

QCM

Apprentissage supervisé

Test
CC du 27/04/2022

Instructions :

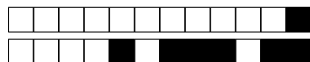
- Le sujet comprend 4 exercices pour 22 questions au total. Les questions faisant apparaître le symbole ♣ peuvent présenter plusieurs bonnes réponses. Les autres ont une unique bonne réponse. Des points négatifs pourront être affectés à de *mauvaises* réponses.
- Seul le questionnaire à la 7ème page est à rendre. Vous commencerez par renseigner votre nom et prénom dans la case prévue ainsi que votre **numéro étudiant Rennes 2** (1 case à colorier par colonne). Un numéro d'étudiant mal renseigné entraînera la note de 0.
- Il faut **colorier** les cases correspondants aux bonnes réponses (sur la page 7), mettre une croix dans la case n'est **pas suffisant**. Les cases devront être **coloriées avec un stylo noir** (pas de crayon papier, de stabilo...).
- Le barème sera effectué de la façon suivante :
 - Aucune case coloriée entraînera une note de 0 sur la question.
 - Pour les questions à une seule bonne réponse (sans le symbole ♣), un nombre de points sera affecté (par exemple +2) si la bonne case est cochée. Un nombre de points sera retranché (par exemple -1) si une mauvaise case est coloriée ou si plusieurs cases sont coloriées.
 - Pour les questions avec plusieurs bonnes réponses (avec le symbole ♣), un nombre de points (par exemple +0.5) sera affecté pour chaque bonne réponse coloriée et pour chaque mauvaise réponse non coloriée. Un nombre de points (par exemple -0.5) sera retranché pour chaque mauvaise réponse coloriée et pour chaque bonne réponse non coloriée.
- La correction étant automatique, un non respect des consignes aura forcément un impact sur la note finale.

Durée : 1 heure 20 minutes.

Exercice 1. On considère $(X_1, Y_1), \dots, (X_n, Y_n)$ un n échantillon i.i.d où X_i est à valeurs dans \mathbb{R} et Y_i dans $\{0, 1\}$. On désigne par g une règle de classification et par $L(g) = \mathbf{P}(g(X) \neq Y)$ son erreur de classification. On note g^* la règle de Bayes (règle optimale pour l'erreur de classification).

Question 1 ♣ Soit $x \in \mathbb{R}$. Cocher la (ou les) assertion(s) vraie(s) :

- | | |
|---|---|
| <input type="checkbox"/> A $g : \{0, 1\} \rightarrow \mathbb{R}$ | <input type="checkbox"/> F $L(g^*) > L(g)$ |
| <input type="checkbox"/> B $g^*(x) = 1$ si $\mathbf{P}(Y = 1 X = x) = 0.45$ | <input type="checkbox"/> G $g^*(x) = 1$ si $\mathbf{P}(Y = 1 X = x) > 0.5$ |
| <input type="checkbox"/> C $g^*(x) = 1$ si $\mathbf{P}(Y = 1 X = x) < 0.45$ | <input type="checkbox"/> H $g^*(x) = 1$ si $\mathbf{P}(Y = 1 X = x) = 0.65$ |
| <input type="checkbox"/> D $g : \mathbb{R} \rightarrow \{0, 1\}$ | <input type="checkbox"/> I $L(g^*) \geq L(g)$ |
| <input type="checkbox"/> E $L(g^*) \leq L(g)$ | <input type="checkbox"/> J <i>Aucune de ces réponses n'est correcte.</i> |



Question 2 ♣ On suppose de plus que X suit une loi normale $\mathcal{N}(0, 1)$ et que, pour $x \in \mathbb{R}$, la loi de $Y|X = x$ est

- une loi de Bernoulli $\mathcal{B}(0.90)$ si $x \geq 0$;
- une loi de Bernoulli $\mathcal{B}(0.25)$ si $x < 0$;

Cocher la (ou les) assertion(s) vraie(s) :

- | | |
|---|--|
| <input type="checkbox"/> A $g^*(-1) = 0$ | <input type="checkbox"/> H $L(g^*) = 0$ |
| <input type="checkbox"/> B $L(g^*) = 7/40$ | <input type="checkbox"/> I $L(g^*) = 1/4$ |
| <input type="checkbox"/> C $g^*(1) = 0$ | <input type="checkbox"/> J $L(g^*) = 7/20$ |
| <input type="checkbox"/> D $L(g^*) = 9/40$ | <input type="checkbox"/> K $L(g^*) = 1$ |
| <input type="checkbox"/> E $L(g^*) = 31/80$ | <input type="checkbox"/> L $g^*(x) = 0$ si $x \geq 0$ |
| <input type="checkbox"/> F $L(g^*) = 9/80$ | <input type="checkbox"/> M <i>Aucune de ces réponses n'est correcte.</i> |
| <input type="checkbox"/> G $g^*(x) = 1$ si $x \geq 0$ | |

Question 3 ♣ On considère g_n un algorithme de prévision règle qui souffre de surapprentissage. Cocher la (ou les) assertion(s) vraie(s) :

- | | |
|--|--|
| <input type="checkbox"/> A g_n possède un biais élevé mais une variance faible | <input type="checkbox"/> C g_n possède un biais faible mais une variance élevée |
| <input type="checkbox"/> B les données d'apprentissage seront très bien ajustées par g_n | <input type="checkbox"/> D les données d'un échantillon test seront très bien classées par g_n |
| | <input type="checkbox"/> E <i>Aucune de ces réponses n'est correcte.</i> |

Question 4 ♣ Soit k un entier plus petit que n . On désigne par $g_{n,k}$ l'algorithme des k plus proches voisins. Cocher la (ou les) assertion(s) vraie(s) :

- | | |
|--|--|
| <input type="checkbox"/> A $g_{n,k}$ possède généralement une variance faible pour de grandes valeurs de k | <input type="checkbox"/> D $g_{n,k}$ aura tendance à surajuster si k est trop petit |
| <input type="checkbox"/> B $g_{n,k}$ possède généralement un biais élevé pour de petites valeurs de k | <input type="checkbox"/> E $g_{n,k}$ possède généralement un biais faible pour de petites valeurs de k |
| <input type="checkbox"/> C $g_{n,k}$ aura tendance à surajuster si k est trop grand | <input type="checkbox"/> F $g_{n,k}$ possède généralement une variance élevée pour de grandes valeurs de k |
| | <input type="checkbox"/> G <i>Aucune de ces réponses n'est correcte.</i> |

Exercice 2. On dispose de $n = 20$ observations $(x_i, y_i), i = 1, \dots, n$ avec $x_i \in \mathbb{R}^2$ et $y_i \in \{0, 1\}$. On souhaite construire un arbre de classification pour expliquer les y_i par les x_i . La figure 1 propose deux coupures différentes du nœud racine. On considèrera comme mesure d'impureté d'un nœud N l'impureté de Gini. On rappelle qu'elle se met sous la forme $\mathcal{I}(N) = 2p(1-p)$ où p désigne un paramètre présenté en cours.

Question 5 On désigne par N le nœud racine. L'impureté $\mathcal{I}(N)$ de N vaut :

- | | | | | | |
|--------------------------------|------------------------------------|---------------------------------|---|----------------------------------|-----------------------------------|
| <input type="checkbox"/> A 0.5 | <input type="checkbox"/> B 11/100 | <input type="checkbox"/> C 9/20 | <input type="checkbox"/> D 11/20 | <input type="checkbox"/> E 9/100 | <input type="checkbox"/> F 99/200 |
| | <input type="checkbox"/> G 101/200 | <input type="checkbox"/> H 0 | <input type="checkbox"/> I <i>Aucune réponse n'est correcte</i> | | |

Question 6 L'impureté $\mathcal{I}(N1)$ de $N1$ vaut :

- | | | | | | |
|--------------------------------|---------------------------------|---|--------------------------------|----------------------------------|---------------------------------|
| <input type="checkbox"/> A 2/5 | <input type="checkbox"/> B 3/25 | <input type="checkbox"/> C 12/50 | <input type="checkbox"/> D 1/4 | <input type="checkbox"/> E 12/25 | <input type="checkbox"/> F 6/25 |
| | <input type="checkbox"/> G 0.5 | <input type="checkbox"/> H <i>Aucune réponse n'est correcte</i> | | | |

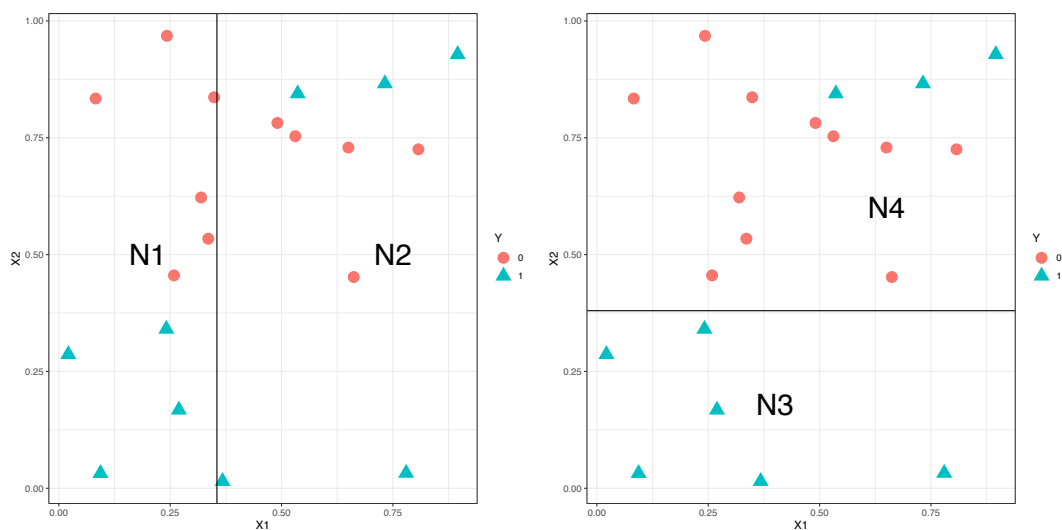


Figure 1: 2 exemples de coupure pour le même échantillon.

Question 7 L'impureté $\mathcal{I}(N3)$ de $N3$ vaut :

- ☐ A 0
 ☐ B 0.25
 ☐ C 0.75
 ☐ D 3/10
 ☐ E 0.5
 ☐ F 0.125
☐ G 1
☐ H Aucune réponse n'est correcte

Pour les questions suivantes, nous considérons les commandes R suivantes.

```
> arbre <- rpart(Y~.,data=df,cp=0.0001,minsplit=2)
> printcp(arbre)

Classification tree:
rpart(formula = Y ~ ., data = df, cp = 1e-04, minsplit = 2)

Variables actually used in tree construction:
[1] X1 X2

Root node error: ????
```

n= 20

	CP	nsplit	rel error	xerror	xstd
1	0.66667	0	1.00000	1.00000	0.24721
2	0.22222	1	aaaaaaa	0.33333	0.17743
3	0.11111	2	bbbbbbb	0.66667	0.22771
4	0.00010	3	ccccccc	0.22222	0.085434

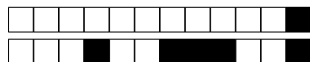
Question 8 La colonne `xstd` permet d'estimer la variance (ou l'écart-type) des termes calculés dans la colonne

- ☐ A `xerror`
☐ B `CP`
☐ C `nsplit`
☐ D `rel error`
☐ E Aucune réponse n'est correcte

Question 9 Retrouver la valeur du terme manquant `????` après `Root node error`.

- ☐ A 0.45
☐ B 0
☐ C 0.75
☐ D 0.6
☐ E 0.125
☐ F 0.5
☐ G 1
☐ H 0.55
☐ I 0.25
☐ J Aucune réponse n'est correcte

On entre ensuite les commandes suivantes.



```
> arbre1 <- prune(arbre,cp=0.2223)
> arbre2 <- prune(arbre,cp=0.0001)
> rpart.plot(arbre1)
> rpart.plot(arbre2)
```

On obtient les arbres représentés sur la figure 2.

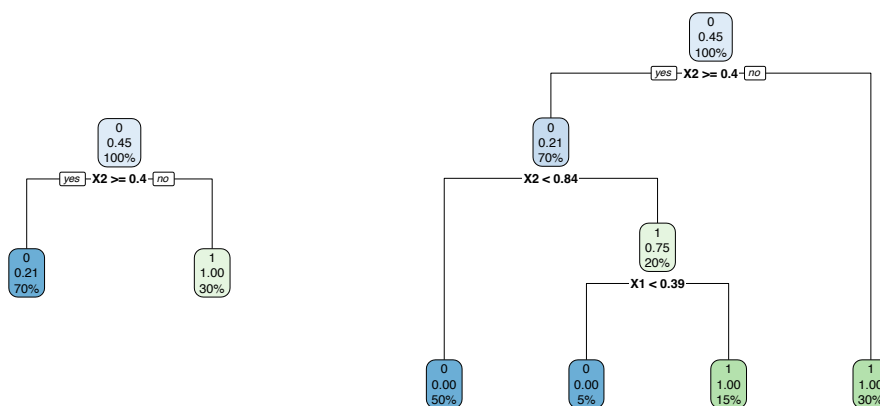


Figure 2: 2 arbres obtenus avec la fonction `rpart.plot` (`arbre1` à gauche et `arbre2` à droite).

Question 10 Retrouver la valeur du terme manquant `aaaaaaa` dans le tableau obtenu avec la commande `printcp`.

- ☐ A 3/4 ☐ B 1 ☐ C 1/4 ☐ D 3/10 ☐ E 1/3 ☐ F 7/20
☐ G 1/2 ☐ H 3/20 ☐ I 3/14 ☐ J Aucune réponse n'est correcte

Question 11 Retrouver la valeur du terme manquant `bbbbbbb` dans le tableau obtenu avec la commande `printcp`. **Pas de point négatif à cette question, on peut donc répondre sans craintes.**

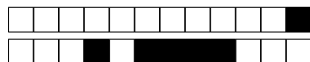
- ☐ A 0 ☐ B 1/4 ☐ C 1/8 ☐ D 1 ☐ E 1/9 ☐ F 1/11
☐ G 3/20 ☐ H 1/2 ☐ I 1/3 ☐ J 1/20
☐ K Aucune réponse n'est correcte

Question 12 Retrouver la valeur du terme manquant `ccccccc` dans le tableau obtenu avec la commande `printcp`.

- ☐ A 1/3 ☐ B 1/2 ☐ C 1/9 ☐ D 0 ☐ E 1/8 ☐ F 1/11 ☐ G 1
☐ H 1/20 ☐ I 3/20 ☐ J 1/4 ☐ K Aucune réponse n'est correcte

Question 13 Parmi les arbres suivants, lequel est le plus pertinent selon vous (selon la méthode d'élagage présentée en cours) ?

- ☐ A arbre ☐ B arbre1 ☐ C arbre2 ☐ D Aucun



Question 14 ♣ Cette question porte sur les arbres en général (elle n'est pas en lien avec les questions précédentes). Cocher la (ou les) assertion(s) vraie(s) :

- | | |
|---|--|
| <input type="checkbox"/> A Un arbre très profond possède généralement peu de biais | <input type="checkbox"/> D Un arbre peu profond possède généralement une forte variance |
| <input type="checkbox"/> B Le risque de surapprentissage augmente avec la profondeur de l'arbre | <input type="checkbox"/> E La qualité de prédiction d'un arbre augmente avec sa profondeur |
| <input type="checkbox"/> C Plus un arbre est profond, mieux il ajuste les données | <input type="checkbox"/> F <i>Aucune de ces réponses n'est correcte.</i> |

Exercice 3. Les questions suivantes portent sur les forêts aléatoires.

Question 15 ♣ Cocher la (ou les) assertion(s) vraie(s) :

- | | |
|--|---|
| <input type="checkbox"/> A Il n'est pas possible de faire des forêts aléatoires en présence de variables explicatives qualitatives | <input type="checkbox"/> D Une forêt aléatoire permet de réduire le biais des arbres qu'elle agrège |
| <input type="checkbox"/> B Une forêt aléatoire surapprend lorsqu'on agrège trop d'arbres | <input type="checkbox"/> E On peut faire de la classification supervisée et de la régression avec des forêts aléatoires |
| <input type="checkbox"/> C On doit choisir un nombre d'arbres le plus grand possible pour une forêt aléatoire | <input type="checkbox"/> F <i>Aucune de ces réponses n'est correcte.</i> |

Question 16 ♣ Cocher la (ou les) assertion(s) vraie(s) :

- | | |
|--|--|
| <input type="checkbox"/> A On doit choisir des arbres avec un biais élevé pour une forêt aléatoire. | <input type="checkbox"/> C On doit choisir des arbres avec une variance élevée pour une forêt aléatoire. |
| <input type="checkbox"/> B On doit choisir des arbres avec un biais faible pour une forêt aléatoire. | <input type="checkbox"/> D Il faut utiliser des arbres peu profonds pour les forêts aléatoires. |
| | <input type="checkbox"/> E <i>Aucune de ces réponses n'est correcte.</i> |

Question 17 Au niveau temps de calcul, l'erreur OOB d'une forêt aléatoire se calcule plus rapidement que l'erreur par validation croisée 10 blocs.

- | | |
|---------------------------------|---------------------------------|
| <input type="checkbox"/> A Faux | <input type="checkbox"/> B Vrai |
|---------------------------------|---------------------------------|

Exercice 4. Les questions suivantes portent sur le logiciel R.

Question 18 ♣ La fonction `tune_grid` du package `tidymodels` permet

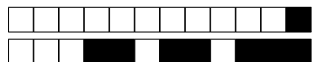
- | | |
|--|--|
| <input type="checkbox"/> A d'estimer des risques pour des algorithmes d'apprentissage supervisé. | <input type="checkbox"/> B de sélectionner les paramètres de certains algorithmes d'apprentissage supervisé. |
| | <input type="checkbox"/> C <i>Aucune de ces réponses n'est correcte.</i> |

Question 19 ♣

- | | |
|--|---|
| <input type="checkbox"/> A L'option <code>minsplit</code> dans <code>rpart</code> permet de modifier la profondeur des arbres construits | <input type="checkbox"/> D L'option <code>cp</code> dans <code>rpart</code> permet de modifier la profondeur des arbres construits |
| <input type="checkbox"/> B L'option <code>cp</code> dans <code>rpart</code> permet de calculer des erreurs de prévision | <input type="checkbox"/> E La fonction <code>prune</code> du package <code>rpart</code> permet de sélectionner un arbre dans une sous-suite d'arbres. |
| <input type="checkbox"/> C L'option <code>cp</code> dans <code>rpart</code> permet de calculer des erreurs d'ajustement | <input type="checkbox"/> F <i>Aucune de ces réponses n'est correcte.</i> |

Question 20 Le paramètre `num.trees` de la fonction `ranger` doit être

- | | | |
|--------------------------------------|----------------------------------|----------------------------------|
| <input type="checkbox"/> A ça dépend | <input type="checkbox"/> B petit | <input type="checkbox"/> C grand |
|--------------------------------------|----------------------------------|----------------------------------|



Question 21 Le paramètre `mtry` de la fonction **ranger** doit être

- ☐ A grand ☐ B ça dépend ☐ C petit

Question 22 Le paramètre `min.node.size` de la fonction **ranger** doit être

- ☐ A petit ☐ B ça dépend ☐ C grand



Feuille de réponses :

0	0	0	0	0	0	0	0
1	1	1	1	1	1	1	1
2	2	2	2	2	2	2	2
3	3	3	3	3	3	3	3
4	4	4	4	4	4	4	4
5	5	5	5	5	5	5	5
6	6	6	6	6	6	6	6
7	7	7	7	7	7	7	7
8	8	8	8	8	8	8	8
9	9	9	9	9	9	9	9

← codez votre numéro d'étudiant ci-contre, et inscrivez votre nom et prénom ci-dessous.

Nom et prénom :

.....

.....

Les réponses aux questions sont à donner exclusivement sur cette feuille : les réponses données sur les feuilles précédentes ne seront pas prises en compte.

- QUESTION 1 :

A	B	C	D	E	F	G	H	I	J
---	---	---	---	---	---	---	---	---	---
- QUESTION 2 :

A	B	C	D	E	F	G	H	I	J	K	L	M
---	---	---	---	---	---	---	---	---	---	---	---	---
- QUESTION 3 :

A	B	C	D	E
---	---	---	---	---
- QUESTION 4 :

A	B	C	D	E	F	G
---	---	---	---	---	---	---
- QUESTION 5 :

A	B	C	D	E	F	G	H	I
---	---	---	---	---	---	---	---	---
- QUESTION 6 :

A	B	C	D	E	F	G	H
---	---	---	---	---	---	---	---
- QUESTION 7 :

A	B	C	D	E	F	G	H
---	---	---	---	---	---	---	---
- QUESTION 8 :

A	B	C	D	E
---	---	---	---	---
- QUESTION 9 :

A	B	C	D	E	F	G	H	I	J
---	---	---	---	---	---	---	---	---	---
- QUESTION 10 :

A	B	C	D	E	F	G	H	I	J
---	---	---	---	---	---	---	---	---	---
- QUESTION 11 :

A	B	C	D	E	F	G	H	I	J	K
---	---	---	---	---	---	---	---	---	---	---
- QUESTION 12 :

A	B	C	D	E	F	G	H	I	J	K
---	---	---	---	---	---	---	---	---	---	---
- QUESTION 13 :

A	B	C	D
---	---	---	---
- QUESTION 14 :

A	B	C	D	E	F
---	---	---	---	---	---
- QUESTION 15 :

A	B	C	D	E	F
---	---	---	---	---	---
- QUESTION 16 :

A	B	C	D	E
---	---	---	---	---
- QUESTION 17 :

A	B
---	---
- QUESTION 18 :

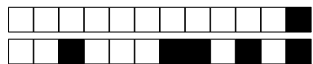
A	B	C
---	---	---
- QUESTION 19 :

A	B	C	D	E	F
---	---	---	---	---	---
- QUESTION 20 :

A	B	C
---	---	---
- QUESTION 21 :

A	B	C
---	---	---
- QUESTION 22 :

A	B	C
---	---	---



+1/8/53+