

## Éléments de correction CC Apprentissage -2017

### Exercice 1 :

classe réelle \ classe prédite	A	B	C	Total
A	35	1	2	38
B	2	30	0	32
C	3	2	25	30
Total	40	33	27	100

Soit la matrice de confusion ci-dessus, obtenue par application d'un modèle M1 de fouille de données. Après avoir rappelé la formule de calcul, déterminer :

1. (1 pt)

- Le taux d'erreur en généralisation : **10%**

- accuracy rate : **90%**

2. La précision pour la classe A ; (1 pt)

$$P_A = (A,A) / (.,A) = 35/40$$

3. Le rappel pour la classe B ; (1 pt)

$$R_B = (B,.B) / (B, .) = 30/32$$

4. Le taux de faux positifs (FP rate) pour la classe C ; (1 pt)

$$FP_C = ((A,C) + (B,C)) / ((A,.) + (B,.)) = (2 + 0) / (38 + 32) = 2/70$$

5 . Le taux de vrais positifs (TP rate) pour la classe A ; (1 pt)

$$VP_A = (A,A) / (A,.) = 35/38$$

6. La sensibilité pour la classe B ; (1 pt)

$$Sensi_B = VP_B = (B,B) / (B,.) = 30/32$$

7. La spécificité pour la classe C ; (1 pt)

$$Speci_C = VN_C = ((A,A) + (B,B)) / ((A,.) + (B,.)) = (35+30)/(38+32) = 65/70$$

8. La F-mesure de la classe A. (1 pt)

$$F\text{-mesure}_A = (2 * R_A * P_A) / (R_A + P_A) = (2 * 35/38 * 35/40) / (35/38 + 35/40)$$

## Exercice 2 :

1. L'attribut **Custmer\_ID** n'est pas pertinent car sa valeur est unique pour chaque enregistrement (**Overfitting**) (1 pt)

2. Détermination de la racine de l'arbre de décision

	Gendre		Car_Type			Shirt_Size			
	M	F	Family	Sports	Luxury	S	M	L	XL
<b>C0</b>	3	2	0	4	1	1	1	2	1
<b>C1</b>	2	3	1	0	4	1	1	1	2
<b>GINI</b>	<b>0,48</b>	<b>0,48</b>	<b>0</b>	<b>0</b>	<b>0,32</b>	<b>0,5</b>	<b>0,5</b>	<b>0,44</b>	<b>0,44</b>
	<b>0,48</b>		<b>0,16</b>			<b>0,464</b>			

$$\text{GINI}_{\text{Gendre}}(\text{M}) = 1 - (3/5)^2 - (2/5)^2 ; \quad \text{GINI}_{\text{Gendre}}(\text{F}) = 1 - (2/5)^2 - (3/5)^2$$

$$\text{GINI}_{\text{split}}(\text{Gendre}) = 5/10 * (1 - (3/5)^2 - (2/5)^2) + 5/10 * (1 - (2/5)^2 - (3/5)^2) = 0,48 \quad (1 \text{ pt})$$

$$\text{GINI}(\text{Gendre}) = 1 - (5/10)^2 - (5/10)^2 = 0,5 \quad (0,5 \text{ pt})$$

$$\text{GINI}_{\text{Car\_Type}}(\text{Family}) = 1 - (0/1)^2 - (1/1)^2 = 0 ; \quad \text{GINI}_{\text{Car\_Type}}(\text{Sports}) = 1 - (4/4)^2 - (0/4)^2 = 0 ;$$

$$\text{GINI}_{\text{Car\_Type}}(\text{Luxury}) = 1 - (1/5)^2 - (4/5)^2 = 0,32$$

$$\text{GINI}_{\text{split}}(\text{Car\_Type}) = 0 + 0 + 5/10 * (1 - (1/5)^2 - (4/5)^2) = 0,16 \quad (1 \text{ pt})$$

$$\text{GINI}(\text{Car\_Type}) = 1 - (5/10)^2 - (5/10)^2 = 0,5 \quad (0,5 \text{ pt})$$

$$\text{GINI}_{\text{Shirt\_Size}}(\text{S}) = 1 - (1/2)^2 - (1/2)^2 ; \quad \text{GINI}_{\text{Shirt\_Size}}(\text{M}) = 1 - (1/2)^2 - (1/2)^2 ;$$

$$\text{GINI}_{\text{Shirt\_Size}}(\text{L}) = 1 - (2/3)^2 - (1/3)^2 ; \quad \text{GINI}_{\text{Shirt\_Size}}(\text{XL}) = 1 - (1/3)^2 - (2/3)^2$$

$$\text{GINI}_{\text{split}}(\text{Shirt\_Size}) = 2/10 * (1 - (1/2)^2 - (1/2)^2) + 2/10 * (1 - (1/2)^2 - (1/2)^2) + 3/10 * (1 - (2/3)^2 - (1/3)^2) + 3/10 * (1 - (1/3)^2 - (2/3)^2) \approx 0,464 \quad (1 \text{ pt})$$

$$\text{GINI}(\text{Shirt\_Size}) = 1 - (5/10)^2 - (5/10)^2 = 0,5 \quad (0,5 \text{ pt})$$

- L'attribut **Car\_Type** a le plus petit index GINI, on choisit **Car\_Type** comme la racine de l'arbre de décision.

- Partitionnement avec les attributs non-choisis (1 pt)

En partant de nœud **Car\_Type**, on obtient les étudiants {2,4,6,8} qui sont de classe C0 en prenant la branche « Sports », {12} qui est de classe C1 en prenant la branche «Family» et {10,14,16,18,20} qui ne sont pas de même classe en prenant la branche «Luxury»,

Nous devons traiter les lignes suivantes :

Ligne	10	14	16	18	20
Gendre	F	M	F	F	F
Shirt_Size	L	XL	S	M	L
Classe	C0	C1	C1	C1	C1

Nous observons les attributs Gendre et Shirt\_Size :

	Gendre		Shirt_Size			
	M	F	S	M	L	XL
C0	0	1	0	0	1	0
C1	1	3	1	1	1	1
GINI	0	0,375	0	0	0,5	0
	0,3		0,2			

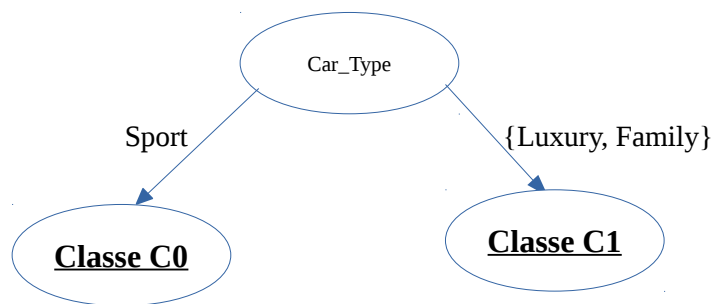
$GINI(\text{Gendre}) = 0,32$  ;  $GINI_{\text{split}}(\text{Gendre}) = 0,3$  ;

$GINI(\text{Shirt\_Size}) = 0,32$  ;  $GINI_{\text{split}}(\text{Shirt\_Size}) = 0,2$

- L'attribut **Shirt\_Size** a le plus petit index de GINI, cependant, ces différentes branches contiennent très peu d'éléments. Il est donc moins pertinent.

- En choisissant l'attribut **Gendre**, le gain n'est pas vraiment important après la division en deux branches (« M » et « F »), car  $GINI(\text{Gendre}) \approx GINI_{\text{split}}(\text{Gendre})$ . (1 pt)

On obtient l'arbre de décision suivant :



3 - Évaluation du modèle M1 sur D2.

(2 pts)

Customer_ID	Classe réelle	M1
1	C0	C1 (faux)
3	C0	C0
5	C0	C0
7	C0	C0
9	C0	C0
11	C1	C1
13	C1	C1
15	C1	C1
17	C1	C1
19	C1	C1

3 – Matrice de confusion sur D2

(2 pts)

	Classes prédites			
Classes réelles		C0	C1	Total
	C0	4	1	5
	C1	0	5	5
	Total	4	6	10

4. Erreur apparente :  $\text{Erreur}_{\text{App}}(\text{M1}) = 0,1$

(1pt)