



Module 7 - Arbres de décisions

Exercices - Corrigé

Exercice 2

L'entropie peut être calculée selon :

Individu	Taille	Couleur	Class
1	P	U	Pure
2	P	U	Pure
3	G	U	Pure
4	G	M	Pure
5	G	M	NPure
6	G	M	NPure
7	P	U	NPure
8	P	M	NPure

$$E(P) = - \sum_{i=1}^n p_i \log p_i \quad (1)$$

avec \log désigne la fonction de logarithme en base 2.

1. L'entropie de la variable couleur est donnée par :

$$E(\text{Couleur}) = \frac{4}{8} \times I(3, 1) + \frac{4}{8} \times I(1, 3)$$

avec

$$I(3, 1) = -\frac{3}{4} \times \log\left(\frac{3}{4}\right) - \frac{1}{4} \times \log\left(\frac{1}{4}\right) = 0.8113$$

$$I(1, 3) = -\frac{1}{4} \times \log\left(\frac{1}{4}\right) - \frac{3}{4} \times \log\left(\frac{3}{4}\right)$$

d'où :

$$E(\text{Couleur}) = 2 \times \left(\frac{4}{8}\right) \times I(3, 1) = 0.8113$$

2. L'entropie de la variable taille est donnée par :

$$E(\text{Taille}) = \frac{4}{8} \times I(2, 2) + \frac{4}{8} \times I(2, 2)$$

avec

$$I(2, 2) = -\frac{2}{4} \times \log\left(\frac{2}{4}\right) - \frac{2}{4} \times \log\left(\frac{2}{4}\right) = 1$$

d'où :

$$E(\text{Taille}) = \frac{4}{8} \times 1 + \frac{4}{8} \times 1 = 1$$

3. Soit un ensemble de données T , le gain d'informations de T par rapport à une partition T_j donnée est la variation d'entropie causée par la partition de T selon T_j .

$$\text{Gain}(X, T) = E(T) - E(X, T) = E(T) - \sum_{j=1}^m \frac{|T_j|}{|T|} E(T_j) \quad (2)$$

Nous pouvons, donc, calculer le gain de chacune des variables :

$$E(T) = E(\text{Race}) = 1$$

$$\text{Gain}(\text{Race}, \text{Couleur}) = 1 - E(\text{Couleur}) = 1 - 0.8113 = 0.19$$

$$\text{Gain}(\text{Race}, \text{Taille}) = 0$$

Le gain de taille étant nul, cette variable ne constitue pas un bon prédicteur de la race.

Exercice 3

1. Calcul de l'entropie de la variable **Couleur** :

```

1 install.packages("DescTools")
2 library("DescTools")
3 Data<- matrix(c(1,1, 1, 1,2, 1, 1, 1,3, 0, 1, 1,4, 0, 0,
4               1,5, 0, 0, 0, 6, 0, 0, 0, 6, 1, 1, 0, 6, 1, 0, 0),
5               nrow=8, ncol=4, byrow=T)
6 # Entropie de la variable Couleur
7 EntrCouleur=Entropy(Data[,3],Data[,4],base=2)-1
8 View(EntrCouleur)
9 # Entropie de la variable classe
10 EntrClasse=Entropy(Data[,4],base=2)-1
11 View(EntrClasse)
12 # Calcul du Gain de la variable Couleur
13 GainCouleur=EntrClasse-EntrCouleur
14 View(GainCouleur)

```

L'entropie de la variable Couleur est 0,8113

2. Calcul de l'entropie de la variable Taille :

```

1 # Entropie de la variable classe
2 EntrClasse=Entropy(Data[,4],base=2)-1
3 View(EntrClasse)
4 # Calcul du Gain de la variable Taille
5 GainTaille=EntrClasse-EntrTaille
6 View(GainTaille)

```

L'entropie de la variable Taille est 1

3. Le gain qui est defini par :

```

1 # Calcul du gain de la variable "Taille"
2 I<-function(p,n)
3 {
4   return(((p/(n+p))*log(p/(p+n)))-(n/(n+p))*log(n/(p+n)))
5 }
6 Entropie<- -(4/8)*log(4/8)-(4/8)*log(4/8)
7 GTaille<-Entropie - EntropieTaille

```

On peut aussi calculer le gain en utilisant la bibliotheque "FSelector" :

```

1 # Installer la bibliotheque "FSelector"
2 install.packages("FSelector")
3
4 # Charger la bibliotheque

```

```

5 library(FSelector)
6
7 weights <- information.gain(Data$class~Data$Taille + Data
8   $Couleur , Data, unit = "log2")

```

Donc, puisque la fonction Gain égale à zéro, alors, la variable Taille ne constitue pas un bon prédicteur de la race.

Exercice 4

1. Installer et charger les bibliothèques rpart et FSelector :

```

1 install.packages("rpart", dep=TRUE)
2 library(rpart)
3 install.packages("FSelector")
4 library(FSelector)

```

2. Importation des deux base de données : iris et kyphosis

```

1 # Importer la base de donnees "iris"
2 data(iris)
3
4 # Importer la base de donnees "kyphosis"
5 data(kyphosis)

```

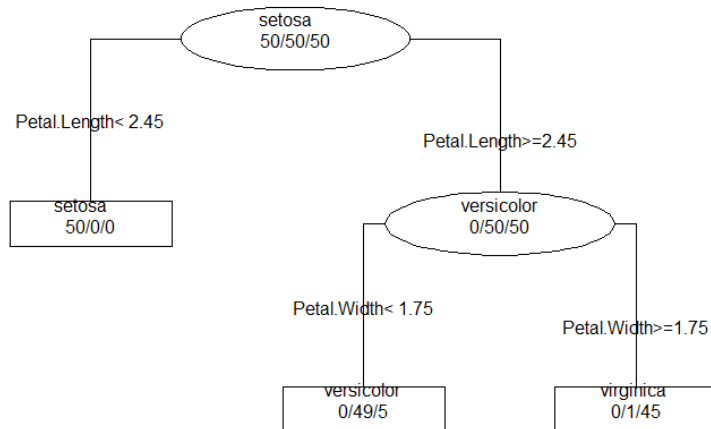
3. Les arbres de décisions des ensembles de données sont :

```

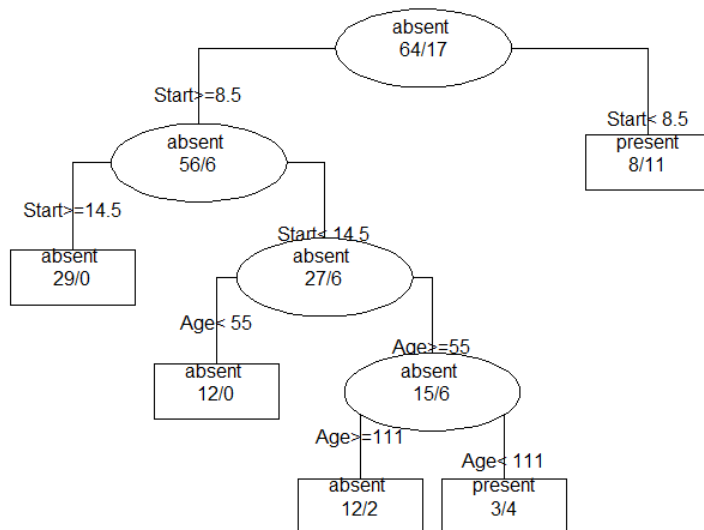
1 # Arbre de decision de la base de donnees "iris"
2 dev.new()
3 arbreIris = rpart(iris$Species ~ ., method="class",
4   minsplit=20, xval=81, data=iris)
5 plot(arbreIris, uniform=TRUE, margin=0.1, main="Arbre de
6   decision 'iris'")
7 text(arbreIris, fancy=TRUE, use.n=TRUE, pretty=0, all=
8   TRUE)
9 # Arbre de decision de la base de donnees "Kyphosis"
10 dev.new()
11 arbreKyphosis = rpart(Kyphosis ~ ., method="class",
12   minsplit=20, xval=81, data= kyphosis)
13 plot(arbreKyphosis, uniform=TRUE, margin=0.1, main="Arbre
14   de decision 'kyphosis'")
15 text(arbreKyphosis, fancy=TRUE, use.n=TRUE, pretty=0, all
16   =TRUE)

```

Arbre de décision 'iris'



Arbre de décision 'kyphosis'



4. Rapport de gain et gain pour chaque attribut :

```

1 # Calcul Gain ration et Gain pour les attributs
2 print("Iris Database")
3 information.gain(Species~., data=iris , unit = "log2")
4 gain.ratio(Species ~. , data=iris , unit = "log2")
5
6 print("kyphosis Database")
7 information.gain(Kyphosis~., data=kyphosis , unit = "log2")
8 gain.ratio(Kyphosis ~. , data=kyphosis , unit = "log2")
  
```

```

> print("Iris Database")
[1] "Iris Database"
> information.gain(Species~., data=iris, unit = "log2")
      attr_importance
Sepal.Length      0.6522837
Sepal.width       0.3855963
Petal.Length      1.3565450
Petal.width       1.3784027
> gain.ratio(Species ~. , data=iris, unit = "log2")
      attr_importance
Sepal.Length      0.4196464
Sepal.width       0.2472972
Petal.Length      0.8584937
Petal.width       0.8713692
> print("kyphosis Database")
[1] "kyphosis Database"
> information.gain(kyphosis~., data=kyphosis, unit = "log2")
      attr_importance
Age              0.0000000
Number          0.0000000
Start           0.1690017
> gain.ratio(kyphosis ~. , data=kyphosis, unit = "log2")
      attr_importance
Age              0.0000000
Number          0.0000000
Start           0.1712875
> |
  
```

5. Les attributs qui ont plus d'importance pour la base de données 'iris' sont :

- Petal.Length
- Petal.With

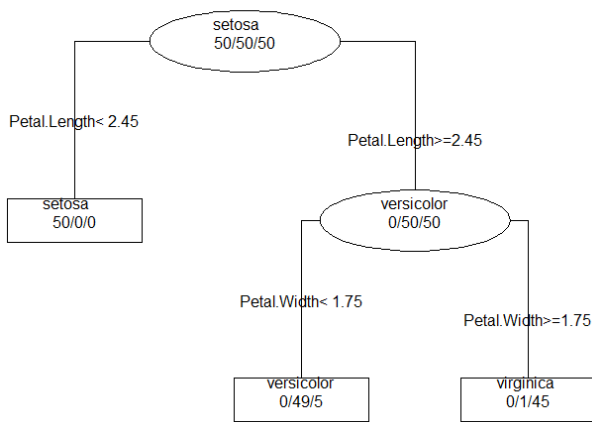
L'attribut qui a plus d'importance pour la base de données 'kyphosis' est :

— Start

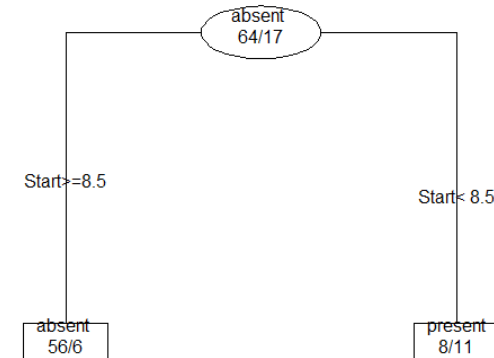
6. Après élimination des attributs non importants, nous obtenons les arbres suivants :

```
1 # Arbre de decision de la base de donnees "iris"
2 dev.new()
3 arbre = rpart(iris$Species ~ Petal.Length +Petal.Width,
4               method="class", minsplit=20, xval=81, data=iris)
5 plot(arbre, uniform=TRUE, margin=0.1, main="Arbre de
6       decision 'iris'")
7 text(arbre, fancy=TRUE, use.n=TRUE, pretty=0, all=TRUE)
8
9 # Arbre de decision de la base de donnees "Kyphosis"
10 dev.new()
11 arbre = rpart(Kyphosis ~ Start, method="class", minsplit
12               =20, xval=81, data= kyphosis)
13 plot(arbre, uniform=TRUE, margin=0.1, main="Arbre de
14       decision 'kyphosis'")
15 text(arbre, fancy=TRUE, use.n=TRUE, pretty=0, all=TRUE)
```

Arbre de décision 'iris'



Arbre de décision 'kyphosis'



Exercice 4

1. Chargement de la base de données "CreditData.txt" en utilisant R :

```
1 # Charger la base de donnees
2 Data<-read.csv("C:/Users/youssef/Dropbox/INF1421-Module7-
3               ArbresDecision/R/Exercice4/CreditData.txt", header =
4               TRUE, sep = ",")
```

2. Divisez la base de données en base d'apprentissage et en base de test :

```
1 #Base d'entrainement
2 Train<-Data[251 :1000,]
3 #Base de Test
4 Test<-Data[1 :250,]
```

3. Affichage de l'arbre de décision complète de la base d'apprentissage.

— Algorithmes CART :

```
1 # Arbre de decision CART
2 library(rpart)
3 install.packages("rpart.plot")
4 library(rpart.plot)
5 dev.new()
6 Arbre1= rpart(Train$Creditability ~ ., data= Train)
7 prp(Arbre1, extra=1)
8 title("Arbre de decision Cart - base d'entrainement")
```

