

# ML - Assignment

①

Name: Ch. Keerthana

Reg no: 22BCE9635

slot: G1

1) Construct Decision tree using ID3 Algorithm.

age	income	student	credit rating	Buy Computer
<=30	high	no	fair	no
<=30	high	no	excellent	no
31-40	high	no	fair	yes
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
31-40	low	yes	excellent	yes
<=30	medium	no	fair	no
<=30	low	yes	fair	yes
>40	medium	yes	fair	yes
<=30	medium	yes	excellent	yes
31-40	medium	no	excellent	yes
31-40	high	yes	fair	yes
>40	medium	no	excellent	no

Attribute : Age

$$\text{Entropy}(S) = - \underbrace{\frac{9}{14} \log_2 \left( \frac{9}{14} \right)}_{\text{"yes" probability}} - \underbrace{\frac{5}{14} \log_2 \left( \frac{5}{14} \right)}_{\text{"No" probability}}$$

(for Buys Computer)

$$\text{Entropy}(S) = - \frac{9}{14} \log_2 \left( \frac{9}{14} \right) - \frac{5}{14} \log_2 \left( \frac{5}{14} \right) = 0.94$$

- age  $< 30$  ( 2 yes and 3 no)
- age 31-40 ( 4 yes and 0 no)
- age  $> 40$  ( 3 yes & 2 no)

Entropy (Age):

(i)  $< 30$

$$\begin{aligned}\text{Entropy} &= \frac{5}{14} \left[ -\frac{2}{5} \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \log_2 \left( \frac{3}{5} \right) \right] \\ &= \frac{5}{14} (0.9709)\end{aligned}$$

(ii) ages 31-40

Entropy =  $\frac{4}{14} (0)$  because, for age 31-40 only 'yes' was there so, entropy value will be '0'.

(iii) age:  $> 40$

$$\begin{aligned}\text{Entropy} &= \frac{5}{14} \left[ -\frac{3}{5} \log_2 \left( \frac{3}{5} \right) - \frac{2}{5} \log_2 \left( \frac{2}{5} \right) \right] \\ &= \frac{5}{14} (0.9709)\end{aligned}$$

$$\begin{aligned}\text{Entropy (age)} &= \frac{5}{14} (0.9709) + 0 + \frac{5}{14} (0.9709) \\ &= 0.6935\end{aligned}$$

$$\begin{aligned}\therefore \text{Gain (age)} &= \underset{\substack{\uparrow \\ \text{total}}}{\text{Entropy (s)}} - \text{Entropy (age)} \\ &= 0.94 - 0.6935 = 0.2465\end{aligned}$$

Attribute : Income

$$\begin{aligned}\text{Entropy (s)} &= 0.94 \\ \downarrow \\ \text{whole data set}\end{aligned}$$

- Income (high)  $\rightarrow$  2 yes & 2 no
- Income (medium)  $\rightarrow$  4 yes & 2 no
- Income (low)  $\rightarrow$  3 yes & 1 no

$$\text{Entropy (Income)} = \frac{4}{14} \left[ -\frac{2}{4} \log_2 \left( \frac{2}{4} \right) - \frac{2}{4} \log_2 \left( \frac{2}{4} \right) \right] \text{ (high)} \\ + \frac{6}{14} \left[ -\frac{4}{6} \log_2 \left( \frac{4}{6} \right) - \frac{2}{6} \log_2 \left( \frac{2}{6} \right) \right] + \frac{4}{14} \left[ -\frac{3}{4} \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right) \right] \\ \hookrightarrow \text{(medium)}$$

$$= \frac{4}{14} (1) + \frac{6}{14} (0.918) + \frac{4}{14} (0.811)$$

$$= 0.285714 + 0.393428 + 0.231429 = 0.9108$$

$$\text{Gain (Income)} = \text{Entropy (S)} - \text{Entropy (Income)} \\ = 0.94 - 0.9108 = 0.0292$$

Attribute : Student

- Student (yes)  $\Rightarrow$  (6 yes & 1 NO)
- Student (NO)  $\Rightarrow$  (3 yes & 4 NO)

$$\text{Entropy (Student)} = \frac{7}{14} \left[ -\frac{6}{7} \log_2 \left( \frac{6}{7} \right) - \frac{1}{7} \log_2 \left( \frac{1}{7} \right) \right] + \\ \frac{4}{14} \left[ -\frac{3}{7} \log_2 \left( \frac{3}{7} \right) - \frac{4}{7} \log_2 \left( \frac{4}{7} \right) \right] \\ = \frac{7}{14} (0.5916) + \frac{4}{14} (0.9852) \\ = 0.2958 + 0.4926 = 0.7884$$

$$\therefore \text{Gain (Student)} = 0.94 - 0.7884 = 0.1516.$$

Attribute : Credit-Rating

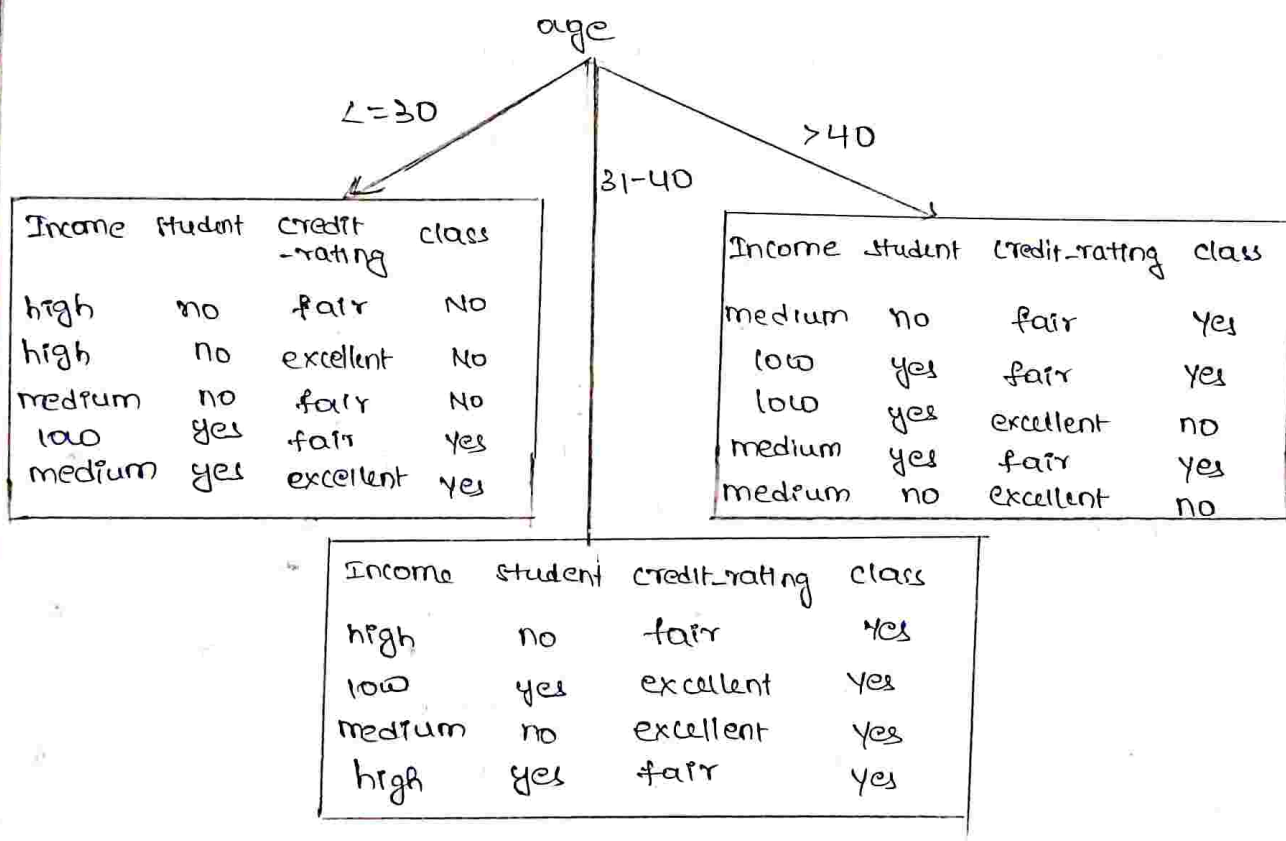
- Credit-rating (fair)  $\Rightarrow$  (6 yes & 2 NO)
- Credit-rating (excellent)  $\Rightarrow$  (3 yes & 3 NO)

$$\text{Entropy (Credit-rating)} = \frac{8}{14} \left[ -\frac{6}{8} \log_2 \left( \frac{6}{8} \right) - \frac{2}{8} \log_2 \left( \frac{2}{8} \right) \right] + \frac{6}{14} \left[ -\frac{3}{6} \log_2 \left( \frac{3}{6} \right) - \frac{3}{6} \log_2 \left( \frac{3}{6} \right) \right] \\ = \frac{8}{14} (0.8112) + \frac{6}{14} (1) = 0.4635 + 0.4285 \\ = 0.8920$$

$$\text{Gain (Credit-rating)} = 0.94 - 0.8920 \\ = 0.048$$

By Comparing Information Gain for Age, Student, Income and Credit-rating, Age is having more Gain so, we consider Age as root node and start constructing the tree.

⇒ under Age (root node) we have 3 more conditions i.e.,  $\leq 20$ ,  $31-40$  &  $> 40$  so, we construct subsets for those.



Attribute: Income

$$E(S_{age \leq 20}) = E(2/3) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) = 0.97$$

$$\text{Entropy}(\text{Income}) = \frac{2}{5}(0) + \frac{2}{5}(1) + \frac{1}{5}(0) = \frac{2}{5}(1) = 0.4$$

$$\text{Gain}(\text{Income}) = 0.97 - 0.4 = 0.57$$

Attribute: Student

$$E(S_{age \leq 20}) = E(2/3) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) = 0.97$$

$$\text{Entropy}(\text{Student}) = \frac{2}{5}(0) + \frac{3}{5}(0) = 0$$

$$\text{Gain}(\text{Student}) = (0.97 - 0) = 0.97$$



Attribute (credit-rating)

$$\text{Entropy (credit)} = \frac{3}{5} \left( -\frac{1}{3} \log_2 \left( \frac{1}{3} \right) - \frac{2}{3} \log_2 \left( \frac{2}{3} \right) \right) + \frac{2}{5} (1) \\ = 0.9508$$

$$\text{Gain (Credit rating)} = 0.97 - 0.9508 = 0.0192$$

for Age > 40

$$E(S_{\text{Age} > 40}) = E(3,2) = -\frac{3}{5} \log_2 \left( \frac{3}{5} \right) - \frac{2}{5} \log_2 \left( \frac{2}{5} \right) = 0.97$$

$$\text{Entropy (Income)} = \frac{3}{5} \left[ -\frac{2}{3} \log_2 \left( \frac{2}{3} \right) - \frac{1}{3} \log_2 \left( \frac{1}{3} \right) \right] + \frac{2}{5} (1) \\ = \frac{3}{5} (0.9182) + \frac{2}{5} = 0.95$$

$$\text{Gain (Income)} = (0.97 - 0.95) = 0.02$$

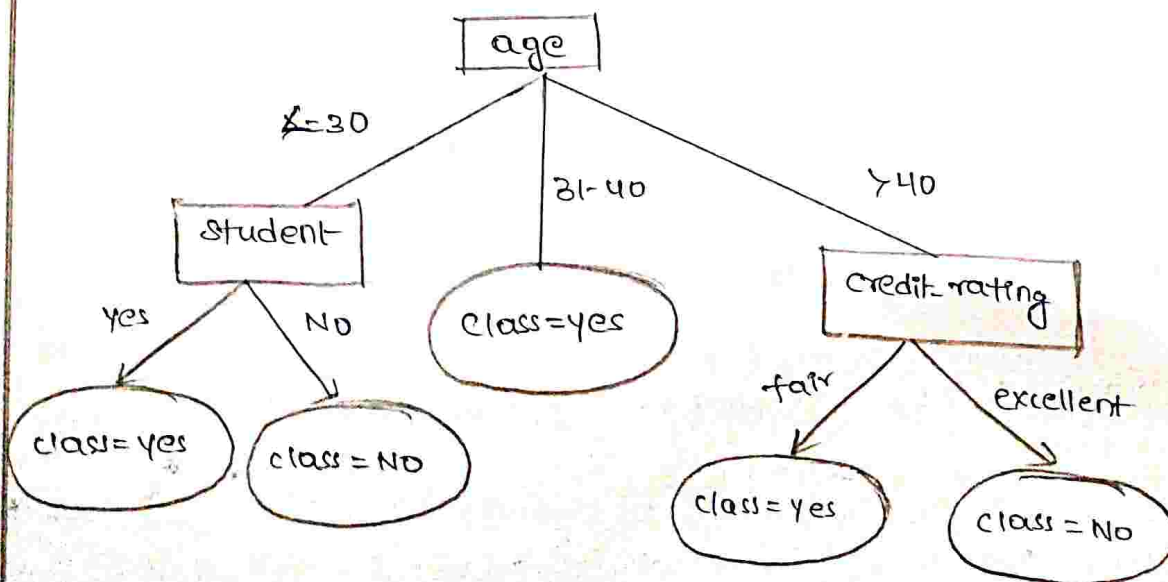
$$\text{Entropy (student)} = \frac{3}{5} \left[ -\frac{2}{3} \log_2 \left( \frac{2}{3} \right) - \frac{1}{3} \log_2 \left( \frac{1}{3} \right) \right] + \frac{2}{5} (1) = 0.95$$

$$\text{Gain (student)} = (0.97 - 0.95) = 0.02$$

$$\text{Entropy (credit rating)} = \frac{3}{5} (0) + \frac{2}{5} (0) = 0$$

$$\text{Gain (credit rating)} = (0.97 - 0) = 0.97$$

Final decision trees:



2) Construct Decision tree using CART algorithm

CGPA	Inter active	practical knowledge	Common skills	Job offer
$\geq 9$	Yes	Very good	Good	Yes
$\geq 8$	No	Good	Moderate	Yes
$\geq 9$	No	Average	poor	No
$< 8$	No	Average	Good	No
$\geq 8$	Yes	Good	Moderate	Yes
$\geq 9$	Yes	Good	Moderate	Yes
$< 8$	Yes	Good	poor	No
$\geq 9$	No	Very good	Good	Yes
$\geq 8$	Yes	Good	Good	Yes
$\geq 8$	Yes	Average	Good	Yes

Step-1: calculate the Gini-index for dataset.

$$\text{Gini-index}(T) = 1 - \sum_{i=1}^n P_i^2$$

$$\begin{aligned} \text{Gini-index}(T) &= 1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2 = 1 - 0.49 - 0.09 = 1 - 0.58 \\ &= 0.42 \quad (\text{Whole dataset}) \end{aligned}$$

Step-2: compute Gini-index for each attribute and each of the subset in the attribute.

i) CGPA: ( $\geq 9$ ,  $\geq 8$ ,  $< 8$ )

possible subsets

$\{ \}, \{ \geq 9 \}, \{ \geq 8 \}, \{ \geq 9, \geq 8 \}, \{ \geq 9, < 8 \}, \{ \geq 8, < 8 \}, \{ \geq 9, \geq 8, < 8 \}$

calculation of best splitting subset:

$$\text{Gini-index}(T, A) = \left| \frac{S_1}{T} \right| \text{Gini}(S_1) + \left| \frac{S_2}{T} \right| \text{Gini}(S_2)$$

$$\text{Gini-index}(T, \text{CGPA} \in \{ \geq 9, \geq 8 \}) = 1 - \left(\frac{7}{10}\right)^2 - \left(\frac{3}{10}\right)^2 = 1 - 0.49 - 0.09 = 0.42$$

$$\text{Gini-index}(T, \text{CGPA} \in \{ < 8 \}) = 1 - \left(\frac{2}{2}\right)^2 - \left(\frac{2}{2}\right)^2 = 1 - 1 = 0$$

$$\text{Gini-Index}(T, \text{CGPA} \in \{\geq 9, \leq 8\}, \{ \leq 8 \}) = (8/10)(0.2194) + (2/10)(0) = 0.17552$$

(ii)

$$\text{Gini-Index}(T, \text{CGPA} \in \{\geq 9, \leq 8\}) = 1 - (3/6)^2 - (3/6)^2 = (1 - 0.5) = 0.5$$

$$\text{Gini-Index}(T, \text{CGPA} \in \{\geq 8\}) = 1 - (4/4)^2 - (0/4)^2 = (1 - 1) = 0$$

$$\text{Gini-Index}(T, \text{CGPA} \in \{\geq 9, \leq 8\}, \{\geq 9, \leq 8\}) = (6/10)0.5 + (4/10)0 = 0.3$$

(iii)

$$\text{Gini-Index}(T, \text{CGPA} \in \{\geq 8, \leq 8\}) = 1 - (4/6)^2 - (2/6)^2 = 1 - 0.555 = 0.445$$

$$\text{Gini-Index}(T, \text{CGPA} \in \{\geq 9\}) = 1 - (3/4)^2 - (1/4)^2 = 1 - 0.625 = 0.375$$

$$\text{Gini-Index}(T, \text{CGPA} \in \{\geq 8, \leq 8\}, \{\geq 9\}) = (6/10)0.445 + (4/10)0.375$$

$$\text{Step-3:} = 0.417$$

In the above 3 cases 1st case subset is having minimum Gini-Index value, so, 1st subset is the best splitting subset

→ The subset  $\text{CGPA} = \{\geq 9, \leq 8\}, \leq 8\}$  has low T value

Step-4:

Compute ( $\Delta \text{Gini}$  / best splitting subset) of that attribute.

$$\begin{aligned} \Delta \text{Gini}(\text{CGPA}) &= \text{Gini}(T) - \text{Gini}(T, \text{CGPA}) \\ &= (0.42 - 0.1755) = 0.2445 \end{aligned}$$

Repeat the same process for remaining attributes in dataset such as for interactivens, practical Knowledge and Communication skills.

Category Categories for interactivens:

$$\text{Gini-Index}(T, \text{interactivens} \in \{\text{yes}\}) = 1 - (5/6)^2 - (1/6)^2 = (1 - 0.72) = 0.28$$

$$\text{Gini-Index}(T, \text{intr} \in \{\text{No}\}) = 1 - (2/4)^2 - (2/4)^2 = (1 - 0.5) = 0.5$$

$$\begin{aligned} \text{Gini-Index}(T, \text{intr} \in \{\text{yes}, \text{No}\}) &= 6/10(0.28) + 4/10(0.5) \\ &= 0.168 + 0.2 = 0.368 \end{aligned}$$

$$\Delta \text{Gini}(\text{Intera} = \text{Gini}(T) - \text{Gini}(T, \text{intr}) = (0.42 - 0.368) = 0.052$$

-ctivens)



## Categories for practical knowledge

$$(i) \text{ Gini-Index}(T, \text{pract know} \in \{\text{very good}, \text{good}\}) = 1 - (4/7)^2 - (1/7)^2 = 0.2456$$

$$\text{Gini-Index}(T, \text{pract know} \in \{\text{avg}\}) = 1 - (1/3)^2 - (2/3)^2 = 1 - 0.555 = 0.445$$

$$\text{Gini-Index}(T, \text{pract know} \in \{\text{VG}, \text{G}, \text{avg}\}) = (7/10) \cdot 0.2456 + (3/10) \cdot 0.445 = 0.3054$$

$$(ii) \text{ G-I}(T, \text{pract know} \in \{\text{VG}, \text{avg}\}) = 1 - (3/5)^2 - (2/5)^2 = 1 - 0.52 = 0.48$$

$$\text{G-I}(T, \text{pract know} \in \{\text{good}\}) = 1 - (4/5)^2 - (1/5)^2 = 1 - 0.68 = 0.32$$

$$\text{G-I}(T, \text{pract know} \in \{\text{VG}, \text{avg}, \text{G}\}) = 5/10(0.48) + 5/10(0.32) = 0.40$$

(iii)

$$\text{G-I}(T, \text{pract know} \in \{\text{G}, \text{avg}\}) = 1 - (5/8)^2 - (3/8)^2 = 1 - 0.5312 = 0.4688$$

$$\text{G-I}(T, \text{pract know} \in \{\text{VG}\}) = 1 - (2/2)^2 - (0/2)^2 = (1-1) = 0$$

$$\text{G-I}(T, \text{pract know} \in \{\text{G}, \text{avg}, \text{VG}\}) = (8/10) \cdot 0.4688 + (2/10) \cdot 0 = 0.3750$$

case - (i) is having low Gini-Index value so, 1st subset is best splitting subset

$$\Delta \text{Gini}(\text{pract know}) = \text{Gini}(T) - \text{Gini}(T, \text{pract know})$$

$$= 0.42 - 0.3054 = 0.1146$$

## Categories for ~~Common~~ <sup>Common</sup> ~~connection~~ <sup>skill</sup>

$$(i) \text{ G-I}(T, \text{comsk} \in \{\text{G}, \text{Mod}\}) = 1 - (3/8)^2 - (1/8)^2 = 1 - 0.4606 = 0.2194$$

$$\text{G-I}(T, \text{comsk} \in \{\text{poor}\}) = 1 - (2/2)^2 - (0/2)^2 = (1-1) = 0$$

$$\text{G-I}(T, \text{comsk} \in \{\text{G}, \text{Mod}, \text{P}\}) = (8/10) \cdot 0.2194 + (2/10) \cdot 0 = 0.1955$$

$$(ii) \text{ G-I}(T, \text{comsk} \in \{\text{G}, \text{P}\}) = 1 - (4/7)^2 - (3/7)^2 = 1 - 0.5101 = 0.4899$$

$$\text{G-I}(T, \text{comsk} \in \{\text{Mod}\}) = 1 - (3/3)^2 - (0/3)^2 = (1-1) = 0$$

$$\text{G-I}(T, \text{comsk} \in \{\text{G}, \text{P}, \text{Mod}\}) = (7/10) \cdot 0.4899 + (3/10) \cdot 0 = 0.3429$$



(iii)

$$G-I(T, comix \in \{Mod, P\}) = 1 - (3/5)^2 - (2/5)^2 = 1 - 0.36 - 0.16 = 0.48$$

$$G-I(T, comix \in \{G\}) = 1 - (4/5)^2 - (1/5)^2 = 1 - 0.64 - 0.04 = 0.32$$

$$G-I(T, comix \in \{Mod, P\}, G\}) = (5/10) \cdot 0.48 + (5/10) \cdot 0.32 = 0.40$$

$$\Delta Gini(comix) = Gini(T) - Gini(T, comix) \quad (\text{1st subset} \rightarrow \text{best splitting subset})$$
$$= 0.42 - 0.1755 = 0.2445$$

Step-5:

After calculating Gini-index and  $\Delta Gini$  for all attributes we need to consider highest  $\Delta Gini$  value attribute as root node here, for CGPA & comix having highest  $\Delta Gini$  value. so, let us consider CGPA as root node

$\Rightarrow$  for subset  $\{ \geq 9, \geq 8 \}$  have both yes and No values so, for deleting that again calculate Gini-Index value

$$Gini-Index(T) = 1 - (7/8)^2 - (1/8)^2 = 1 - 0.766 - 0.0156 = 0.2184$$

(i) Interactiveness

$$G-I(T, intr \in \{yes\}) = 1 - (5/5)^2 - (0/5)^2 = 0$$

$$G-I(T, intr \in \{No\}) = 1 - (2/3)^2 - (1/3)^2 = 1 - 0.44 - 0.11 = 0.449$$

$$G-I(T, intr \in \{yes, No\}) = 2/8(0) + 1/8(0.449) = 0.056$$

$$\Delta Gini(Int) = Gini(T) - Gini(T, intr)$$
$$= (0.2184 - 0.056) = 0.1624$$

(ii) practical knowledge

$$G-I(T, intr \in \{V, G, G\}, G\}) = 0.125$$

$$\Delta Gini(Pract Know) = Gini(T) - Gini(T, Pract Know)$$
$$= 0.2184 - 0.125$$
$$= 0.0934$$

(iii) Communication skills

subsets

Gini-index

{ {G, Hod}, P }

0

{ {G, P}, Hod }

0.2

{ {Hod, P}, G }

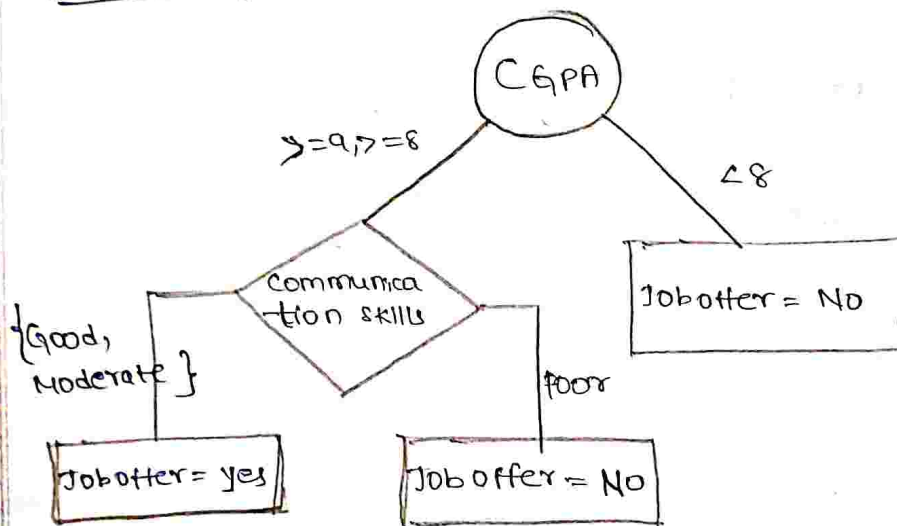
0.1875

$$\Delta \text{Gini}(\text{CommSK}) = \text{Gini}(T) - \text{Gini}(T, \text{CommSK})$$

$$= 0.2184 - 0 = 0.2184$$

∴ Communication skills having highest  $\Delta \text{Gini}$  value (high priority)

Decision trees:



3) Construct decision tree using C4.5 Algorithm

CGPA	Interactive	Practical Knowledge	Common Skills	Job offer
$\geq 9$	Yes	very good	Good	yes
$\geq 8$	No	Good	Moderate	yes
$\geq 9$	No	Average	poor	No
$< 8$	No	Average	Good	No
$\geq 8$	Yes	Good	Moderate	yes
$\geq 9$	Yes	Good	Moderate	yes
$< 8$	Yes	Good	poor	No
$\geq 9$	No	very good	Good	yes
$\geq 8$	Yes	Good	Good	yes
$\geq 8$	Yes	Average	Good	yes

Step-1:

calculate class-Entropy for target class 'job offer'.

$$\text{Entropy-Info}(T) = - \sum_{i=1}^n P_i \log_2 P_i$$

$$\begin{aligned} \text{Entropy-Info}(\text{job offer}) &= \left[ \frac{7}{10} \log_2 \left( \frac{7}{10} \right) + \frac{3}{10} \log_2 \left( \frac{3}{10} \right) \right] \\ &= (-0.3599 + (-0.5208)) = 0.8807 \end{aligned}$$

Step-2:

calculate the Entropy-Info, Gain (Info-gain), Split-Info, Gain-ratio for each attribute in training dataset.

$$\text{Entropy-Info}(T, A) = \sum_{i=1}^n \left| \frac{A_i}{T} \right| \times \text{Entropy-Info}(A_i)$$

(1) CGPA:

$$\begin{aligned} \text{Entropy-Info}(T, \text{CGPA}) &= \frac{4}{10} \left[ -\frac{3}{4} \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right) \right] + \frac{4}{10} \left[ -\frac{4}{4} \log_2 \left( \frac{4}{4} \right) \right] \\ &\quad + \frac{2}{10} \left[ -\frac{6}{2} \log_2 \left( \frac{6}{2} \right) - \frac{2}{2} \log_2 \left( \frac{2}{2} \right) \right] = \frac{4}{10} [0.3111 + 0.4999] + 0 + 0 \\ &= 0.3243 \end{aligned}$$



$$\text{Info\_Gain}(A) = \text{Entropy\_Info}(T) - \text{Entropy\_Info}(T, A)$$

$$\text{Gain}(CSPA) = (0.8807 - 0.3243) = 0.5564$$

$$\text{split\_info}(T, A) = - \sum_{i=1}^K \frac{|A_i|}{|T|} \times \log_2 \left( \frac{|A_i|}{|T|} \right)$$

$$\begin{aligned} \text{split\_info}(CSPA) &= -\frac{4}{10} \log_2(4/10) - \frac{4}{10} \log_2(4/10) - \frac{2}{10} \log_2(2/10) \\ &= (0.5285 + 0.5285 + 0.4641) = 1.5211 \end{aligned}$$

$$\text{Gain ratio}(CSPA) = \frac{0.5564}{1.5211} = 0.3658$$

(2) Interactiveness:

$$\begin{aligned} \text{Entropy\_Info}(T, \text{inter}) &= \frac{6}{10} \left[ -5/6 \log_2(5/6) - 1/6 \log_2(1/6) \right] + \frac{4}{10} \left[ -2/4 \log_2(2/4) - \frac{2}{4} \log_2(2/4) \right] \\ \text{interactive} &= \frac{6}{10} [0.291 + 0.4306] + \frac{4}{10} [0.4999 + 0.4999] \\ &= (0.3898 + 0.3998) = 0.7896 \end{aligned}$$

$$\text{Gain}(\text{Inter}) = (0.8807 - 0.7896) = 0.0911$$

$$\text{split\_info}(T, \text{inter}) = -\frac{6}{10} \log_2(6/10) - \frac{4}{10} \log_2(4/10) = 0.9704$$

$$\text{Gain ratio}(\text{inter}) = \frac{0.0911}{0.9704} = 0.0939$$

(3) Practical Knowledge:

$$\begin{aligned} \text{Entropy\_Info}(T, \text{pract kn}) &= \frac{2}{10} \left[ -2/2 \log_2(2/2) \right] + \frac{3}{10} \left[ -1/3 \log_2(1/3) - 2/3 \log_2(2/3) \right] \\ \text{practical knowledge} &\quad + \frac{5}{10} \left[ -4/5 \log_2(4/5) - 1/5 \log_2(1/5) \right] \\ &= 2/10(0) + 3/10(0.5280 + 0.3897) + 5/10(0.2574 + 0.4641) \\ &= 0 + 0.2573 + 0.3608 = 0.6361 \end{aligned}$$

$$\text{Gain}(\text{pract kn}) = 0.8807 - 0.6361 = 0.2448$$

$$\begin{aligned} \text{Split info}(T, \text{pract kn}) &= -\frac{2}{10} \log_2(2/10) - \frac{5}{10} \log_2(5/10) - \frac{3}{10} \log_2(3/10) \\ &= 1.4853 \end{aligned}$$

$$\text{Gain ratio}(\text{pract kn}) = \frac{0.2448}{1.4853} = 0.1648$$

(4)

Communication skills

$$\begin{aligned} \text{Entropy-info (T, comsk)} &= \frac{5}{10} \left[ -4/5 \log_2(4/5) - 1/5 \log_2(1/5) \right] + \\ &\quad \frac{3}{10} \left[ -3/3 \log_2(3/3) \right] + \frac{2}{10} \left[ -2/2 \log_2(2/2) \right] \\ &= 5/10 (0.5280 + 0.3897) + 3/10(0) + 2/10(0) \\ &= 0.3609 \end{aligned}$$

$$\text{Gain(comsk)} = (0.8813 - 0.3609) = 0.5202$$

$$\begin{aligned} \text{split-info (T, comsk)} &= -5/10 \log_2(5/10) - 3/10 \log_2(3/10) - 2/10 \log_2(2/10) \\ \text{common skill} &= 1.4853 \end{aligned}$$

$$\text{Gain-ratio (comsk)} = \left( \frac{0.5202}{1.4853} \right) = 0.3502$$

CGPA is root node

for CGPA ( $\geq 9$ )

$$\begin{aligned} \text{Entropy-info(job offer)} &= -3/4 \log_2(3/4) - 1/4 \log_2(1/4) \\ &= 0.3112 + 0.5 = 0.8112 \end{aligned}$$

$$\text{Entropy-info(T, inter)} = 2/4 \left[ -2/2 \log_2(2/2) \right] + 2/4 \left[ 2 \left( \frac{1}{2} \right) \log_2(1/2) \right] = 0.4997$$

$$\text{Gain(inter)} = 0.8112 - 0.4997 = 0.3111$$

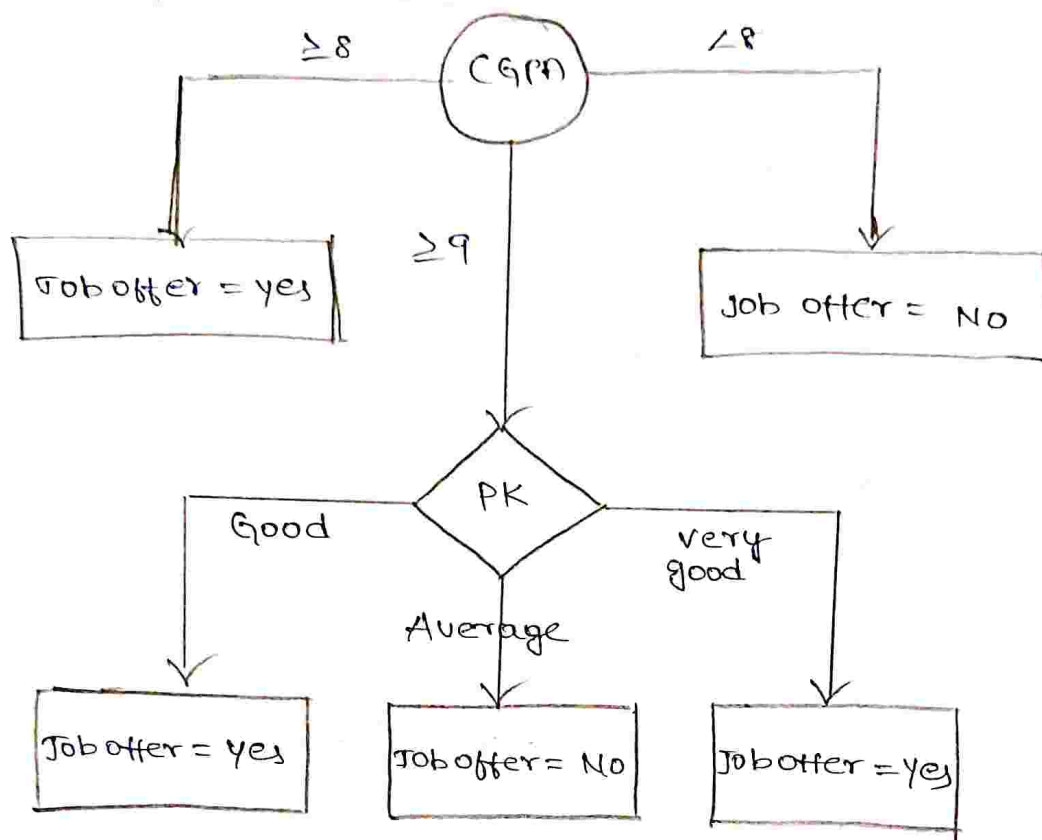
$$\text{split-info (T, inter)} = -2/4 \log_2(2/4) - 2/4 \log_2(2/4) = (0.5 + 0.5) = 1$$

$$\text{Gain ratio} = \frac{0.3112}{1} = 0.3112$$

Gain-ratio

<u>Attributes</u>	<u>Gain-ratio</u>
interactiveness	0.3112
practical Knowledge	0.5408
Communication skills	0.5408

Decision tree :





4)

Implement KNN classifier (only for classification)

Example:

Givendata query  $\Rightarrow x = (\text{Maths} = 6, \text{CS} = 8)$ and  $K = 3$  - nearest neighbour

classification - pass/fail

S.No	Maths	CS	Result
1)	4	3	F
2)	6	7	P
3)	7	8	P
4)	5	5	F
5)	8	8	P

Step-1:

Calculate Euclidean distance (d)

$$d = \sqrt{(x_{O1} - x_{A1})^2 + (x_{O2} - x_{A2})^2}$$

O  $\rightarrow$  observed valueA  $\rightarrow$  Actual Value.

$$d_1 = \sqrt{(6-4)^2 + (8-3)^2} = \sqrt{4+25} = \sqrt{29} = 5.38$$

$$d_2 = \sqrt{(6-6)^2 + (8-7)^2} = \sqrt{0+1} = 1$$

$$d_3 = \sqrt{(6-7)^2 + (8-8)^2} = \sqrt{1+0} = 1$$

$$d_4 = \sqrt{(6-5)^2 + (8-5)^2} = \sqrt{1+9} = \sqrt{10} = 3.16$$

$$d_5 = \sqrt{(6-8)^2 + (8-8)^2} = \sqrt{4+0} = \sqrt{4} = 2$$

Step-2: The distances that are closer and less than K value are 1, 1, 2 i.e.,  $d_2$ ,  $d_3$  and  $d_5$  their, result values are Pass

$\therefore$  classification - pass

5)

Implement Linear discriminant analysis (LDA) for suitable example.

$$x_1 = \{ (1,1), (2,4), (2,3), (3,6), (4,4) \}$$

$$x_2 = \{ (9,10), (6,8), (9,5), (8,7), (10,8) \}$$

Step - 1:

Calculate means

$$\mu_1 = \begin{bmatrix} \frac{1+2+2+3+4}{5} \\ \frac{1+4+3+6+4}{5} \end{bmatrix} = \begin{bmatrix} 3 \\ 3.6 \end{bmatrix}$$

$$\mu_2 = \begin{bmatrix} \frac{9+6+9+8+10}{5} \\ \frac{10+8+5+7+8}{5} \end{bmatrix} = \begin{bmatrix} 8.4 \\ 7.6 \end{bmatrix}$$

Step - 2:

Scatter Matrix

$$S_i = \sum_{j=1}^n (x_{ij} - \text{mean } x_i) (x_{ij} - \text{mean } x_i)^T$$

$$S_1 = \left[ (1-3, 1-3.6) (1-3, 1-3.6)^T + (2-3, 4-3.6) (2-3, 4-3.6)^T + \right. \\ \left. (-1) (-0.6) (2-3) (3-3.6)^T + (3-3) (6-3.6) (6-3.6)^T + \right. \\ \left. (4-3) (4-3.6)^T + (4-3) (4-3.6) \right]$$

$$S_1 = \begin{bmatrix} 0.8 & -0.4 \\ -0.4 & 2.6 \end{bmatrix}$$

$$S_2 = \left[ (9-8.4) (10-7.6) (9-8.4) (10-7.6)^T + (6-8.4) (8-7.6) (6-8.4) \right. \\ \left. (8-9.6)^T + (9-8.4) (5-7.6) (9-8.4) (5-7.6)^T + (8-8.4) \right. \\ \left. (7-7.6) (8-8.4) (7-7.6)^T + (10-8.4) (8-7.6) (8-7.6)^T \right]$$

Q-

$$S_2 = \begin{bmatrix} 1.84 & -0.04 \\ -0.04 & 2.64 \end{bmatrix}$$

$$S_W = (S_1 + S_2)$$

$$S_W = \begin{bmatrix} 0.8 & -0.4 \\ -0.4 & 2.6 \end{bmatrix} + \begin{bmatrix} 1.84 & -0.04 \\ -0.04 & 2.64 \end{bmatrix}$$

$$S_W = \begin{bmatrix} 2.64 & -0.44 \\ -0.44 & 5.28 \end{bmatrix}$$

Step - 3:

Scatter matrix  $S_B$

$$\begin{aligned} S_B &= (\mu_1 - \mu_2) (\mu_1 - \mu_2)^T \\ &= \begin{bmatrix} 3-8.4 \\ 3.6-7.6 \end{bmatrix} \begin{bmatrix} 3-8.4 \\ 3.6-7.6 \end{bmatrix}^T \\ &= \begin{bmatrix} -5.4 \\ -4 \end{bmatrix} \begin{bmatrix} -5.4 & -4 \end{bmatrix} \\ &= \begin{bmatrix} 29.16 & 21.6 \\ 21.6 & 16 \end{bmatrix} \end{aligned}$$

$$S_W^{-1} = \frac{\text{Adj}(S_W)}{|S_W|}$$

$$|S_W| = (2.64)(5.28) - (0.44)^2$$

$$|S_W| = 13.93 - 0.19$$

$$|S_W| = 13.73$$

$$\text{Adj}(S_W) = \begin{bmatrix} 16 & -21.6 \\ -21.6 & 29.16 \end{bmatrix}$$



$$S_{\omega}^{-1} S_B = \frac{1}{13.73} \begin{bmatrix} 16 & -21.6 \\ -21.6 & 29.16 \end{bmatrix} \begin{bmatrix} 29.16 & 21.6 \\ 21.6 & 16 \end{bmatrix}$$

$$S_{\omega}^{-1} S_B = \begin{bmatrix} 11.89 & 8.81 \\ 5.08 & 3.76 \end{bmatrix}$$

Step - 4:

The LDI projection is then obtained as the solution of the generalized eigen value problem.

$$S_{\omega}^{-1} S_B v = \lambda v$$

$$| S_{\omega}^{-1} S_B - \lambda I | = 0$$

$$\begin{vmatrix} (11.89 - \lambda) & 8.81 \\ 5.08 & (3.76 - \lambda) \end{vmatrix} = 0$$

$$(11.89 - \lambda)(3.76 - \lambda) - (8.81)(5.08) = 0$$

$$\lambda = 15.85, 0$$

Step - 5:

$$\omega^* = \begin{bmatrix} \omega_1 \\ \omega_2 \end{bmatrix}$$

$$\begin{bmatrix} \omega_1 \\ \omega_2 \end{bmatrix} = S_{\omega}^{-1} [\mu_1, \mu_2]$$

$$= \begin{bmatrix} 16 & -21.6 \\ -21.6 & 29.16 \end{bmatrix} \begin{bmatrix} -5.4 \\ -4 \end{bmatrix} = \begin{bmatrix} -0.91 \\ -0.39 \end{bmatrix}$$

$$\left\{ \begin{array}{l} \therefore \omega_1 = -0.91 \\ \omega_2 = -0.39 \end{array} \right\}$$