

**Examen final - 9 novembre 2016**  
**Durée 2h00 - Documents interdits**

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

1) (2 pts) Expliquer les étapes de construction d'un arbre de classification (CART).

2) (2 pts) Expliquer les étapes de construction d'une forêt aléatoire (random forest).

3) (2 pts) Supposons que le jeu de données  $\mathbf{d}$  contient le facteur binaire  $y$  et un ensemble de variables explicatives. Donner le code R permettant de comparer la qualité des classifieurs CART et random forest.

**Exercice 2 (8 pts)** Soient  $Y$  un facteur binaire à expliquer  $Y \in \{0, 1\}$ ,  $X_1$  un facteur explicatif prenant les modalités  $\mathbf{a}$  et  $\mathbf{b}$  et  $X_2$  un deuxième facteur explicatif prenant les modalités  $\mathbf{c}$  et  $\mathbf{d}$ . Nous disposons d'un échantillon de 100 observations dont les résultats sont synthétisés dans les deux tableaux ci-dessous

— lorsque $y = 0$		$x_1 = a$	$x_1 = b$
	$x_2 = c$	15	5
	$x_2 = d$	5	25

— lorsque $y = 1$		$x_1 = a$	$x_1 = b$
	$x_2 = c$	0	30
	$x_2 = d$	10	10

1) (2 pts) Donner le classifieur obtenu à l'aide de la méthode de l'analyse discriminante non paramétrique.

2) (2 pts) Proposer un classifieur alternatif dont la construction n'est pas basée sur une hypothèse d'indépendance conditionnelle.

3) (4 pts) Nous supposons que

$$\mathbb{P}(Y = 0) = 1/2, \mathbb{P}((X_1, X_2) = (i, j) | Y = 0) = \mathbb{P}(X_1 = i | Y = 0)\mathbb{P}(X_2 = j | Y = 0),$$

$$\mathbb{P}(X_1 = a | Y = 0) = \mathbb{P}(X_2 = c | Y = 0) = 1/3 \text{ et enfin que}$$

$$\mathbb{P}((X_1, X_2) = (i, j) | Y = 1) = \begin{cases} 0 & (i, j) = (a, c) \\ 2/3 & (i, j) = (b, c) \\ 1/6 & (i, j) = (a, d) \\ 1/6 & (i, j) = (b, d) \end{cases} .$$

Pour la fonction de coût élémentaire  $h(y, d) = \mathbf{1}_{y \neq d}$ , donner le coût moyen des deux classifieurs mis en évidence questions 1) et 2).

**Exercice 3 (6 pts)** Commenter en détails le fichier R Markdown fourni en annexes. Les données sont celles du challenge kaggle sur le naufrage du Titanic.

①

Correction examen  
final 09/11/2016  
HPPA303

Exercice 1

Voir le cours

Exercice 2

1) La technique de l'analyse discriminante paramétrique consiste à estimer les lois

$$(X_1, X_2) | Y=0 \text{ et } (X_1, X_2) | Y=1$$

sous les hypothèses  $X_1 \perp X_2 | Y=0$

$$X_1 \perp X_2 | Y=1$$

D'après la table de contingence pour  $Y=0$ , nous avons

$$\hat{P}(X_1=a | Y=0) = \frac{2}{5} \text{ et } \hat{P}(X_2=c | Y=0) = \frac{2}{5}$$

News obtenons,

$$\widehat{P}((X_1, X_2) = (i, j) | Y = 0) = \begin{cases} \frac{4}{25} & (i, j) = (a, c) \\ \frac{6}{25} & (i, j) = (b, c) \\ \frac{6}{25} & (i, j) = (a, d) \\ \frac{9}{25} & (i, j) = (b, d) \end{cases}$$

De même,

$$\widehat{P}(X_1 = a | Y = 1) = \frac{1}{5} \quad \text{et} \quad \widehat{P}(X_2 = c | Y = 1) = \frac{3}{5}$$

$$\widehat{P}((X_1, X_2) = (i, j) | Y = 1) = \begin{cases} \frac{3}{25} & (i, j) = (a, c) \\ \frac{12}{25} & (i, j) = (b, c) \\ \frac{2}{25} & (i, j) = (a, d) \\ \frac{8}{25} & (i, j) = (b, d) \end{cases}$$

En velleurs,

$$\widehat{P}(Y = 0) = \frac{1}{2} = \widehat{P}(Y = 1)$$

News of them,

(3)

$$\widehat{P}(Y=0 | (X_1, X_2) = (u, c))$$

$$= \frac{\widehat{P}(Y=0) \widehat{P}((X_1, X_2) = (u, c) | Y=0)}{\widehat{P}(Y=0) \widehat{P}((X_1, X_2) = (u, c) | Y=0) + \widehat{P}(Y=1) \widehat{P}((X_1, X_2) = (u, c) | Y=1)}$$
$$= \frac{\frac{4}{25}}{\frac{4}{25} + \frac{3}{25}} = \frac{4}{7}$$

De même,

$$\widehat{P}(Y=0 | (X_1, X_2) = (t, c)) = \frac{\frac{6}{25}}{\frac{6}{25} + \frac{12}{25}} = \frac{1}{3}$$

$$\widehat{P}(Y=0 | (X_1, X_2) = (u, d)) = \frac{\frac{6}{25}}{\frac{6}{25} + \frac{2}{25}} = \frac{3}{4}$$

$$\widehat{P}(Y=0 | (X_1, X_2) = (t, d)) = \frac{\frac{9}{25}}{\frac{9}{25} + \frac{8}{25}} = \frac{9}{17}$$

On m'écrit la classification

$$g_2(N_1, N_2) = \begin{cases} 0 & (N_1, N_2) = (u, c) \\ 1 & (N_1, N_2) = (t, c) \\ 0 & (N_1, N_2) = (u, d) \\ 0 & (N_1, N_2) = (t, d) \end{cases}$$

2] On peut estimer les

(4)

lois  $(X_1, X_2) | Y=0$  et

$(X_1, X_2) | Y=1$  sans utiliser  
l'hypothèse d'indépendance  
conditionnelle.

$$\widehat{\mathbb{P}}((X_1, X_2) = (a, c) | Y=0) = \frac{15}{50} = \frac{3}{10}$$

$$\widehat{\mathbb{P}}((X_1, X_2) = (b, c) | Y=0) = \frac{5}{50} = \frac{1}{10}$$

$$\widehat{\mathbb{P}}((X_1, X_2) = (a, d) | Y=0) = \frac{5}{50} = \frac{1}{10}$$

$$\widehat{\mathbb{P}}((X_1, X_2) = (b, d) | Y=0) = \frac{25}{50} = \frac{1}{2}$$

De même,

$$\widehat{\mathbb{P}}((X_1, X_2) = (i, j) | Y=1) = \begin{cases} 0 & (i, j) = (a, c) \\ \frac{3}{5} & (i, j) = (b, c) \\ \frac{1}{5} & (i, j) = (a, d) \\ \frac{1}{5} & (i, j) = (b, d) \end{cases}$$

Comme, pour toutes  $n$ ,

$$\widehat{\mathbb{P}}(Y=1) = \widehat{\mathbb{P}}(Y=0) = \frac{1}{2}$$

mess of errors

(5)

$$\widehat{\mathbb{P}}(Y=0 | (X_1, X_2) = (i, j)) = \begin{cases} 1 & (i, j) = (u, c) \\ \frac{1}{7} & (i, j) = (k, c) \\ \frac{1}{3} & (i, j) = (u, d) \\ \frac{5}{7} & (i, j) = (k, d) \end{cases}$$

Im in absolut & definition

$$g_2(\pi_1, \pi_2) = \begin{cases} 0 & (\pi_1, \pi_2) = (u, c) \\ 1 & (\pi_1, \pi_2) = (k, c) \\ 1 & (\pi_1, \pi_2) = (u, d) \\ 0 & (\pi_1, \pi_2) = (k, d) \end{cases}$$

$$3] C(g_1) = \mathbb{E}[h(Y, g_1(X))] \quad X = (X_1, X_2)$$

$$C(g_1) = \mathbb{P}(Y \neq g_1(X))$$

$$C(g_1) = \mathbb{P}(\{Y=0\} \cap \{g_1(X)=1\}) \\ + \mathbb{P}(\{Y=1\} \cap \{g_1(X)=0\})$$

$$C(g_1) = \mathbb{P}(\{Y=0\} \cap \{(X_1, X_2) = (b, c)\}) \quad (6)$$

$$+ \mathbb{P}(\{Y=1\} \cap \{(X_1, X_2) = (a, c) \cup (X_1, X_2) = (a, d) \cup (X_1, X_2) = (b, d)\})$$

$$C(g_1) = \mathbb{P}((X_1, X_2) = (b, c) | Y=0) \mathbb{P}(Y=0)$$

$$+ \mathbb{P}((X_1, X_2) = (a, c) | Y=1) \mathbb{P}(Y=1)$$

$$+ \mathbb{P}((X_1, X_2) = (a, d) | Y=1) \mathbb{P}(Y=1)$$

$$+ \mathbb{P}((X_1, X_2) = (b, d) | Y=1) \mathbb{P}(Y=1)$$

$$= \frac{2}{9} \times \frac{1}{2} + \frac{1}{2} \left[ 0 + \frac{1}{6} + \frac{1}{6} \right]$$

$$= \frac{1}{2} \left[ \frac{2}{9} + \frac{2}{6} \right] = \left( \frac{15}{54} \right)$$

Di min

$$C(g_2) = \mathbb{P}(\{Y=0\} \cap \{(X_1, X_2) = (b, c) \cup (X_1, X_2) = (a, d)\})$$

$$+ \mathbb{P}(\{Y=1\} \cap \{(X_1, X_2) = (a, c) \cup (X_1, X_2) = (b, d)\})$$

7

$$C(g_2) = \frac{1}{2} \left[ \frac{2}{9} + \frac{1}{4} \right] + \frac{1}{2} \left[ 0 + \frac{1}{6} \right]$$

$$\begin{aligned} C(g_2) &= \left[ \frac{1}{9} + \frac{1}{8} + \frac{1}{12} \right] \\ &= \left[ \frac{8}{72} + \frac{9}{72} + \frac{6}{72} \right] = \frac{23}{72} \end{aligned}$$

Au final,  $C(g_1) < C(g_2)$

# 0,28      # 0,32

### Exercice 3

Dans une première partie, les données sont pré-traitées avec notamment la création de la variable "title" qui contient le titre des passagers. Cette information est extraite de leurs noms.



Dans un second temps, les données manquantes imputées.

Pour l'âge, il est utilisé l'âge médian correspondant

au titre, le calcul sur la

base de données fusionnée (apprentissage

et test) pour la classe

de titre de l'individu pour

lequel l'âge est manquant.

Dans une dernière partie, un modèle de classification par arbre dont le paramètre d'élagage est fixé par validation croisée est mis en œuvre.

Enfin, un modèle de forêt aléatoire est mis en œuvre -