
Examen

On dispose de la base d'exemples suivante représentant les résultats d'analyse de 15 patients concernant l'infection par le corona virus COVID19 :

N°	AGE	DR	FV	TX	FAT	COVID19
1	V	Oui	Non	Oui	Oui	+
2	J	Oui	Oui	Non	Non	+
3	V	Non	Oui	Non	Non	-
4	E	Non	Non	Oui	Oui	+
5	J	Oui	Non	Oui	Non	-
6	E	Oui	Oui	Non	Oui	-
7	V	Non	Non	Oui	Oui	+
8	V	Oui	Oui	Non	Non	-
9	J	Non	Oui	Oui	Non	-
10	J	Non	Non	Oui	Oui	+
11	V	Oui	Oui	Non	Non	+
12	J	Oui	Non	Oui	Oui	-
13	V	Non	Oui	Oui	Non	-
14	E	Non	Non	Oui	Oui	+
15	J	Oui	Oui	Non	Non	-

- N° : Numéro du patient
- DR : Difficulté respiratoire (Oui, Non)
- TX : Toux (Oui, Non)
- AGE : l'âge du patient (J :Jeune, V :Vieux, E :enfant)
- FV : Fièvre (Oui, Non)
- COVID 19 : Résultat du test PCR (infecté par le virus COVID19 : + (oui), - (non))

On vous demande de :

1. Règles d'association
 - (a) Construire un modèle de décision en utilisant la méthode OneR sur les exemples de 1 à 10.
 - (b) Calculer sa précision et son rappel sur la base d'entraînement (composée des exemples de 1 à 10) et sur la base de test (composée des exemples de 11 à 15).
 - (c) Commenter le résultat.
2. Arbre de décision
 - (a) Construire, en utilisant les exemples de 1 à 10, l'arbre de décision permettant de prédire l'infection par le virus COVID-19 en utilisant l'algorithme ID3 avec le Gain d'information pour le choix des attributs.
 - (b) Calculer la précision et le rappel du modèle sur la base d'entraînement puis sur la base de test composée des exemples de 11 à 15.
 - (c) Commenter le résultat.
3. Comparer les deux modèles.

Bonne Chance

Corrigé type

1. Règles d'association

(7 pt)

(a) Construire un modèle de décision en utilisant la méthode OneR sur les exemples de 1 à 10.

– Modèle OneR

(2.5 pts)

AGE	+	-	DR	+	-	FV	+	-	TX	+	-	FAT	+	-
V	2	2	Oui	2	3	Oui	1	4	Oui	4	2	Oui	4	1
J	2	2	Non	3	2	Non	4	1	Non	1	3	Non	1	4
E	1	1	Err		4	Err		2	Err		3	Err		2
Err		5												

– Le modèle :

(2 pts)

- Si FV = Oui Alors COVID19 = -
- Sinon COVID19 = +

(b) Précision

– Sur la base d'entraînement = $8/10 = 80\%$

(0.25 pt)

– Sur la base de test = $3/5 = 60\%$

(0.25 pt)

(c) Rappel

–

– Sur la base d'entraînement = $CP/(CP + FN) = 4/(4+1) = 80\%$

(0.5 pt)

– Sur la base de test = $CP/(CP + FN) = 1/(1+1) = 50\%$

(0.5 pt)

(d) Commenter le résultat : Le modèle représente une modèle surentrainé (sur-apprentissage) (1 pt)

2. Arbre de décision

(12 pts)

(a) Construction de l'arbre de décision

– $H(\text{covid}) = H(5,5) = 1$

(0.5 pt)

$H(\text{Age}) = 4/10 H(V) + 4/10 H(J) + 2/10 H(E) = 4/10 H(2,2) + 4/10 H(2,2) + 2/10 H(1,1)$
 $= 0.4 + 0.4 + 0.2 = 1$

$\text{Gain}(\text{Age}) = 1 - 1 = 0$

(0.5 pt)

$H(\text{DR}) = 5/10 H(\text{oui}) + 5/10 H(\text{non}) = 5/10 H(2,3) + 5/10 H(2,3) = 0.971$

$\text{Gain}(\text{DR}) = 1 - 0.971 = 0.029$

(0.5 pt)

$H(\text{FV}) = 5/10 H(\text{oui}) + 5/10 H(\text{non})$

$= 5/10 H(4,1) + 5/10 H(4,1) = 0.722$ **Variable de segmentation**

$\text{Gain}(\text{FV}) = 1 - 0.722 = 0.278$

(0.5 pt)

$H(\text{TX}) = 6/10 H(\text{oui}) + 4/10 H(\text{non}) = 6/10 H(4,2) + 4/10 H(3,1) = 0.875$

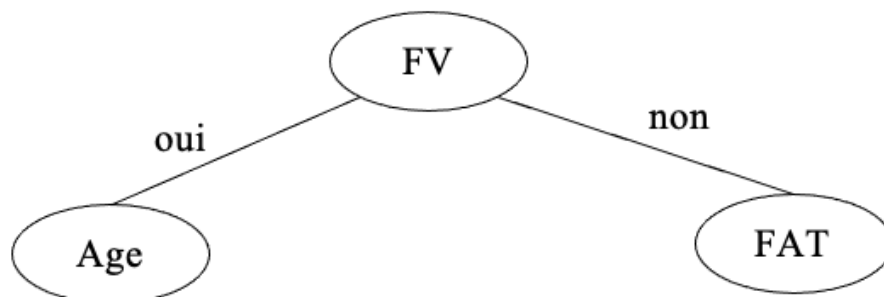
$\text{Gain}(\text{TX}) = 1 - 0.875 = 0.125$

(0.5 pt)

$H(\text{FAT}) = 5/10 H(\text{oui}) + 5/10 H(\text{non}) = 5/10 H(4,1) + 5/10 H(4,1) = 0.722$

$\text{Gain}(\text{FAT}) = 0.278$

(0.5 pt)



(0.5 pt)

- **FV = Non**

$H(\text{COVID10}) = H(4,1) = 0.722$ (0.5 pt)

$H(\text{AGE}) = 2/5H(J) + 2/5H(V) + 1/5H(E) = 2/5H(1,1) + 2/5H(2,0) + 1/5H(1,0) = 0.4$

$\text{Gain}(\text{AGE}) = 0.722 - 0.4 = 0.322$ (0.5 pt)

$H(\text{DR}) = 3/5H(\text{Oui}) + 2/5H(\text{Non}) = 3/5H(2,1) + 2/5H(2,0) = 0.55$

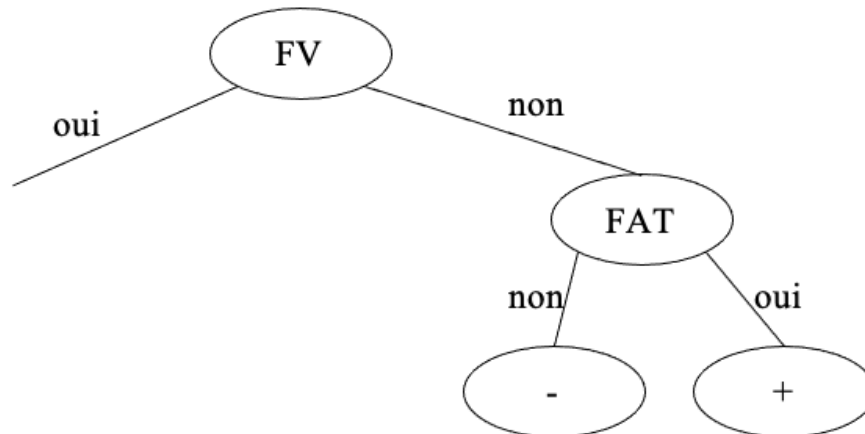
$\text{Gain}(\text{DR}) = 0.722 - 0.55 = 0.17$ (0.5 pt)

$H(\text{TX}) = 4/5H(3,1) + 1/5H(1,0) = 0.64$

$\text{Gain}(\text{TX}) = 0.722 - 0.64 = 0.08$ (0.5 pt)

$H(\text{FAT}) = 4/5H(\text{Oui}) + 1/5H(\text{Non}) = 4/5H(4,0) + 1/5H(1,0) = 0$

$\text{Gain}(\text{FAT}) = 0.722 - 0 = 0.722$ **Variable de segmentation** (0.5 pt)



- **FV = Oui**

$H(\text{COVID10}) = H(4,1) = 0.722$ (0.5 pt)

$H(\text{AGE}) = 2/5H(J) + 2/5H(V) + 1/5H(E) = 2/5H(1,1) + 2/5H(2,0) + 1/5H(1,0) = 0.4$

$\text{Gain}(\text{AGE}) = 0.322$ **Variable de segmentation**

$H(\text{DR}) = 3/5H(\text{Oui}) + 2/5H(\text{Non}) = 3/5H(2,1) + 2/5H(2,0) = 0.55$

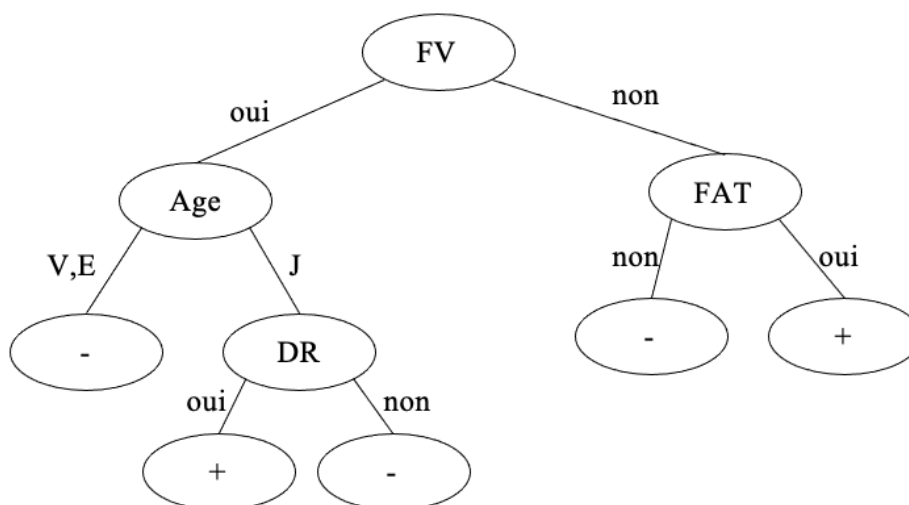
$\text{Gain}(\text{DR}) = 0.171$ (0.5 pt)

$H(\text{TX}) = 4/5H(3,1) = 0.722$

$\text{Gain}(\text{TX}) = 0$ (0.5 pt)

$H(\text{FAT}) = 1/5H(\text{Oui}) + 4/5H(\text{Non}) = 1/5H(1,0) + 4/5H(3,1) = 0.64$

$\text{Gain}(\text{FAT}) = 0.0722$ (0.5 pt)



(1.5 pt)

(b) Précision

– Sur la base d’entraînement = $10/10 = 100\%$ **(0.25 pt)**

– Sur la base de test = $2/5 = 40\%$ **(0.25 pt)**

(c) Rappel

– Sur la base d’entraînement = $CP/(CP + FN) = 5/5 = 100\%$ **(0.5 pt)**

– Sur la base de test = $CP/(CP + FN) = 1/(1+1) = 50\%$ **(0.5 pt)**

(d) Commenter le résultat : Le modèle représente une modèle surentrainé (sur-apprentissage) **(1 pt)**

3. Le deuxième modèle est plus précis que le premier. **(1 pt)**