

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

Exercice 3:

1. Déroulement de l'algorithme ID3 :

➤ **L'entropie de la population:**

$$E(C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}) = -\left(\frac{9}{14} \log_2 \left(\frac{9}{14}\right) + \frac{5}{14} \log_2 \left(\frac{5}{14}\right)\right)$$

$$= -\left(0.64 \frac{\ln\left(\frac{9}{14}\right)}{\ln(2)} + 0.35 \frac{\ln\left(\frac{5}{14}\right)}{\ln(2)}\right) = 0.92$$

➤ **Calcul du gain de chaque attribut:**

A. Calcul de gain du Ciel:

1. Ciel= Ensoleillé:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{2}{5} \\ C_{\text{Jouer= Non}} = \frac{3}{5} \end{array} \right. \quad \frac{| \text{Ciel= Ensoleillé} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{5}{14}$$

$$E(\text{Ciel} = \text{Ensoleillé}) = -\left(\frac{2}{5} \log_2 \left(\frac{2}{5}\right) + \frac{3}{5} \log_2 \left(\frac{3}{5}\right)\right) = 0.96 \quad /* \text{ Entropie de l'attribut}$$

2. Ciel= Pluie:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{3}{5} \\ C_{\text{Jouer= Non}} = \frac{2}{5} \end{array} \right. \quad \frac{| \text{Ciel= Pluie} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{5}{14}$$

$$E(\text{Ciel} = \text{Pluie}) = -\left(\frac{3}{5} \log_2 \left(\frac{3}{5}\right) + \frac{2}{5} \log_2 \left(\frac{2}{5}\right)\right) = 0.96$$

3. Ciel=Couvert:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{4}{4} = 1 \\ C_{\text{Jouer= Non}} = 0 \end{array} \right. \quad \frac{| \text{Ciel= Couvert} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{4}{14}$$

$$E(\text{Ciel} = \text{Couvert}) = 0$$

D'où

$$\text{Gain}(\text{Ciel}) = 0.92 - \left(0.96 * \frac{5}{14} + 0.96 * \frac{5}{14} + 0 * \frac{5}{14}\right)$$

$$= 0.24$$

B. Calcul de gain de la Température:

1. Température= Chaude:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{2}{4} \\ C_{\text{Jouer= Non}} = \frac{2}{4} \end{array} \right. \quad \frac{| \text{Température= Chaude} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{4}{14}$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

$$E(\text{Température} = \text{Chaude}) = -\left(\frac{1}{2}\log_2\left(\frac{1}{2}\right) + \frac{1}{2}\log_2\left(\frac{1}{2}\right)\right) = 1$$

2. Température = Tiède:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{4}{6} \\ C_{\text{Jouer= Non}} = \frac{2}{6} \end{array} \right. \quad \frac{|\text{Température} = \text{Tiède}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{6}{14}$$

$$E(\text{Température} = \text{Tiède}) = -\left(\frac{4}{6}\log_2\left(\frac{4}{6}\right) + \frac{2}{6}\log_2\left(\frac{2}{6}\right)\right) = 0.92$$

3. Température = Fraîche:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{3}{4} \\ C_{\text{Jouer= Non}} = \frac{1}{4} \end{array} \right. \quad \frac{|\text{Température} = \text{Fraîche}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{4}{14}$$

$$E(\text{Température} = \text{Fraîche}) = -\left(\frac{3}{4}\log_2\left(\frac{3}{4}\right) + \frac{1}{4}\log_2\left(\frac{1}{4}\right)\right) = 0.8$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}) &= 0.92 - \left(1 * \frac{4}{14} + 0.92 * \frac{6}{14} + 0.8 * \frac{4}{14}\right) \\ &= 0.03 \end{aligned}$$

C. Calcul de gain d'Humidité:

1. Humidité = Elevée:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{3}{7} \\ C_{\text{Jouer= Non}} = \frac{4}{7} \end{array} \right. \quad \frac{|\text{Humidité} = \text{Elevée}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{7}{14}$$

$$E(\text{Humidité} = \text{Elevée}) = -\left(\frac{3}{7}\log_2\left(\frac{3}{7}\right) + \frac{4}{7}\log_2\left(\frac{4}{7}\right)\right) = 0.98$$

2. Humidité = Normale:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{6}{7} \\ C_{\text{Jouer= Non}} = \frac{1}{7} \end{array} \right. \quad \frac{|\text{Humidité} = \text{Normale}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{7}{14}$$

$$E(\text{Humidité} = \text{Normale}) = -\left(\frac{6}{7}\log_2\left(\frac{6}{7}\right) + \frac{1}{7}\log_2\left(\frac{1}{7}\right)\right) = 0.58$$

D'où

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

$$\begin{aligned} \text{Gain(Humidité)} &= 0.92 - (0.98 * \frac{7}{14} + 0.58 * \frac{7}{14}) \\ &= 0.14 \end{aligned}$$

D. Calcule de gain du Vent:

1. Vent =Faible:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{6}{8} \\ C_{\text{Jouer= Non}} = \frac{2}{8} \end{array} \right. \quad \frac{|\text{Vent = Faible}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{8}{14}$$

$$E(\text{Vent} = \text{faible}) = -\left(\frac{6}{8} \log_2\left(\frac{6}{8}\right) + \frac{2}{8} \log_2\left(\frac{2}{8}\right)\right) = 0.8$$

2. Vent =Fort:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{3}{6} \\ C_{\text{Jouer= Non}} = \frac{3}{6} \end{array} \right. \quad \frac{|\text{Vent = Fort}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{8}{14}$$

$$E(\text{Vent} = \text{Fort}) = -\left(\frac{3}{6} \log_2\left(\frac{3}{6}\right) + \frac{3}{6} \log_2\left(\frac{3}{6}\right)\right) = 1$$

D'où

$$\begin{aligned} \text{Gain(Vent)} &= 0.92 - (0.8 * \frac{8}{14} + 1 * \frac{6}{14}) \\ &= 0.05 \end{aligned}$$

On a :

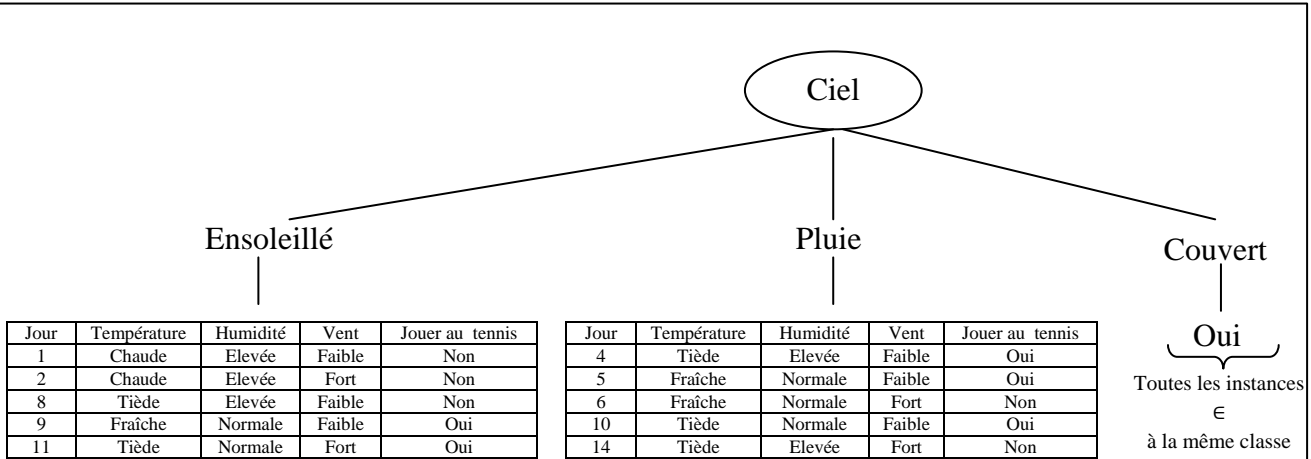
$$\left\{ \begin{array}{l} \text{Gain(Température)} = 0.03 \\ \text{Gain(Ciel)} = 0.24 \\ \text{Gain(Humidité)} = 0.14 \\ \text{Gain(Vent)} = 0.05 \end{array} \right.$$

Donc :

$$\text{Gain(Ciel)} > \text{Gain(Humidité)} > \text{Gain(Vent)} > \text{Gain(Température)}$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD



➤ **Raffinement de la branche « Ensoleillé »:**

$$E(C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}) = -\left(\frac{2}{5} \log_2 \left(\frac{2}{5}\right) + \frac{3}{5} \log_2 \left(\frac{3}{5}\right)\right) = 0.96$$

Calcul du gain pour chaque attribut:

A. Calcul de gain de la Température:

1. Température = Chaude:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{2}{2} = 0 \\ C_{\text{Jouer= Non}} = 1 \end{array} \right. \quad \frac{| \text{Température= Chaude} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{2}{5}$$

$$E(\text{Température} = \text{Chaude}) = 0$$

2. Température = Tiède:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{1}{2} \\ C_{\text{Jouer= Non}} = \frac{1}{2} \end{array} \right. \quad \frac{| \text{Température= Tiède} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{2}{5}$$

$$E(\text{Température} = \text{Tiède}) = 1$$

3. Température = Fraîche:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = 1 \\ C_{\text{Jouer= Non}} = 0 \end{array} \right. \quad \frac{| \text{Température= Fraîche} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{1}{5}$$

$$E(\text{Température} = \text{Fraîche}) = 0$$

D'où

$$\text{Gain}(\text{Température}) = 0.96 - (0 * \frac{2}{5} + 1 * \frac{2}{5} + 0 * \frac{1}{5}) = 0.56$$

B. Calcule de gain d'Humidité:

1. Humidité= Elevée:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = 0 \\ C_{\text{Jouer= Non}} = \frac{3}{3} = 1 \end{array} \right. \quad \frac{| \text{Humidité = Elevée} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{3}{5}$$

$$E(\text{Humidité} = \text{Elevée}) = 0$$

2. Humidité= Normale:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{2}{2} = 1 \\ C_{\text{Jouer= Non}} = 0 \end{array} \right. \quad \frac{| \text{Humidité = Normale} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{2}{5}$$

$$E(\text{Humidité} = \text{Normale}) = 0$$

D'où

$$\begin{aligned} \text{Gain}(\text{Humidité}) &= 0.96 - (0 * \frac{3}{5} + 0 * \frac{2}{5}) \\ &= 0.96 \end{aligned}$$

C. Calcule de gain du Vent:

1. Vent =Faible:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{1}{3} \\ C_{\text{Jouer= Non}} = \frac{2}{3} \end{array} \right. \quad \frac{| \text{Vent = Faible} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{3}{5}$$

$$E(\text{Vent} = \text{faible}) = -\left(\frac{1}{3} \log_2 \left(\frac{1}{3}\right) + \frac{2}{3} \log_2 \left(\frac{2}{3}\right)\right) = 0.92$$

2. Vent =Fort:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{1}{2} \\ C_{\text{Jouer= Non}} = \frac{1}{2} \end{array} \right. \quad \frac{| \text{Vent = Fort} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{2}{5}$$

$$E(\text{Vent} = \text{Fort}) = -\left(\frac{1}{2} \log_2 \left(\frac{1}{2}\right) + \frac{1}{2} \log_2 \left(\frac{1}{2}\right)\right) = 1$$

D'où

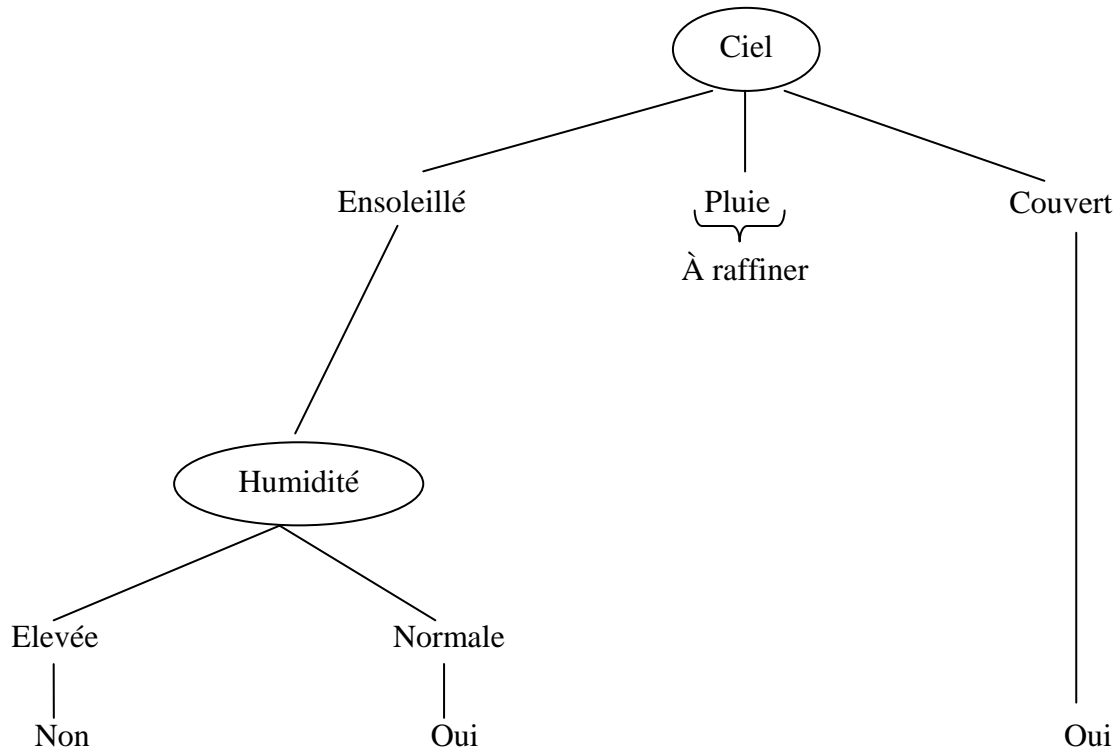
$$\begin{aligned} \text{Gain}(\text{Vent}) &= 0.96 - (0.92 * \frac{3}{5} + 1 * \frac{2}{5}) \\ &= 0.008 \end{aligned}$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

On a :

$$\left\{ \begin{array}{l} \text{Gain(Température)} = 0.56 \\ \text{Gain(Humidité)} = 0.96 \\ \text{Gain(Vent)} = 0.008 \end{array} \right. \Rightarrow \text{Gain maximal est de l'humidité}$$



➤ **Raffinement de la branche « Pluie »:**

$$E(C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}) = -\left(\frac{3}{5} \log_2 \left(\frac{3}{5}\right) + \frac{2}{5} \log_2 \left(\frac{2}{5}\right)\right) = 0.96$$

Calcul du gain pour chaque attribut:

A. Calcul de gain de la Température :

1. Température = Tiède:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{2}{3} \\ C_{\text{Jouer= Non}} = \frac{1}{3} \end{array} \right. \quad \frac{| \text{Température= Tiède} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{3}{5}$$

$$E(\text{Température} = \text{Tiède}) = 0.92$$

2. Température = Fraîche:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{1}{2} \\ C_{\text{Jouer= Non}} = \frac{1}{2} \end{array} \right. \quad \frac{| \text{Température= Fraîche} |}{| C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}} |} = \frac{2}{5}$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

$$E(\text{Température} = \text{Fraîche}) = 1$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}) &= 0.96 - (0.92 * \frac{3}{5} + 1 * \frac{2}{5}) \\ &= 0.01 \end{aligned}$$

B. Calcul de gain d'Humidité:

1. Humidité= Elevée:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{1}{2} \\ C_{\text{Jouer= Non}} = \frac{1}{2} \end{array} \right. \quad \frac{|\text{Humidité} = \text{Elevée}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{2}{5}$$

$$E(\text{Humidité} = \text{Elevée}) = 1$$

2. Humidité= Normale:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{2}{3} \\ C_{\text{Jouer= Non}} = \frac{1}{3} \end{array} \right. \quad \frac{|\text{Humidité} = \text{Normale}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{3}{5}$$

$$E(\text{Humidité} = \text{Normale}) = 0.92$$

D'où

$$\begin{aligned} \text{Gain}(\text{Humidité}) &= 0.96 - (1 * \frac{2}{5} + 0.92 * \frac{3}{5}) \\ &= 0.008 \end{aligned}$$

C. Calcul de gain du Vent:

1. Vent = Faible:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = \frac{3}{3} = 1 \\ C_{\text{Jouer= Non}} = 0 \end{array} \right. \quad \frac{|\text{Vent} = \text{Faible}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{3}{5}$$

$$E(\text{Vent} = \text{faible}) = 0$$

2. Vent = Fort:

$$\left\{ \begin{array}{l} C_{\text{Jouer= Oui}} = 0 \\ C_{\text{Jouer= Non}} = \frac{2}{2} = 1 \end{array} \right. \quad \frac{|\text{Vent} = \text{Fort}|}{|C_{\text{Jouer= Oui}}, C_{\text{Jouer= Non}}|} = \frac{2}{5}$$

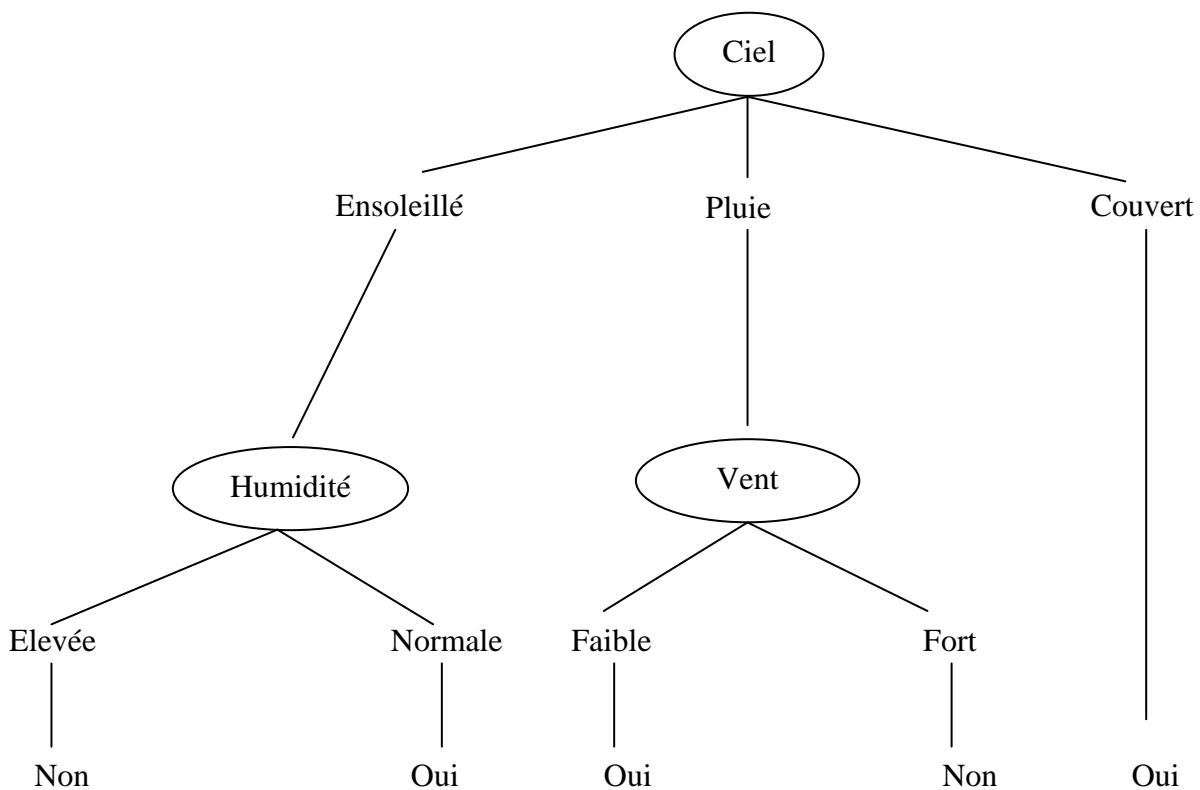
$$E(\text{Vent} = \text{Fort}) = 0$$

D'où

$$\begin{aligned} \text{Gain}(\text{Vent}) &= 0.96 - (0 * \frac{3}{5} + 0 * \frac{2}{5}) \\ &= 0.96 \end{aligned}$$

On a:

$$\left\{ \begin{array}{l} \text{Gain}(\text{Température}) = 0.01 \\ \text{Gain}(\text{Humidité}) = 0.008 \\ \text{Gain}(\text{Vent}) = 0.96 \end{array} \right. \Rightarrow \text{Vent a un gain maximal}$$



Arbre final

2. Déroulement de l'Algorithme KNN

- Distance euclidienne

$$\begin{cases} d(A,A)=0 \\ d(A,B)=1 \end{cases}$$

- Le paramètre $K=4+1=5$

$$\text{d}(y, x_i) / i \in \{1..14\} = \begin{cases} 0 & \text{si } y=x_i \\ 1 & \text{Sinon} \end{cases}$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

N°	Ciel	Température	Humidité	Vent	Résultat	Classe
1	0	1	0	0	$\sqrt{\sum_{i=1}^n (xi - yi)^2} = 1$	Non
2	0	1	0	1	1.41	Non
3	1	1	0	0	1.41	Oui
4	1	1	0	0	1.41	Oui
5	1	0	1	0	1.41	Oui
6	1	0	1	1	1.73	Non
7	1	0	1	1	1.73	Oui
8	0	1	0	0	1	Non
9	0	0	1	0	1	Oui
10	1	1	1	0	1.73	Oui
11	0	1	1	1	1.73	Oui
12	1	1	0	1	1.73	Oui
13	1	1	1	0	1.73	Oui
14	1	1	0	1	1.73	Non
y	Ensoleillé	Fraîche	Elevée	Faible	Choisir 5 voisins	Non

Exercice 4:

1. e est un nouvel individu avec (Ciel=Ensoleillé, Vent = Faible)

La règle de décision est de classer l'individu e dans la classe jouer au tennis = oui SSI

$$P(\text{Jouer au tennis} = \text{Oui} | \text{Ciel} = \text{Ensoleillé}, \text{Vent} = \text{Faible})$$

>

$$P(\text{Jouer au tennis} = \text{Non} | \text{Ciel} = \text{Ensoleillé}, \text{Vent} = \text{Faible})$$

En appliquant la règle de Bayes, la règle devient:

$$P(\text{Ciel} = \text{Ensoleillé}, \text{Vent} = \text{Faible} | \text{Jouer au tennis} = \text{Oui}) P(\text{Jouer au tennis} = \text{Oui})$$

>

$$P(\text{Ciel} = \text{Ensoleillé}, \text{Vent} = \text{Faible} | \text{Jouer au tennis} = \text{Non}) P(\text{Jouer au tennis} = \text{Non})$$

$$P(\text{Jouer au tennis} = \text{Oui}) = \frac{9}{14}$$

$$P(\text{Jouer au tennis} = \text{Non}) = \frac{5}{14}$$

* $P(\text{Ciel} = \text{Ensoleillé}, \text{Vent} = \text{Faible} | \text{Jouer au tennis} = \text{Oui}) =$

$$P(\text{Ciel} = \text{Ensoleillé} | \text{Jouer au tennis} = \text{Oui}) P(\text{Vent} = \text{Faible} | \text{Jouer au tennis} = \text{Oui})$$

$$= \frac{2}{9} * \frac{6}{9} = \frac{12}{81}$$

$$= 0.14 * \frac{9}{14} = 0.09$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

$$\begin{aligned}
 & \textcircled{**} \quad \mathbf{P}(\text{Ciel}=\text{Ensoleillé}, \text{Vent}=\text{Faible} \mid \text{Jouer au tennis} = \text{Non}) = \\
 & \mathbf{P}(\text{Ciel}=\text{Ensoleillé} \mid \text{Jouer au tennis} = \text{Non}) \mathbf{P}(\text{Vent}=\text{Faible} \mid \text{Jouer au tennis} = \text{Non}) \\
 & = \frac{3}{5} * \frac{2}{5} = \frac{6}{25} \\
 & = 0.24 + \frac{5}{14} = 0.08
 \end{aligned}$$

On a : $\textcircled{*} > \textcircled{**}$

Donc :
 $e \in \text{jouer au tennis} = \text{oui}$.

2. $e = (\text{Ciel}=\text{Ensoleillé}, \text{temps} = 23^\circ\text{C}, \text{Humidité} = 70\%, \text{Vent}=\text{Faible})$

$$(0^\circ\text{C} * \frac{9}{5}) + 32 = ?? \text{ F}^\circ$$

$$(23^\circ\text{C} * \frac{9}{5}) + 32 = 73.4 \text{ F}^\circ$$

$e \in \text{jouer au tennis} = \text{Oui SSI}$

$$\mathbf{P}(\text{Ciel}=\text{Ensoleillé}, \text{Température} = 73.4 \text{ F}^\circ, \text{Humidité} = 70\%, \text{Vent} = \text{Faible} \mid \text{Jouer au tennis} = \text{Oui}) \mathbf{P}(\text{Jouer au tennis} = \text{Oui}) \textcircled{*}$$

>

$$\mathbf{P}(\text{Ciel}=\text{Ensoleillé}, \text{Température} = 73.4 \text{ F}^\circ, \text{Humidité} = 70\%, \text{Vent} = \text{Faible} \mid \text{Jouer au tennis} = \text{Non}) \mathbf{P}(\text{Jouer au tennis} = \text{Non}) \textcircled{**}$$

$$\begin{aligned}
 & \textcircled{*} \Leftrightarrow \mathbf{P}(\text{Ciel} = \text{Ensoleillé} \mid \text{Jouer au tennis} = \text{Oui}) \mathbf{P}(\text{Température} = 73.4 \text{ F}^\circ \mid \text{Jouer au tennis} \\
 & \quad = \text{Oui}) \\
 & \quad \mathbf{P}(\text{Humidité} = 70\% \mid \text{Jouer au tennis} = \text{Oui}) \mathbf{P}(\text{Vent} = \text{faible} \mid \text{Jouer au tennis} = \text{Oui}) \\
 & = \frac{2}{9} * 0 * \frac{2}{9} * \frac{5}{9} = 0 * \frac{9}{14} = 0
 \end{aligned}$$

$$\textcircled{**} = 0 \rightarrow e \text{ est un outlier (bruit)}$$

3. $\mathbf{P}(\text{Température} = 23^\circ\text{C} \mid \text{Jouer au tennis} = \text{Oui}) = 0$

4. $\mathbf{P}(60\% < \text{Humidité} < 75\% \mid \text{Jouer au tennis} = \text{Oui}) = \frac{3}{9} = 0.33$

Exercice 5:

- Déroulement de l'algorithme C4.5
- L'entropie de la population:

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

$$E(C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}) = -\left(\frac{5}{14} \log_2\left(\frac{5}{14}\right) + \frac{9}{14} \log_2\left(\frac{9}{14}\right)\right) = 0.94$$

A. Calcul de gain d'enseillement

1. Enseillement = Soleil

$$\left\{ \begin{array}{l} C_{\text{Jouer=Oui}} = \frac{2}{5} \\ C_{\text{Jouer=Non}} = \frac{3}{5} \end{array} \right. \quad \frac{|\text{Soleil}|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{5}{14}$$

$$E(\text{Enseillé}) = -\left(\frac{2}{5} \log_2\left(\frac{2}{5}\right) + \frac{3}{5} \log_2\left(\frac{3}{5}\right)\right) = 0.97$$

2. Enseillement = Couvert

$$\left\{ \begin{array}{l} C_{\text{Jouer=Oui}} = 1 \\ C_{\text{Jouer=Non}} = 0 \end{array} \right. \quad \frac{|\text{Couvert}|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{4}{14}$$

$$E(\text{Couvert}) = 0$$

3. Enseillement = Pluie

$$\left\{ \begin{array}{l} C_{\text{Jouer=Oui}} = \frac{3}{5} \\ C_{\text{Jouer=Non}} = \frac{2}{5} \end{array} \right. \quad \frac{|\text{Pluie}|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{5}{14}$$

$$E(\text{Pluie}) = -\left(\frac{3}{5} \log_2\left(\frac{3}{5}\right) + \frac{2}{5} \log_2\left(\frac{2}{5}\right)\right) = 0.97$$

D'où

$$\text{Gain(Enseillement)} = 0.94 - (0.97 * \frac{5}{14} + 0 * \frac{4}{14} + 0.97 * \frac{5}{14}) = 0.24$$

B. Calcul de gain d'Humidité

➤ Trier les exemples dans l'ordre croissant (ou décroissant)

Humidité	65	70	70	70	75	78	80	80	80	85	90	90	95	96
N°	8	1	5	11	9	7	10	12	13	3	2	6	4	14
Jouer	Oui	Oui	Oui	Non	Oui	Oui	Non	Oui	Oui	Non	Non	Oui	Non	Oui

Nous coupons entre les exemples:

$$\left\{ \begin{array}{l} x_5 \text{ et } x_{11} \rightarrow S_1 = 70 \\ x_{11} \text{ et } x_9 \rightarrow S_2 = 75 \\ x_7 \text{ et } x_{10} \rightarrow S_3 = 80 \\ x_{10} \text{ et } x_{12} \rightarrow S_4 = 80 \\ x_{13} \text{ et } x_3 \rightarrow S_5 = 85 \\ x_2 \text{ et } x_6 \rightarrow S_6 = 90 \\ x_6 \text{ et } x_4 \rightarrow S_7 = 95 \\ x_4 \text{ et } x_{14} \rightarrow S_8 = 96 \end{array} \right.$$

$$\begin{aligned} \text{➤ } E(\text{Humidité} \leq 70) &= -\left(\frac{3}{4} \log_2\left(\frac{3}{4}\right) + \frac{1}{4} \log_2\left(\frac{1}{4}\right)\right) \\ &= 0.81 \end{aligned}$$

$$\frac{|\text{Humidité} \leq 70|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{4}{14}$$

$$\begin{aligned} \text{➤ } E(\text{Humidité} > 70) &= -\left(\frac{6}{10} \log_2\left(\frac{6}{10}\right) + \frac{4}{10} \log_2\left(\frac{4}{10}\right)\right) \\ &= 0.97 \end{aligned}$$

$$\frac{|\text{Humidité} > 70|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{10}{14}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Humidité}, S_1=70) &= 0.94 - (0.81 * \frac{4}{14} + 0.97 * \frac{10}{14}) \\ &= 0.94 - 0.92 \\ &= 0.02 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Humidité} \leq 75) &= -\left(\frac{4}{5} \log_2\left(\frac{4}{5}\right) + \frac{1}{5} \log_2\left(\frac{1}{5}\right)\right) \\ &= 0.72 \end{aligned}$$

$$\frac{|\text{Humidité} \leq 75|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{5}{14}$$

$$\begin{aligned} \text{➤ } E(\text{Humidité} > 75) &= -\left(\frac{5}{9} \log_2\left(\frac{5}{9}\right) + \frac{4}{9} \log_2\left(\frac{4}{9}\right)\right) \\ &= 0.99 \end{aligned}$$

$$\frac{|\text{Humidité} > 75|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{9}{14}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Humidité}, S_2=75) &= 0.94 - (0.72 * \frac{5}{14} + 0.99 * \frac{9}{14}) \\ &= 0.94 - 0.89 \\ &= 0.05 \end{aligned}$$

$$\begin{aligned} \text{E}(\text{Humidité} \leq 80) &= -\left(\frac{7}{9} \log_2\left(\frac{7}{9}\right) + \frac{2}{9} \log_2\left(\frac{2}{9}\right)\right) & \frac{|\text{Humidité} \leq 80|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{9}{14} \\ &= 0.76 \end{aligned}$$

$$\begin{aligned} \text{E}(\text{Humidité} > 80) &= -\left(\frac{2}{5} \log_2\left(\frac{2}{5}\right) + \frac{3}{5} \log_2\left(\frac{3}{5}\right)\right) & \frac{|\text{Humidité} > 80|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{5}{14} \\ &= 0.97 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Humidité}, S_3=80) &= 0.94 - (0.76 * \frac{9}{14} + 0.97 * \frac{5}{14}) \\ &= 0.94 - 0.83 \\ &= 0.11 \end{aligned}$$

$$\begin{aligned} \text{E}(\text{Humidité} \leq 85) &= -\left(\frac{7}{10} \log_2\left(\frac{7}{10}\right) + \frac{3}{10} \log_2\left(\frac{3}{10}\right)\right) & \frac{|\text{Humidité} \leq 85|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{10}{14} \\ &= 0.88 \end{aligned}$$

$$\begin{aligned} \text{E}(\text{Humidité} > 85) &= -\left(\frac{2}{4} \log_2\left(\frac{2}{4}\right) + \frac{2}{4} \log_2\left(\frac{2}{4}\right)\right) & \frac{|\text{Humidité} > 85|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{4}{14} \\ &= 1 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Humidité}, S_5=85) &= 0.94 - (0.88 * \frac{10}{14} + 1 * \frac{4}{14}) \\ &= 0.94 - 0.91 \\ &= 0.03 \end{aligned}$$

$$\begin{aligned} \text{E}(\text{Humidité} \leq 90) &= -\left(\frac{8}{12} \log_2\left(\frac{8}{12}\right) + \frac{4}{12} \log_2\left(\frac{4}{12}\right)\right) & \frac{|\text{Humidité} \leq 90|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{12}{14} \\ &= 0.92 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Humidité} > 90) &= -\left(\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) & \frac{|\text{Humidité} > 90|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{2}{14} \\ &= 1 \end{aligned}$$

D'où

$$\begin{aligned} \text{➤ } \text{Gain}(\text{Humidité}, S_6=90) &= 0.94 - \left(0.92 * \frac{12}{14} + 1 * \frac{2}{14}\right) \\ &= 0.94 - 0.92 \\ &= 0.02 \end{aligned}$$

$$\begin{aligned} \text{➤ } \text{Gain}(\text{Humidité}, S_7=95) &= 0.94 - \left(0.96 * \frac{13}{14} + 0 * \frac{1}{14}\right) \\ &= 0.94 - 0.89 \\ &= 0.05 \end{aligned}$$

$$\begin{aligned} \text{➤ } \text{Gain}(\text{Humidité}, S_8=96) &= 0.94 - (0.94 * 1 + 0 * 0) \\ &= 0.94 - 0.94 \\ &= 0 \end{aligned}$$

Le gain maximal d'humidité est avec le seuil $S_3=80$

C. Calcul de gain de la Température

Trier les exemples dans l'ordre croissant

Température	64	65	68	69	70	71	72	72	75	75	80	81	83	85
N°	8	11	13	5	14	10	4	6	1	12	2	9	7	3
Jouer	Oui	Non	Oui	Oui	Oui	Non	Non	Oui	Oui	Oui	Non	Oui	Oui	Non

Nous coupons entre les exemples :

$$\left\{ \begin{array}{l} x_8 \text{ et } x_{11} \xrightarrow{S_1} S_1 = 65 \\ x_{11} \text{ et } x_{13} \xrightarrow{S_2} S_2 = 68 \\ x_{14} \text{ et } x_{10} \xrightarrow{S_3} S_3 = 71 \\ x_4 \text{ et } x_6 \xrightarrow{S_4} S_4 = 72 \\ x_{12} \text{ et } x_2 \xrightarrow{S_5} S_5 = 80 \\ x_2 \text{ et } x_9 \xrightarrow{S_6} S_6 = 81 \\ x_7 \text{ et } x_3 \xrightarrow{S_7} S_7 = 85 \end{array} \right.$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 65) &= -\left(\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) & \frac{|\text{Température} \leq 65|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{2}{14} \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 65) &= -\left(\frac{8}{12} \log_2\left(\frac{8}{12}\right) + \frac{4}{12} \log_2\left(\frac{4}{12}\right)\right) & \frac{|\text{Température} > 65|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{12}{14} \\ &= 0.92 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_1 = 65) &= 0.94 - \left(1 * \frac{2}{14} + 0.92 * \frac{12}{14}\right) \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 68) &= -\left(\frac{2}{3} \log_2\left(\frac{2}{3}\right) + \frac{1}{3} \log_2\left(\frac{1}{3}\right)\right) & \frac{|\text{Température} \leq 68|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{3}{14} \\ &= 0.92 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 68) &= -\left(\frac{7}{11} \log_2\left(\frac{7}{11}\right) + \frac{4}{11} \log_2\left(\frac{4}{11}\right)\right) & \frac{|\text{Température} > 68|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{11}{14} \\ &= 0.93 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_2 = 68) &= 0.94 - \left(0.92 * \frac{3}{14} + 0.93 * \frac{11}{14}\right) \\ &= 0.03 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 71) &= -\left(\frac{4}{6} \log_2\left(\frac{4}{6}\right) + \frac{2}{6} \log_2\left(\frac{2}{6}\right)\right) & \frac{|\text{Température} \leq 71|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{6}{14} \\ &= 0.92 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 71) &= -\left(\frac{5}{8} \log_2\left(\frac{5}{8}\right) + \frac{3}{8} \log_2\left(\frac{3}{8}\right)\right) & \frac{|\text{Température} > 71|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{8}{14} \\ &= 0.95 \end{aligned}$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

D'où

$$\text{Gain}(\text{Température}, S_3 = 71) = 0.94 - (0.92 * \frac{6}{14} + 0.95 * \frac{8}{14})$$

$$= 0.01$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 72) &= -\left(\frac{5}{8} \log_2\left(\frac{5}{8}\right) + \frac{3}{8} \log_2\left(\frac{3}{8}\right)\right) & \frac{|\text{Température} \leq 72|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{8}{14} \\ &= 0.95 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 72) &= -\left(\frac{4}{6} \log_2\left(\frac{4}{6}\right) + \frac{2}{6} \log_2\left(\frac{2}{6}\right)\right) & \frac{|\text{Température} > 72|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{6}{14} \\ &= 0.92 \end{aligned}$$

D'où

$$\text{Gain}(\text{Température}, S_4 = 72) = 0.94 - (0.95 * \frac{8}{14} + 0.92 * \frac{6}{14})$$

$$= 0.01$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 80) &= -\left(\frac{7}{11} \log_2\left(\frac{7}{11}\right) + \frac{4}{11} \log_2\left(\frac{4}{11}\right)\right) & \frac{|\text{Température} \leq 80|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{11}{14} \\ &= 0.93 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 80) &= -\left(\frac{2}{3} \log_2\left(\frac{2}{3}\right) + \frac{1}{3} \log_2\left(\frac{1}{3}\right)\right) & \frac{|\text{Température} > 80|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{3}{14} \\ &= 0.92 \end{aligned}$$

D'où

$$\text{Gain}(\text{Température}, S_5 = 80) = 0.94 - (0.93 * \frac{11}{14} + 0.92 * \frac{3}{14})$$

$$= 0.03$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 81) &= -\left(\frac{8}{12} \log_2\left(\frac{8}{12}\right) + \frac{4}{12} \log_2\left(\frac{4}{12}\right)\right) & \frac{|\text{Température} \leq 81|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{12}{14} \\ &= 0.92 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 81) &= -\left(\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) & \frac{|\text{Température} > 81|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{2}{14} \\ &= 1 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_6 = 81) &= 0.94 - \left(1 * \frac{2}{14} + 0.92 * \frac{12}{14}\right) \\ &= 0.01 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 85) &= -\left(\frac{9}{14} \log_2\left(\frac{9}{14}\right) + \frac{5}{14} \log_2\left(\frac{5}{14}\right)\right) & \frac{|\text{Température} \leq 85|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= 1 \\ &= 0.94 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 85) &= 0 & \frac{|\text{Température} > 85|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= 0 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_7 = 80) &= 0.94 - \left(0 * 0 + 0.94 * \frac{14}{14}\right) \\ &= 0 \end{aligned}$$

D. Calcul de gain du Vent

1. Vent = Oui

$$\begin{cases} C_{\text{Jouer=Oui}} = \frac{3}{6} \\ C_{\text{Jouer=Non}} = \frac{3}{6} \end{cases} \quad \frac{|\text{Vent} = \text{Oui}|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{6}{14}$$

$$E(\text{Vent} = \text{Oui}) = -\left(\frac{3}{6} \log_2\left(\frac{3}{6}\right) + \frac{3}{6} \log_2\left(\frac{3}{6}\right)\right) = 1$$

2. Vent = Non

$$\begin{cases} C_{\text{Jouer=Oui}} = \frac{6}{8} \\ C_{\text{Jouer=Non}} = \frac{2}{8} \end{cases} \quad \frac{|\text{Vent} = \text{Non}|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{8}{14}$$

$$E(\text{Vent} = \text{Non}) = -\left(\frac{6}{8} \log_2\left(\frac{6}{8}\right) + \frac{2}{8} \log_2\left(\frac{2}{8}\right)\right) = 0.81$$

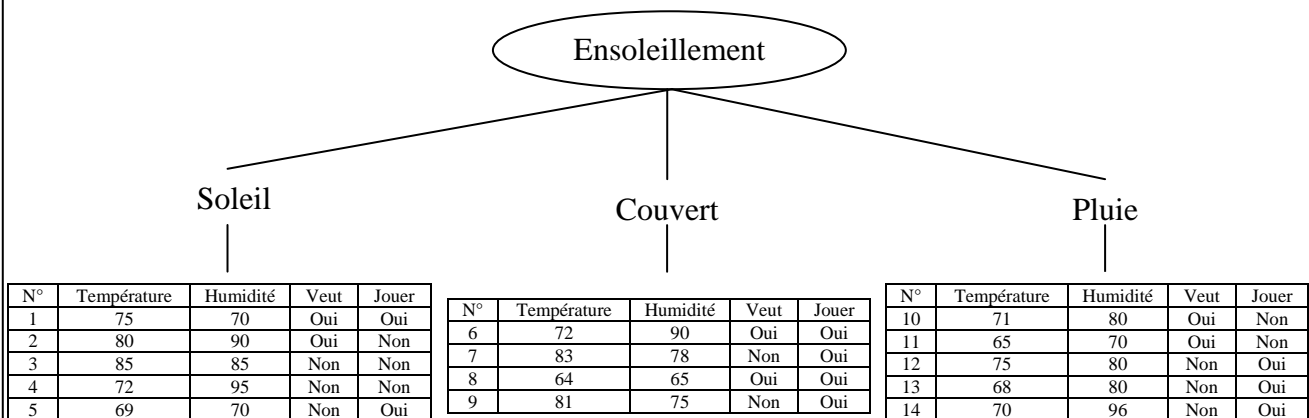
Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

D'où

$$\text{Gain}(\text{Vent}) = 0.94 - \left(1 * \frac{6}{14} + 0.81 * \frac{8}{14}\right) = 0.05$$

L'attribut ensoleillement a la plus grande valeur du gain en information. Donc, c'est le nœud racine de l'arbre:



??????

Jouer

??????

➤ **Raffinement de la branche « Soleil »:**

L'entropie de la population:

$$E(C_{\text{Jouer}=\text{Oui}}, C_{\text{Jouer}=\text{Non}}) = -\left(\frac{2}{5} \log_2 \left(\frac{2}{5}\right) + \frac{3}{5} \log_2 \left(\frac{3}{5}\right)\right) = 0.97$$

A. Calcul de gain de la Température :

Température	69	72	75	80	85
N°	5	4	1	2	3
Jouer	Oui	Non	Oui	Oui	Non

Nous coupons entre les exemples :

$$\left\{ \begin{array}{l} x_5 \text{ et } x_4 \longrightarrow S_1 = 72 \\ x_4 \text{ et } x_1 \longrightarrow S_2 = 75 \\ x_2 \text{ et } x_3 \longrightarrow S_3 = 85 \end{array} \right.$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023

N.BERMAD

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 72) &= -\left(\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) \\ &= 1 \end{aligned}$$

$$\frac{|\text{Température} \leq 72|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{2}{5}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 72) &= -\left(\frac{2}{3} \log_2\left(\frac{2}{3}\right) + \frac{1}{3} \log_2\left(\frac{1}{3}\right)\right) \\ &= 0.92 \end{aligned}$$

$$\frac{|\text{Température} > 72|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{3}{5}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_1 = 72) &= 0.97 - \left(1 * \frac{2}{5} + 0.92 * \frac{3}{5}\right) \\ &= 0.018 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 75) &= -\left(\frac{2}{3} \log_2\left(\frac{2}{3}\right) + \frac{1}{3} \log_2\left(\frac{1}{3}\right)\right) \\ &= 0.92 \end{aligned}$$

$$\frac{|\text{Température} \leq 75|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{3}{5}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 75) &= -\left(\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) \\ &= 1 \end{aligned}$$

$$\frac{|\text{Température} > 75|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{2}{5}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_2 = 75) &= 0.97 - \left(0.92 * \frac{3}{5} + 1 * \frac{2}{5}\right) \\ &= 0.018 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 85) &= -\left(\frac{3}{5} \log_2\left(\frac{3}{5}\right) + \frac{2}{5} \log_2\left(\frac{2}{5}\right)\right) \\ &= 0.97 \end{aligned}$$

$$\frac{|\text{Température} \leq 85|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = 1$$

$$\text{➤ } E(\text{Température} > 85) = 0$$

$$\frac{|\text{Température} > 85|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = 0$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_3 = 85) &= 0.97 - (0.97 * 1 + 0 * 0) \\ &= 0 \end{aligned}$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

B. Calcul de gain d'Humidité:

Humidité	70	70	85	90	95
N°	1	5	3	2	4
Jouer	Oui	Oui	Non	Non	Non

Nous coupons entre les exemples: x_5 et $x_3 \rightarrow S = \frac{85+70}{2} = 77.5$ (Car l'écart est grand)

$$\begin{aligned} \text{➤ } E(\text{Humidité} \leq 77.5) &= -\left(\frac{2}{2} \log_2\left(\frac{2}{2}\right) + 0 \log_2(0)\right) & \frac{|\text{Humidité} \leq 77.5|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{2}{5} \\ &= 0 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Humidité} > 77.5) &= -(0 \log_2(0) + 1 \log_2(1)) & \frac{|\text{Humidité} > 77.5|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{3}{5} \\ &= 0 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Humidité}, S=77.5) &= 0.97 - (0 * \frac{2}{5} + 0 * \frac{3}{5}) \\ &= 0.97 \end{aligned}$$

C. Calcul de gain du Vent

1. Vent = Oui

$$\begin{cases} C_{\text{Jouer=Oui}} = \frac{1}{2} \\ C_{\text{Jouer=Non}} = \frac{1}{2} \end{cases} \quad \frac{|\text{Vent} = \text{Oui}|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{2}{5}$$

$$E(\text{Vent} = \text{Oui}) = -\left(\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) = 1$$

2. Vent = Non

$$\begin{cases} C_{\text{Jouer=Oui}} = \frac{1}{3} \\ C_{\text{Jouer=Non}} = \frac{2}{3} \end{cases} \quad \frac{|\text{Vent} = \text{Non}|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{3}{5}$$

$$E(\text{Vent} = \text{Non}) = -\left(\frac{1}{3} \log_2\left(\frac{1}{3}\right) + \frac{2}{3} \log_2\left(\frac{2}{3}\right)\right) = 0.92$$

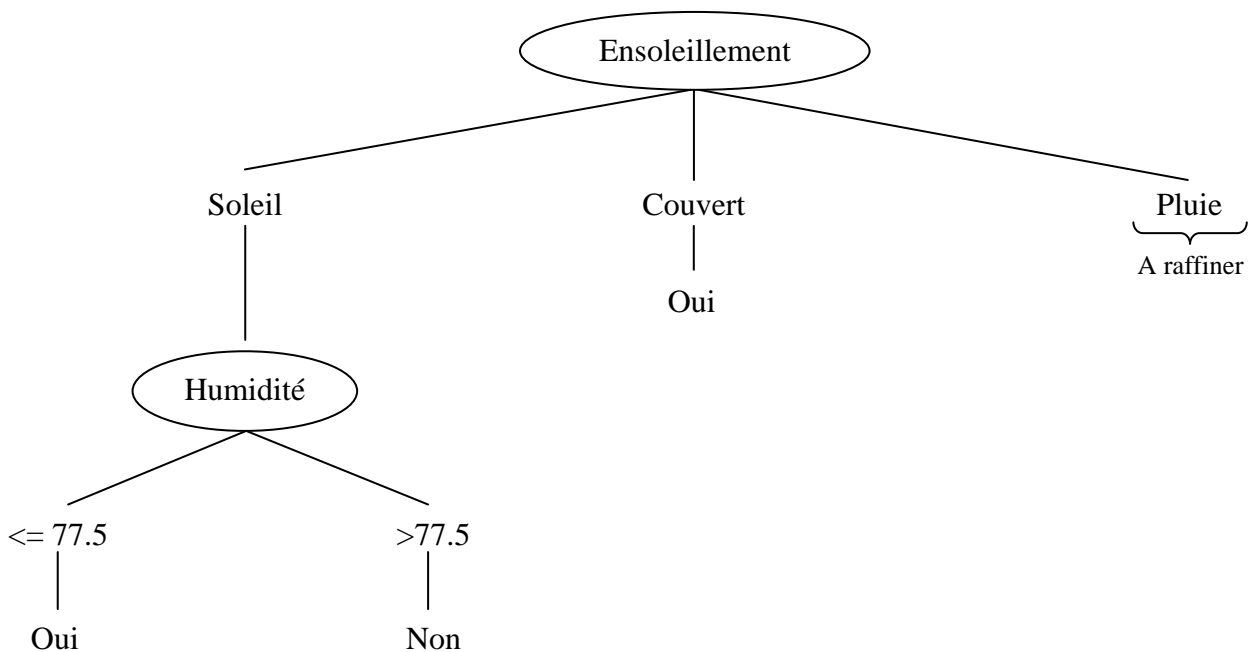
Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

D'où

$$\text{Gain}(\text{Vent}) = 0.97 - \left(1 * \frac{2}{5} + 0.92 * \frac{3}{5}\right) = 0.018$$

Nous choisissons l'attribut humidité ayant un gain maximal. Donc, l'arbre de décision est:



➤ **Raffinement de la branche « Pluie » :**

L'entropie de la population:

$$E(C_{\text{Jouer}=\text{Oui}}, C_{\text{Jouer}=\text{Non}}) = -\left(\frac{3}{5} \log_2 \left(\frac{3}{5}\right) + \frac{2}{5} \log_2 \left(\frac{2}{5}\right)\right) = 0.97$$

A. Calcul de gain de la Température :

Température	65	68	70	71	75
N°	11	13	14	10	12
Jouer	Non	Oui	Oui	Non	Oui

Nous coupons entre:

$$\left\{ \begin{array}{l} x_{11} \text{ et } x_{13} \longrightarrow S_1 = 68 \\ x_{14} \text{ et } x_{10} \longrightarrow S_2 = 71 \\ x_{10} \text{ et } x_{12} \longrightarrow S_3 = 75 \end{array} \right.$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023

N.BERMAD

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 68) &= -\left(\frac{1}{2} \log_2\left(\frac{1}{2}\right) + \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) & \frac{|\text{Température} \leq 68|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{2}{5} \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 68) &= -\left(\frac{2}{3} \log_2\left(\frac{2}{3}\right) + \frac{1}{3} \log_2\left(\frac{1}{3}\right)\right) & \frac{|\text{Température} > 68|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{3}{5} \\ &= 0.92 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_1 = 68) &= 0.97 - \left(1 * \frac{2}{5} + 0.92 * \frac{3}{5}\right) \\ &= 0.018 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 71) &= -\left(\frac{2}{4} \log_2\left(\frac{2}{4}\right) + \frac{2}{4} \log_2\left(\frac{2}{4}\right)\right) & \frac{|\text{Température} \leq 71|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{4}{5} \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 71) &= -(1 \log_2(1) + 0 \log_2(0)) & \frac{|\text{Température} > 71|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{1}{5} \\ &= 0 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_2 = 71) &= 0.97 - \left(1 * \frac{4}{5} + 0 * \frac{1}{5}\right) \\ &= 0.17 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} \leq 75) &= -\left(\frac{3}{5} \log_2\left(\frac{3}{5}\right) + \frac{2}{5} \log_2\left(\frac{2}{5}\right)\right) & \frac{|\text{Température} \leq 75|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{5}{5} = 1 \\ &= 0.97 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Température} > 75) &= 0 & \frac{|\text{Température} > 75|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= 0 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Température}, S_3 = 75) &= 0.97 - (0.97 * 1 + 0 * 0) \\ &= 0 \end{aligned}$$

Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

B. Calcul de gain d'humidité:

Humidité	70	80	80	80	96
N°	11	10	12	13	14
Jouer	Non	Non	Oui	Oui	Oui

Nous coupons entre x_{10} et $x_{12} \rightarrow S = 80$

$$\begin{aligned} \text{➤ } E(\text{Humidité} \leq 80) &= -\left(\frac{2}{4} \log_2\left(\frac{2}{4}\right) + \frac{2}{4} \log_2\left(\frac{2}{4}\right)\right) & \frac{|\text{Humidité} \leq 80|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{4}{5} \\ &= 1 \end{aligned}$$

$$\begin{aligned} \text{➤ } E(\text{Humidité} > 80) &= -\left(\frac{1}{1} \log_2(1) + 0 \log_2(0)\right) & \frac{|\text{Humidité} > 80|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} &= \frac{1}{5} \\ &= 0 \end{aligned}$$

D'où

$$\begin{aligned} \text{Gain}(\text{Humidité}, S = 80) &= 0.97 - \left(1 * \frac{4}{5} + 0 * \frac{1}{5}\right) \\ &= 0.17 \end{aligned}$$

C. Calcul de gain du Vent

1. Vent = Oui

$$\begin{cases} C_{\text{Jouer=Oui}} = \frac{0}{2} = 0 \\ C_{\text{Jouer=Non}} = \frac{2}{2} = 1 \end{cases} \quad \frac{|\text{Vent} = \text{Oui}|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{2}{5}$$

$$E(\text{Vent} = \text{Oui}) = -(0 \log_2(0) + 1 \log_2(1)) = 0$$

2. Vent = Non

$$\begin{cases} C_{\text{Jouer=Oui}} = \frac{3}{3} = 1 \\ C_{\text{Jouer=Non}} = 0 \end{cases} \quad \frac{|\text{Vent} = \text{Non}|}{|C_{\text{Jouer=Oui}}, C_{\text{Jouer=Non}}|} = \frac{3}{5}$$

$$E(\text{Vent} = \text{Non}) = -(1 \log_2(1) + 0 \log_2(0)) = 0$$

D'où

$$\begin{aligned} \text{Gain}(\text{Vent}) &= 0.97 - \left(0 * \frac{3}{5} + 0 * \frac{2}{5}\right) \\ &= 0.97 \end{aligned}$$

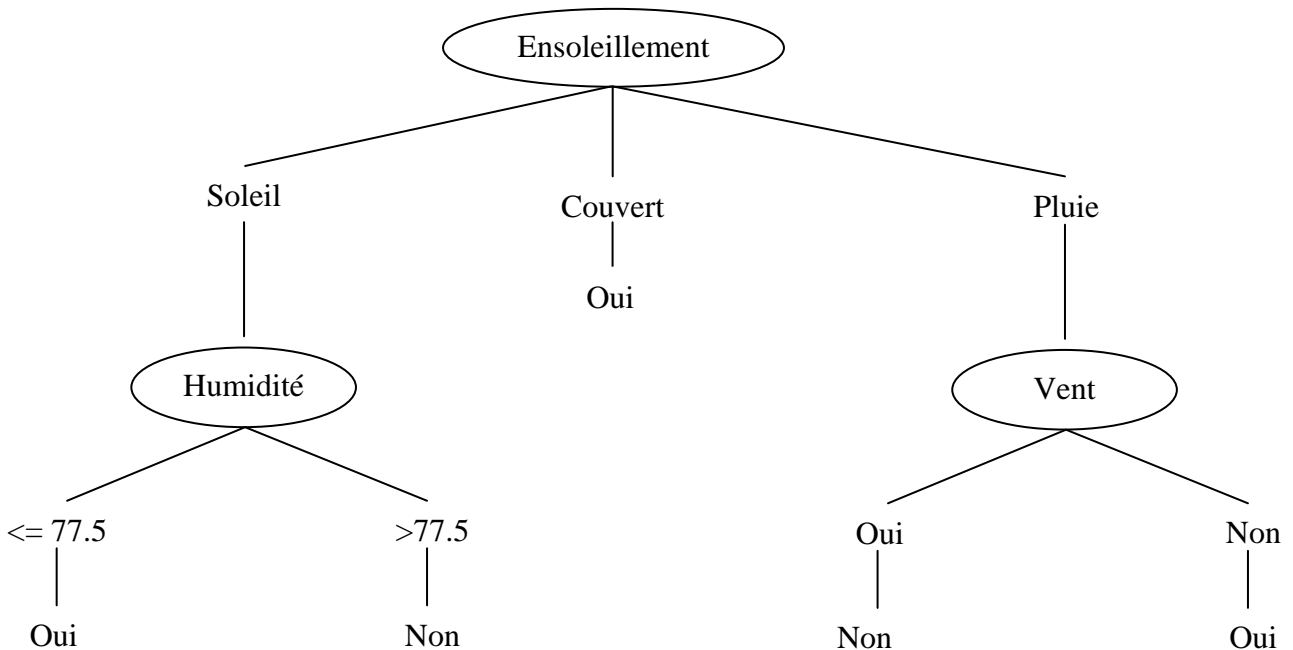
Module: Data Mining
 Nature de document: TD1-Solution
 Niveau: L3-STID

Année: 2022-2023
 N.BERMAD

Débriefing : Nous avons:

$\left\{ \begin{array}{l} \text{Gain (Vent)} = 0.97 \\ \text{Gain (Humidité, S)} = 0.17 \rightarrow \text{ nous choisissons vent ayant un gain maximal} \\ \text{Gain (Température, S}_2) = 0.17 \end{array} \right.$

L'arbre final est:



3. Ici la précision est calculée sur l'ensemble d'apprentissage et non pas de test.

Pour cela, nous devons construire la matrice de confusion:

Classe actuelle	Classe prédite	
	Jouer=Oui	Jouer=Non
Jouer=Oui	F11=9	F10=00
Jouer=Non	F01=00	F00 =5

Nb: $F_{ij}=c$ est le nombre des instances de la classe i prédite dans la classe j .

La prédiction de la classe des données d'apprentissages avec lesquelles on a construit l'arbre de décision.

Donc :

$$\begin{aligned} \text{Précision} &= \frac{\text{Nombre des prédictions correctes}}{\text{Nombre totale des prédictions}} \\ &= \frac{F11 + F00}{F11 + F10 + F01 + F00} \\ &= \frac{9 + 5}{9 + 0 + 0 + 5} \\ &= \frac{14}{14} \\ &= (1 * 100)\% \\ &= 100\% \end{aligned}$$

Une précision de 100% sur l'ensemble d'apprentissage illustre bien le phénomène de surapprentissage qui ne signifie pas forcément que le modèle obtenu généralise bien sur des nouvelles données.