate : 0.6471
P-Value [Acc > NIR] : 2.076e-12

          Kappa   : 0.8744

Mcnemar's Test P-Value : 0.2207

          Sensitivity : 0.9242
          Specificity : 0.9722
   Pos Pred Value   : 0.9839
   Neg Pred Value   : 0.8750
          Prevalence : 0.6471
   Detection Rate   : 0.5980
Detection Prevalence : 0.6078
   Balanced Accuracy : 0.9482

'Positive' Class : M

```

(1)

Correction examen
de contrôle continu
HPPA 303
09/10/2020

Exercice 1

Voir cours

Exercice 2

$(y_1, x_1), \dots, (y_m, x_m)$

$$M_0 = \sum_{i=1}^m \prod_{\{y_i=0\}}$$

$$M_1 = \sum_{i=1}^m \prod_{\{y_i=1\}}$$

(2)

$$\hat{\mu}_0 = \frac{1}{M_0} \sum_{i=1}^M \mathbb{1}_{\{Y_i=0\}} X_i$$

$$\hat{\mu}_1 = \frac{1}{M_1} \sum_{i=1}^M \mathbb{1}_{\{Y_i=1\}} X_i$$

$$\hat{M}_0 = 0, \hat{M}_1 = 2, M_0 = M_1$$

$$\hat{\pi}_0 = \hat{\pi}_1 = \frac{1}{2}$$

$$\text{Einfim, } \hat{\sigma}^2 = \left[\frac{1}{M-2} \right] \sum_{i=1}^M \mathbb{1}_{\{Y_i=j\}} (X_i - \hat{\mu}_j)^2$$

$$\hat{\sigma}^2 = 1$$

$$\hat{\mathbb{P}}(Y=1|X=1) = \frac{e^{-\frac{1}{2}(1-2)^2}}{e^{-\frac{1}{2}1^2} + e^{-\frac{1}{2}(1-2)^2}}$$

Le classifieur LDA, en l'absence
de fonction de perte est tel
que

$$\hat{g}^{LDA}(x) = \begin{cases} 1 & \text{if } \hat{\mathbb{P}}(Y=1|X=x) \geq \frac{1}{2} \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

$$\hat{\mathbb{P}}(Y=1|X=x) \geq \frac{1}{2}$$

$$\Leftrightarrow e^{-\frac{1}{2}x^2} \leq e^{-\frac{1}{2}(x-2)^2}$$

$$\Leftrightarrow x^2 \geq (x-2)^2$$

$$\Leftrightarrow -2x + 4 \geq -4x + 4$$

$$\Leftrightarrow 4x \geq 4$$

$$\Leftrightarrow x \geq 1$$

$$\hat{g}^{LDA}(x) = \begin{cases} 1 & \text{if } x \geq 1 \\ 0 & \text{if } x < 1 \end{cases}$$

Exercice 3

(4)

$$1) f_X(x) = \begin{cases} \frac{2}{3} \left(\frac{1}{2} e^{-x/2} \right) + \frac{1}{3} (2e^{-2x}) & x \geq 0 \\ 0 & x < 0 \end{cases}$$

$$2) P(Y=1 | X=1) = \frac{\frac{2}{3} e^{-2}}{\frac{1}{3} e^{-1/2} + \frac{2}{3} e^{-2}}$$

$$\text{si } x \geq 0$$

$$3) g^*(x) = \begin{cases} 1 & \text{si } P(Y=1 | X=x) \geq \frac{1}{3} \\ 0 & \text{sinon} \end{cases}$$

$$P(Y=1 | X=x) \geq \frac{1}{3}$$

$$\Leftrightarrow \frac{4}{9} e^{-2x} \geq \frac{1}{9} e^{-x/2}$$

$$\Leftrightarrow 4e^{-2x} \geq e^{-x/2}$$

$$\Leftrightarrow -\frac{x}{2} \leq -2x + \ln(4)$$

(5)

$$\Leftrightarrow \frac{3\kappa}{6} \leq \log(k)$$

$$\Leftrightarrow \kappa \leq \frac{2 \log(k)}{3}$$

$$g^*(11) = \begin{cases} 1 & \text{si } \kappa \leq \frac{2 \log(k)}{3} \\ 0 & \text{sinon} \end{cases}$$

4) Tous les ordres moyens
≠ sont moyens

$$\mathbb{P}(Y \neq g^*(x))$$

$$= \mathbb{P}(\{Y=0\} \cap \{g^*(x)=1\}) \\ + \mathbb{P}(\{Y=1\} \cap \{g^*(x)=0\})$$

$$= \frac{2}{3} \mathbb{P}(g^*(x)=1 | Y=0) \\ + \frac{1}{3} \mathbb{P}(g^*(x)=0 | Y=1)$$

(6)

$$P(Y \neq g^*(X))$$

$$= \frac{2}{3} P\left(X \leq \frac{2 \log(4)}{3} \mid Y=0\right)$$

$$+ \frac{1}{3} P\left(X > \frac{2 \log(4)}{3} \mid Y=1\right)$$

$$= \frac{2}{3} \int_0^{\frac{2 \log(4)}{3}} \frac{1}{2} e^{-x/2} dx$$

$$+ \frac{1}{3} \int_{\frac{2 \log(4)}{3}}^{\infty} 2 e^{-2x} dx$$

$$= \frac{2}{3} \left[1 - e^{-\log(4)/3} \right]$$

$$+ \frac{1}{3} \left[e^{-\frac{4 \log(4)}{3}} \right]$$

$$= \frac{2}{3} - \frac{2(4)^{-1/3}}{3} + \frac{(4)^{-4/3}}{3} \quad \# 0.499$$

5

7

$$g^1(\pi) = \prod \{U \leq \frac{1}{2}\}$$

$$g^2(\pi) = \prod \{U \leq \frac{1}{3}\}$$

$$\mathbb{P}(Y \neq g^1(x))$$

$$= \frac{2}{3} \times \frac{1}{2} + \frac{1}{3} \times \frac{1}{2}$$

$$= \frac{1}{2}$$

$$\mathbb{P}(Y \neq g^2(x)) = \binom{2}{3} \mathbb{P}(U \leq \frac{1}{3} | Y=0)$$

$$+ \binom{1}{3} \mathbb{P}(U \geq \frac{1}{3} | Y=1)$$

$$= \binom{2}{3} \binom{1}{3} + \binom{1}{3} \binom{2}{3}$$

$$= \frac{4}{9} \quad \# \quad 0.444$$

Exercice 3

(2)

Prise en compte d'une méthode LDA sur les données de training

Comparaison de trois matrices de confusion :

- la première sur les données de training Accuracy: 97,29%

- la deuxième estimée par validation croisée à 10 ensemble répétée 100 fois

Accuracy: 96,3%

- la troisième sur les données de test Accuracy: 94,12%

Dis-ent les intervalles de confiance.