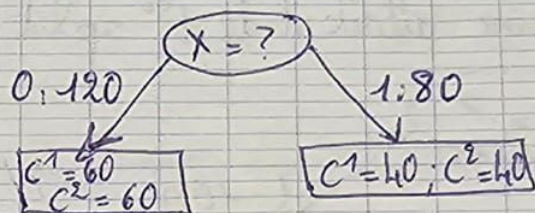


Nom : Sghiri Samir } TP Data Mining

Ex 7:

a) X comme premier attribut pour le Split:



⇒ D'où le tableau de contingence associé à X :

| Class \ X | 0 | 1 | Total |
|----------------|-----|----|-------|
| C ¹ | 60 | 40 | 100 |
| C ² | 60 | 40 | 100 |
| Total | 120 | 80 | 200 |

→ Calculons d'abord l'erreur avant Split:

$$\begin{aligned}
 \text{Erreur}_{\text{origine}} &= 1 - \max \{ P(C^1), P(C^2) \} \\
 &= 1 - \max \left\{ \frac{100}{200}, \frac{100}{200} \right\} \\
 &= 0,5
 \end{aligned}$$

Notons bien que : $\Delta(X) = \text{Erreur}_{\text{origine}} - \text{Erreur}(X)_{\text{split}}$

→ Calculons donc $\text{Erreur}(X)_{\text{split}}$:

$$\text{Erreur}(X)_{\text{split}} = \frac{120}{200} \cdot \text{Erreur}(X=0) + \frac{80}{200} \cdot \text{Erreur}(X=1)$$

$$\text{Erreur}(X=0) = 1 - \max \left\{ \frac{60}{120}, \frac{60}{120} \right\} = 0,5$$

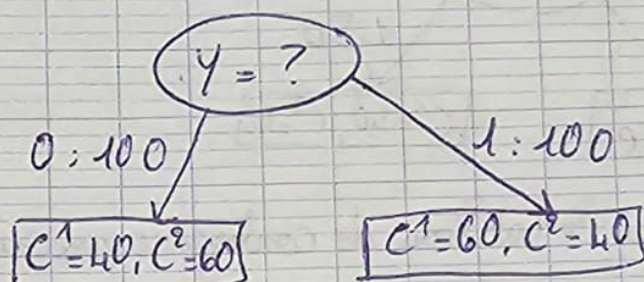
$$\text{Erreur}(X=1) = 1 - \max \left\{ \frac{40}{80}, \frac{40}{80} \right\} = 0,5$$

Page 2

$$\text{Donc : } \text{Erreur}(X)_{\text{split}} = (0,6 \cdot 0,5) + (0,4 \cdot 0,5) = 0,5$$

D'où $\Delta(X) = 0 \Rightarrow$ Aucune amélioration

$\rightarrow Y$ comme attribut du 1^{er} split



D'où le tableau de contingence associé à Y :

| Class \ Y | 0 | 1 | Total |
|-----------|-----|-----|-------|
| C¹ | 40 | 60 | 100 |
| C² | 60 | 40 | 100 |
| Total | 100 | 100 | 200 |

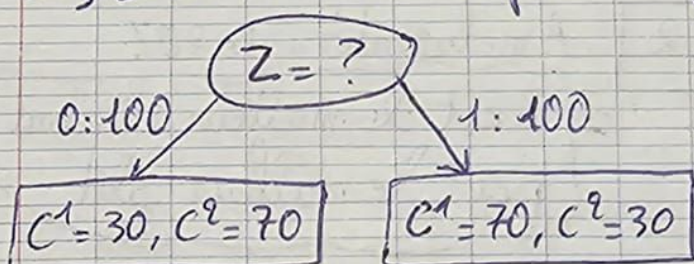
$$\begin{aligned} \text{Erreur}(Y=0) &= 1 - \max\{P(C^1), P(C^2)\} \\ &= 1 - \max\left\{\frac{40}{100}, \frac{60}{100}\right\} \\ &= 0,4 \end{aligned} \quad \left| \quad \begin{aligned} \text{Erreur}(Y=1) &= 1 - \max\{P(C^1), P(C^2)\} \\ &= 1 - \max\left\{\frac{60}{100}, \frac{40}{100}\right\} \\ &= 0,4 \end{aligned} \right.$$

$$\begin{aligned} \text{Erreur}(Y)_{\text{split}} &= \frac{100}{200} \cdot \text{Erreur}(Y=0) + \frac{100}{200} \cdot \text{Erreur}(Y=1) \\ &= (0,5 \cdot 0,4) + (0,5 \cdot 0,4) = 0,4 \end{aligned}$$

D'où: $\Delta(Y) = E_{\text{neur origine}} - E_{\text{neur}(Y)}|_{\text{split}}$
 $= 0,5 - 0,4$

$\Rightarrow \Delta(Y) = 0,1$ Amélioration de 10% au niveau
 du taux d'erreur de Classification
 $\approx Y \text{ est mieux que } X$

$\rightarrow Z$ comme attribut pour le premier split



D'où le tableau de contingence associé à l'attribut Z:

| Class \ Z | 0 | 1 | Total |
|----------------|-----|-----|-------|
| C ¹ | 30 | 70 | 100 |
| C ² | 70 | 30 | 100 |
| Total | 100 | 100 | 200 |

$$E_{\text{neur}}(Z=0) = 1 - \max \{P(C^1), P(C^2)\}$$

$$= 1 - \max \left\{ \frac{30}{100}, \frac{70}{100} \right\}$$

$$= 0,3$$

$$E_{\text{neur}}(Z=1) = 1 - \max \{P(C^1), P(C^2)\}$$

$$= 1 - \max \left\{ \frac{70}{100}, \frac{30}{100} \right\}$$

$$= \frac{30}{100} = 0,3$$

Page 4

$$E_{\text{meu}}(Z)_{\text{split}} = \frac{100}{200} \cdot E_{\text{meu}}(Z=0) + \frac{100}{200} \cdot E_{\text{meu}}(Z=1)$$

$$= (0,5 \cdot 0,3) + (0,5 \cdot 0,3)$$

$$= 0,3$$

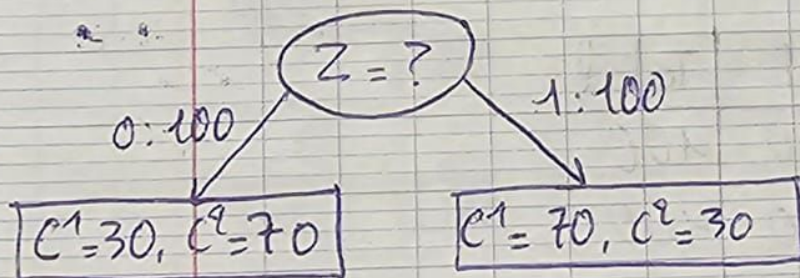
D'où :

$$\Delta(Z) = E_{\text{meu}}_{\text{origine}} - E_{\text{meu}}(Z)_{\text{split}}$$

$$= 0,5 - 0,3$$

$$\Rightarrow \boxed{\Delta(Z) = 0,2}$$
 , Amélioration de 20%
 $\approx \left\{ \begin{array}{l} Z \text{ est le meilleur attribut} \\ \text{pour le meilleur split} \end{array} \right.$

Donc, l'arbre de décision après le 1^{er} split est :

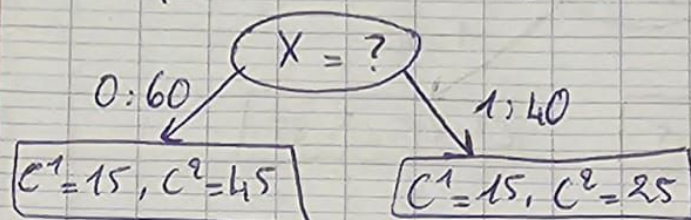


\Rightarrow 2^{ème} Split

→ le sous-Tableau obtenu sous l'hypothèse $Z=0$:

| X | Y | C^1 | C^2 |
|---|---|-------|-------|
| 0 | 0 | 5 | 40 |
| 0 | 1 | 10 | 5 |
| 1 | 0 | 10 | 5 |
| 1 | 1 | 5 | 20 |

→ X pour le deuxième split, sachant $Z=0$
~~Eneur(X=0) | split, z=0 = 1 - max~~



$$\begin{aligned} \text{Eneur}(X=0) |_{z=0} &= 1 - \max \{ P(C^1), P(C^2) \} \\ &= 1 - \max \left\{ \frac{15}{60}, \frac{45}{60} \right\} \\ &= \frac{1}{4} \end{aligned}$$

$$\begin{aligned} \text{Eneur}(X=1) |_{z=0} &= 1 - \max \left\{ \frac{15}{40}, \frac{25}{40} \right\} \\ &= \frac{3}{8} \end{aligned}$$

$$\begin{aligned} \text{Eneur}(X) |_{\text{split}} |_{z=0} &= \frac{60}{100} \cdot \text{Eneur}(X=0) |_{z=0} + \frac{40}{100} \cdot \text{Eneur}(X=1) |_{z=0} \\ &= 0,6 \cdot \frac{2}{8} + 0,4 \cdot \frac{3}{8} \\ &= 0,3 \end{aligned}$$

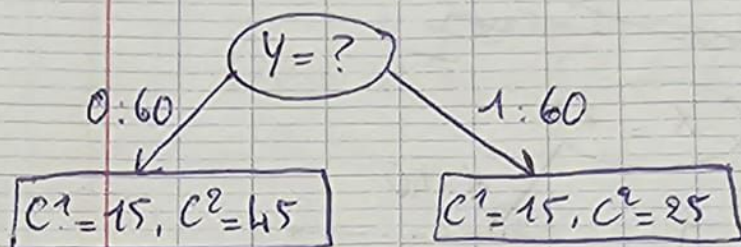
$$\text{Eneur}_{\text{origine}} |_{z=0} = \text{Eneur}(Z=0) = 0,3$$

$$\begin{aligned} \text{D'où } \Delta(X) |_{z=0} &= \text{Eneur}(Z=0) - \text{Eneur}(X) |_{\text{split}} |_{z=0} \\ &= 0,3 - 0,3 \end{aligned}$$

$$\Rightarrow \boxed{\Delta(X) |_{z=0} = 0} : \text{Aucune amélioration}$$

Page 6

→ Y pour le 2^{ème} split sachant $z=0$



$$\begin{aligned} E_{\text{neur}}(Y=0)_{z=0} &= 1 - \max\{P(C^1), P(C^2)\} \\ &= 1 - \max\left\{\frac{15}{60}, \frac{45}{60}\right\} \\ &= \frac{1}{4} \end{aligned}$$

$$\begin{aligned} E_{\text{neur}}(Y=1)_{z=0} &= 1 - \max\left\{\frac{15}{40}, \frac{25}{40}\right\} \\ &= \frac{3}{8} \end{aligned}$$

$$\begin{aligned} E_{\text{neur}}(Y)_{\text{split}}_{z=0} &= \left(\frac{60}{100} \cdot \frac{2}{8}\right) + \left(\frac{40}{100} \cdot \frac{3}{8}\right) \\ &= 0,3 \end{aligned}$$

D'où

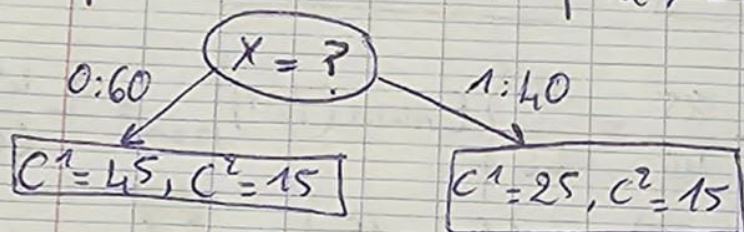
$$\begin{aligned} \Delta(Y)_{z=0} &= E_{\text{neur}}(Z=0) - E_{\text{neur}}(Y)_{\text{split}}_{z=0} \\ &= 0,3 - 0,3 \end{aligned}$$

$$\Rightarrow \boxed{\Delta(Y)_{z=0} = 0} \text{, Aucune Amélioration}$$

→ le sous-Tableau obtenu sous l'hypothèse $z=1$

| X | Y | C^1 | C^2 |
|---|---|-------|-------|
| 0 | 0 | 0 | 15 |
| 0 | 1 | 45 | 0 |
| 1 | 0 | 25 | 0 |
| 1 | 1 | 0 | 15 |

→ X pour le deuxième split, Sachant $z=1$



$$\begin{aligned}
 E_{\text{neur}}(X=0)_{z=1} &= 1 - \max \{ P(C^1), P(C^2) \} \\
 &= 1 - \max \left\{ \frac{45}{60}, \frac{15}{60} \right\} \\
 &= \frac{1}{4}
 \end{aligned}$$

$$\begin{aligned}
 E_{\text{neur}}(X=1)_{z=1} &= 1 - \max \left\{ \frac{25}{40}, \frac{15}{40} \right\} \\
 &= \frac{15}{40}
 \end{aligned}$$

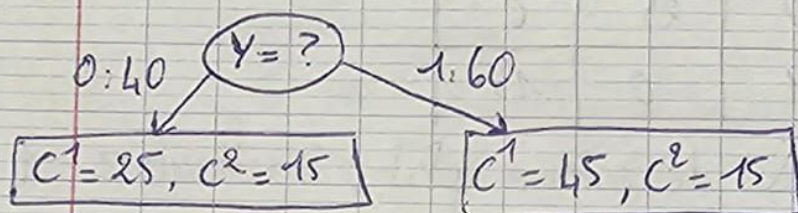
$$\begin{aligned}
 E_{\text{neur}}(X)_{\text{split}}_{z=1} &= \frac{60}{100} \cdot E_{\text{neur}}(X=0)_{z=1} + \frac{40}{100} \cdot E_{\text{neur}}(X=1)_{z=1} \\
 &= \left(\frac{60}{100} \cdot \frac{1}{4} \right) + \left(\frac{40}{100} \cdot \frac{15}{40} \right) \\
 &= 0,3
 \end{aligned}$$

Page 8

D'où : $\Delta(X)|_{z=1} = \text{Erreur}(Z=1) - \text{Erreur}(X)|_{\text{split}}|_{z=1}$
 $= 0,3 - 0,3$

$\Delta(X)|_{z=1} = 0$: Aucune Amélioration

→ Y pour le deuxième split, sachant $z=1$



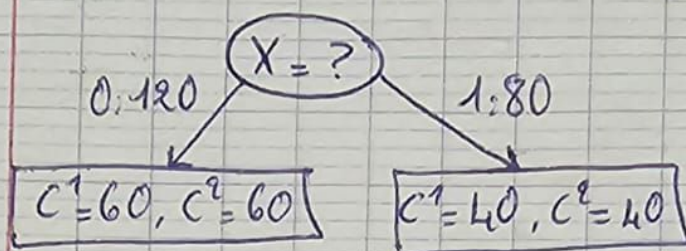
$$\text{Erreur}(Y=0)|_{z=1} = 1 - \frac{25}{40} = \frac{15}{40} \quad \left\{ \quad \text{Erreur}(Y=1)|_{z=1} = 1 - \frac{45}{60} = \frac{15}{60} \right.$$

$$\text{Erreur}(Y)_{\text{split}}|_{z=1} = \left(\frac{40}{100} \cdot \frac{15}{40} \right) + \left(\frac{60}{100} \cdot \frac{15}{60} \right)$$
$$= 0,3$$

$\Delta(Y)|_{z=1} = 0,3 - 0,3 = 0$: Aucune Amélioration

Donc : Le Taux d'erreur global de l'arbre induit est 30%

b) Utilisons X comme attribut pour le premier split
Donc, l'arbre de décision après le premier split est,

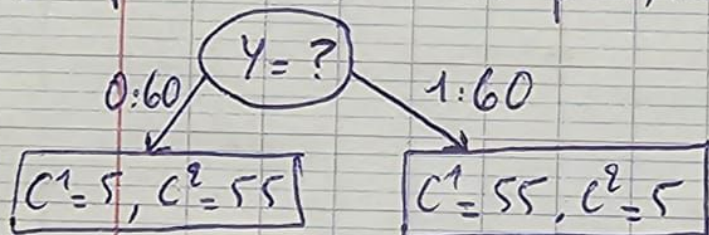


→ Passons au 2^{ème} split

→ le sous-tableau obtenu sous l'hypothèse $X=0$

| Y | Z | C^1 | C^2 |
|-----|-----|-------|-------|
| 0 | 0 | 5 | 40 |
| 0 | 1 | 0 | 15 |
| 1 | 0 | 10 | 5 |
| 1 | 1 | 45 | 0 |

→ Y pour le deuxième split, sachant $X=0$



$$E_{\text{neur}}(X=0)_{X=0} = 1 - \max \left\{ \frac{5}{60}, \frac{55}{60} \right\} \quad E_{\text{neur}}(X=1)_{X=0} = 1 - \max \left\{ \frac{55}{60}, \frac{5}{60} \right\}$$

$$= \frac{1}{12} \quad = \frac{1}{12}$$

Page 10

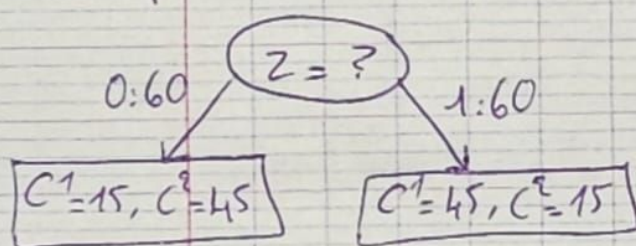
$$E_{\text{neur}}(Y)_{\text{split}}|_{x=0} = \frac{60}{120} \cdot E_{\text{neur}}(Y=0)|_{x=0} + \frac{60}{120} \cdot E_{\text{neur}}(Y=1)|_{x=0}$$
$$= \left(\frac{1}{2} \cdot \frac{1}{12}\right) + \left(\frac{1}{2} \cdot \frac{1}{2}\right) = \frac{1}{12} \approx 0,08$$

$$E_{\text{neur}}_{\text{origine}}|_{x=0} = E_{\text{neur}}(X=0) = 0,5$$

$$\Rightarrow \Delta(Y)|_{x=0} = E_{\text{neur}}(X=0) - E_{\text{neur}}(Y)_{\text{split}}|_{x=0}$$
$$= 0,5 - 0,08$$

$$\Rightarrow \boxed{\Delta(Y)|_{x=0} = 0,42}, \text{ Amélioration de } 42\%$$

→ Z pour le 2^{ème} split, sachant $x=0$



$$E_{\text{neur}}(Z=0)|_{x=0} = 1 - \max\left\{\frac{15}{60}, \frac{45}{60}\right\}$$
$$= \frac{15}{60} = \frac{1}{4}$$

$$E_{\text{neur}}(Z=1)|_{x=0} = 1 - \max\left\{\frac{45}{60}, \frac{15}{60}\right\}$$
$$= \frac{1}{4}$$

$$E_{\text{neur}}(Z)_{\text{split}}|_{x=0} = \left(\frac{60}{120} \cdot \frac{1}{4}\right) + \left(\frac{60}{120} \cdot \frac{1}{4}\right)$$
$$= \frac{1}{4} = 0,25$$

Page 11

$$\Delta(Z)_{x=0} = 0,5 - 0,25 = 0,25$$

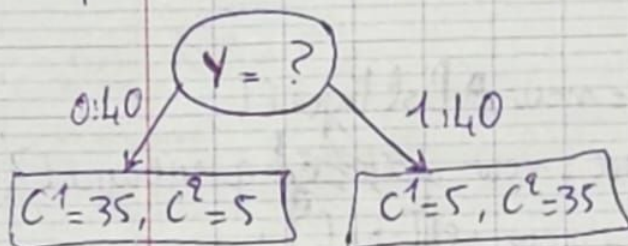
\Rightarrow Amélioration de 25%

\Rightarrow Donc Y est mieux suivant $X=0$

\rightarrow le sous-tableau obtenu sous l'hypothèse $X=1$

| Y | Z | C^1 | C^2 |
|-----|-----|-------|-------|
| 0 | 0 | 10 | 5 |
| 0 | 1 | 25 | 0 |
| 1 | 0 | 5 | 20 |
| 1 | 1 | 0 | 15 |

$\rightarrow Y$ pour la 2^{ème} split, sachant $X=1$



$$E_{\text{meur}}(Y=0)_{X=1} = 1 - \max \left\{ \frac{35}{40}, \frac{5}{40} \right\}$$

$$= \frac{1}{8}$$

$$E_{\text{meur}}(Y=1)_{X=1} = 1 - \max \left\{ \frac{5}{40}, \frac{35}{40} \right\}$$

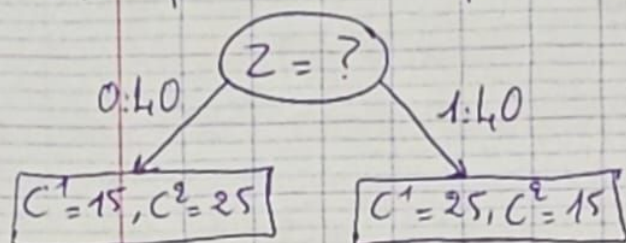
$$= \frac{1}{8}$$

$$E_{\text{meur}}(Y)_{\text{split}}_{X=1} = \left(\frac{40}{80} \cdot \frac{1}{8} \right) + \left(\frac{40}{80} \cdot \frac{1}{8} \right) = \frac{1}{8} \approx 0,125$$

$$\Delta(Y)_{X=1} = 0,5 - 0,125 = 0,375 \Rightarrow \text{Donc Amélioration de } 37,5\%$$

Page 12

→ Z pour le 2^{ème} split, sachant $x=1$



$$\text{Erreur}(Z=0)_{x=1} = \frac{1}{15} \max \left\{ \frac{15}{40}, \frac{25}{40} \right\}$$
$$= \frac{15}{40}$$

$$\text{Erreur}(Z=1)_{x=1} = \frac{1}{15} \max \left\{ \frac{25}{40}, \frac{15}{40} \right\}$$
$$= \frac{15}{40}$$

$$\text{Erreur}(Z)_{\text{split}}_{x=1} = \frac{40}{80} \left(2 \cdot \frac{5}{40} \right)$$
$$= \frac{1}{2} \cdot 2 \cdot \frac{15}{40} = \frac{15}{40} = 0,375$$

$$\Delta(Z)_{x=1} = 0,5 - 0,375$$
$$= 0,125 \approx 12,5\% \text{ d'amélioration}$$

Donc y est mieux suivant $x=1$

D'où le taux d'erreur de l'arbre induit est 37,5%

- c) L'heuristique gourmande, bien qu'efficace pour la construction d'arbres de décision en minimisant le taux d'erreur à chaque étape, n'est pas exempte de limitations. Comme le montre le résultat obtenu dans la question (b), le choix initial de l'attribut peut avoir un impact significatif sur le taux d'erreur global de l'arbre. En effet, si l'attribut sélectionné n'est pas optimal, cela peut entraîner une augmentation du taux d'erreur final.