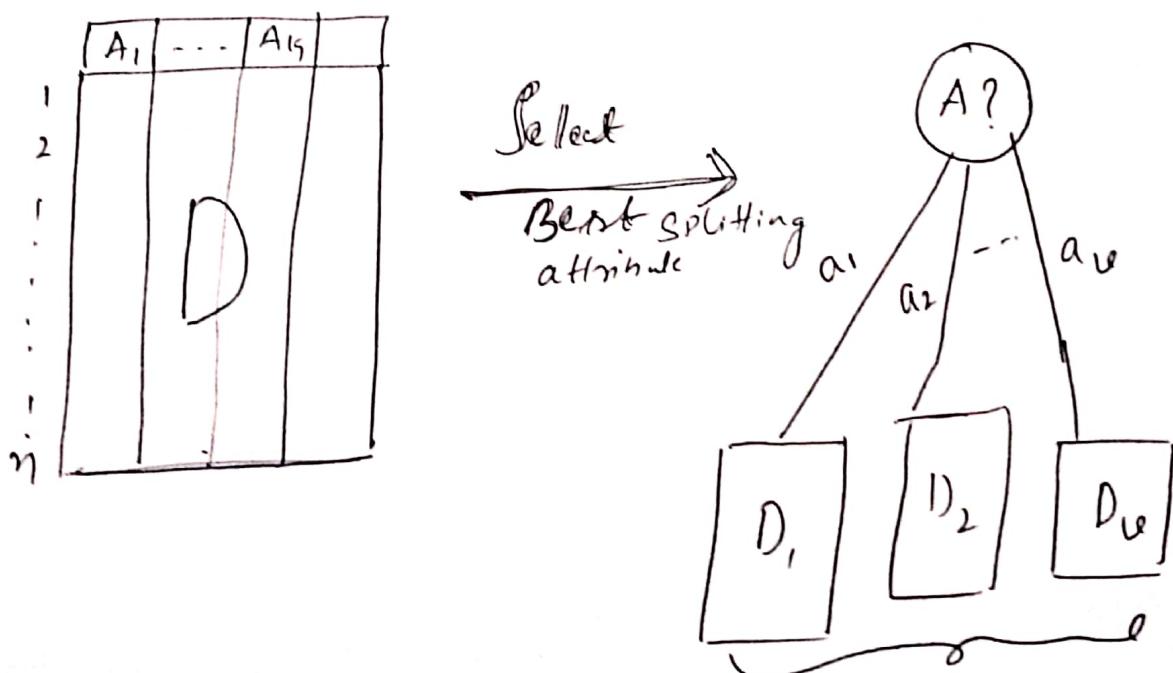


Attribute Selection Measures

An attribute Selection Measure is
a heuristic for selecting criterion
that best separates a given Data
partition, D of class-labeled training
tuples in individual classes.

Selecting Criterion \rightarrow Best Splitting attribute or Best Splitting point.



Ideally, we expect resultant partition to be pure i.e., after splitting.

→ The attribute Selection Measure provides a ranking for each attribute describing given training tuples.

~~The attribute Selection measure provides a ranking for each attribute~~

- The attribute having the best score is chosen as the splitting attribute for the given tuples.
- The tree node created for partition D is labelled with Splitting Criterion & branches grown for each outcome of criterion & tuples are partitioned accordingly.

Attribute Selection Measures

- (1) Information Gain
- (2) Gain ratio
- (3) Gini index.

Notation :

Let D , Data partition be a training set of class-labeled tuples

Let m be the no. of distinct classes

$$C_i \quad (i = 1 \dots m)$$

Let C_{iD} be the set of tuples of class C_i in D .

Let $|D|$ be the no. of tuples in D

Let $|C_{iD}|$ be no. of tuples in C_{iD}

Information Gain

→ ID3 Algorithm uses Information Gain as its Attribute Selection Measure.

→ This Measure is based on Shannon's Information Theory.

→ The attribute with the highest information gain is chosen as the Splitting attribute for Node N.

This attribute minimizes the information needed to classify tuples in the resulting partitions. & reflects least impurity in ~~these~~ partitions.

The expected information needed to classify a tuple in D is given by

$$\boxed{\text{Info}(D) = - \sum_{i=1}^m p_i \log_2(p_i)}$$

where

P_i = probability of a tuple in D
belongs to class C_i

i.e.,

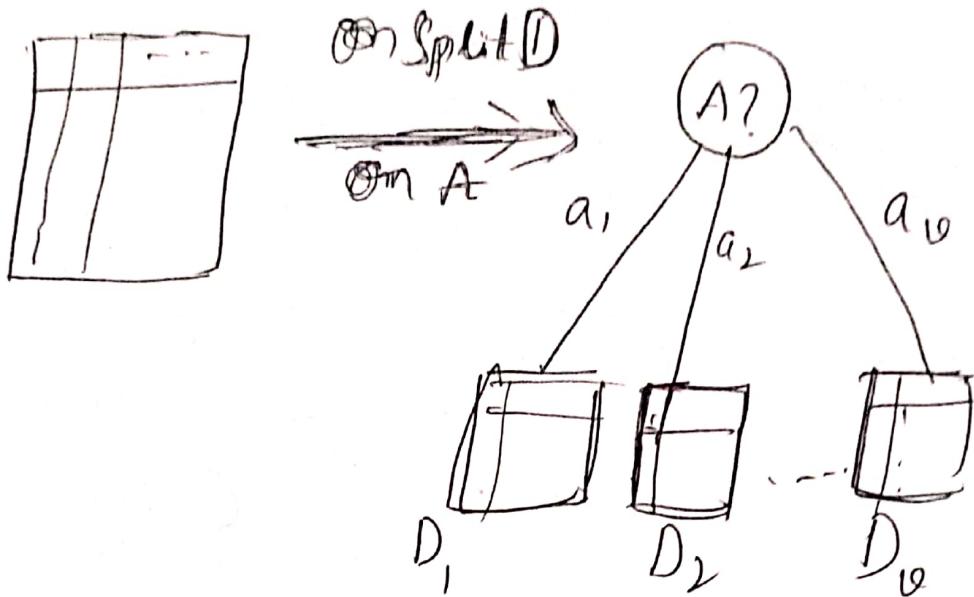
$$P_i = \frac{|C_i D|}{|D|} \quad \rightarrow ②$$

Then, partition the tuples in D on
some attribute A having v distinct
values, $\{a_1, a_2, \dots, a_v\}$, as observed
from the training Data.

Attribute A can be used to split
 D into v partitions $\{D_1, D_2, \dots, D_v\}$

where D_j = contains
tuples in D
that have
outcome a_j of A .

i.e.



Then Information required after.

partitioning is given by
on A

$$\text{Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

where D_j = no. of ϕ -tuples
in D_j

$$\text{Info}(D_j) = - \sum_{i=1}^m p_i \log_2(p_i)$$

$$p_i = \frac{|C_i D_j|}{|D_j|}$$

Then

$$\text{Gain}(A) = \underline{\text{Info}(D)} - \underline{\text{Info}_A(D)}$$

Where $\text{Gain}(A)$ is Expected
reduction in the information
requirement on partition D
by attribute (A)

Note: The attribute A with
highest information gain
is chosen as the splitting
attribute at node N.

Ex:-

TID	age	income	student	credit rating	Class buys-computer
1	youth	high	no	fair	no
2	youth	high	no	Excellent	no
3.	middle-aged	high	no	fair	yes
4.	Senior	medium	no	fair	yes
5	Senior	low	yes	fair	yes
6	Senior	low	yes	Excellent	no
7	middle-aged	low	yes	Excellent	yes
8	youth	medium	no	fair	no
9	youth	low	yes	fair	yes
10	Senior	medium	yes	fair	yes
11	youth	medium	yes	Excellent	yes
12	middle-aged	medium	no	Excellent	yes
13	Middle-aged	high	yes	fair	yes
14	Senior	medium	no	Excellent	no

$$\text{Info}(D) = - \sum_{i=1}^2 p_i \log_2 p_i$$

$$= - [p_1 \log_2 p_1 + p_2 \log_2 p_2]$$

where $p_1 = \frac{|C_{1D}|}{|D|}$ $p_2 = \frac{|C_{2D}|}{|D|}$

$$= \frac{9}{14} \quad p_2 = \frac{5}{14}$$

$$= - \frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

$$\boxed{\text{Info}(D) = 0.940 \text{ bits}}$$

④ Let partition D on attribute age

$$\begin{aligned}
 \text{Info}(D) &= \sum_{\substack{\text{age} \\ j=1}}^v \frac{|D_j|}{|D|} \times \text{Info}(D_j) \\
 &= \frac{|D_{\text{age}=\text{young}}|}{|D|} \text{Info}\left(D_{\text{age}=\text{young}}\right) \\
 &\quad + \frac{|D_{\text{age}=\text{middle-age}}|}{|D|} \text{Info}\left(D_{\text{age}=\text{middle-age}}\right) \\
 &\quad + \frac{|D_{\text{age}=\text{senior}}|}{|D|} \text{Info}\left(D_{\text{age}=\text{senior}}\right) \\
 &= \frac{5}{14} \times \left[-\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \right] \\
 &\quad + \frac{4}{14} \times \left[-\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \right] \\
 &\quad + \frac{5}{14} \times \left[-\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \right] \\
 &= 0.694 \text{ bits}
 \end{aligned}$$

$$Gain(\text{age}) = \text{Info}(D) - \frac{\text{Info}(D)}{\text{age}}$$

$$= 0.940 - 0.694$$

$$= 0.246 \text{ bits}$$

Similarly

$$Gain(\text{income}) = 0.029 \text{ bits}$$

$$Gain(\text{student}) = -0.151 \text{ bits}$$

$$Gain(\text{credit-rating}) = 0.048 \text{ bits.}$$

∴ The attribute age is
Selected as splitting attribute

age ?

income	student	credit-rating	class-label
high	no	fair	no
high	no	Excellent	no
medium	no	fair	yes
low	yes	fair	yes
medium	yes	Excellent	yes

income	student	credit-rating	class-label
medium	no	fair	yes
low	yes	fair	yes
low	yes	Excellent	no
medium	yes	fair	yes
medium	no	Excellent	no

income	student	credit-rating	class-label
high	no	fair	yes
low	yes	Excellent	yes
medium	no	Excellent	yes
high	yes	fair	yes

resultant first
 Since The ↑ partition is impure,
 Split partition by Best Attribute
 from remaining attribute list
 $= \{ \text{income, student, credit-rating} \}$

Then \oplus find Information gain \Rightarrow
 for each remain attribute.
 & Select \oplus attribute with highest
 Information gain.

$$Info(D) = - \sum_{i=1}^2 p_i \log_2 p_i = -[P_1 \log_2 P_1 + P_2 \log_2 P_2]$$

$$P_1 = \frac{|C_1 D|}{|D|} = \frac{2}{5} = 0.4$$

$$P_2 = \frac{|C_2 D|}{|D|} = \frac{3}{5} = 0.6$$

$$\begin{aligned} &= - \left[\frac{2}{5} \log_2 \frac{2}{5} + \frac{3}{5} \log_2 \frac{3}{5} \right] \\ &= - [0.4(-1.321) + 0.6(-0.7369)] \\ &= -[-0.5284 - 0.4421] \\ &= \underline{0.9705 \text{ bits}} \end{aligned}$$

$$Info(D)_{income} = \sum_{j=1}^3 \frac{|D_j|}{|D|} Info(D_j)$$

$$= \frac{|D_{high}|}{|D|} \times Info(D_{high})$$

$$+ \frac{|D_{medium}|}{|D|} \times Info(D_{medium})$$

$$+ \frac{|D_{low}|}{|D|} \times Info(D_{low})$$

$$\begin{aligned}
 &= \frac{2}{5} \times \left[-\frac{0}{2} \log_2 \cancel{\bar{1}}^0 (\cancel{\bar{1}}^0) - \frac{2}{2} \log_2 \cancel{\bar{1}}^0 (\cancel{\bar{1}}^0) \right] \\
 &\quad + \frac{2}{5} \times \left[-\frac{1}{2} \log_2 (\cancel{\bar{1}}^0) - \frac{1}{2} \log_2 (\cancel{\bar{1}}^0) \right] \\
 &\quad + \frac{1}{5} \times \left[-\frac{1}{1} \log_2 \cancel{\bar{1}}^0 (\cancel{\bar{1}}^0) - \frac{0}{1} \log_2 \cancel{\bar{1}}^0 (\cancel{\bar{1}}^0) \right] \\
 &= 0.4 \times \left[-Y_2(-1) - \frac{1}{2}(-1) \right] \\
 &= 0.4 \times \left[Y_2 + Y_2 \right] \\
 &= 0.4 [.] \\
 &= \underline{0.4 \text{ bits}}
 \end{aligned}$$

~~Info~~ Gain(income)

$$\begin{aligned}
 &= \text{Info}(D) - \text{Info}_{\text{income}}(D) \\
 &= 0.9705 - 0.4
 \end{aligned}$$

$$\boxed{\text{Info} \underline{\text{Gain}}(\text{income}) = \underline{\underline{0.5705 \text{ bits}}}}$$

$$\text{Info}(D)_{\text{student}} = \sum_{j=1}^{|D|} \frac{|D_j|}{|D|} \text{Info}(D_j)$$

$$= \frac{|D_{yes}|}{|D|} \text{Info}(D_{yes})$$

$$+ \frac{|D_{no}|}{|D|} \text{Info}(D_{no})$$

$$= \frac{2}{5} \text{Info}(D_{yes})$$

$$+ \frac{3}{5} \text{Info}(D_{no})$$

$$= \frac{2}{5} \left[-\frac{2}{2} \cancel{\log_2(\frac{2}{2})} - \frac{0}{2} \cancel{\log_2(\frac{0}{2})} \right]$$

$$+ \frac{3}{5} \left[-\frac{0}{3} \cancel{\log_2(\frac{0}{3})} - \frac{3}{3} \cancel{\log_2(\frac{3}{3})} \right]$$

$$= 0 \text{ bits}$$

~~Info~~ Gain (student)

$$= \text{Info}(D) - \text{Info}_{\text{student}}(D)$$

$$= 0.9705 - 0$$

$$\boxed{\text{InfoGain}(\text{student}) = \underline{0.9705 \text{ bits}}}$$

$$\begin{aligned}
 \text{Info}(D) &= \sum_{j=1}^v \frac{|D_j|}{|D|} \text{Info}(D_j) \\
 \text{Credit-rating} &= \frac{|D_{\text{fair}}|}{|D|} \text{Info}(D_{\text{fair}}) \\
 &\quad + \frac{|D_{\text{uncollectible}}|}{|D|} \text{Info}(D_{\text{uncollectible}}) \\
 &= \frac{3}{5} \text{Info}(D_{\text{fair}}) \\
 &\quad + \frac{2}{5} \text{Info}(D_{\text{uncollectible}}) \\
 &= \frac{3}{5} \left[-\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) \right] \\
 &\quad + \frac{2}{5} \left[-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right] \\
 &= \frac{3}{5} \left[-(-0.5331) - (-0.399) \right] \\
 &= \frac{3}{5} (0.9321) + \frac{2}{5} (1) \\
 &= \underline{0.5592 \text{ bits}} + 0.4 = 0.9592
 \end{aligned}$$

$$\begin{aligned}
 \cancel{\text{Info}} \text{Gain}(\text{credit_ratio}) &= \text{Info}(D) - \underbrace{\text{Info}(D)}_{\text{credit_ratio}} \\
 &= 0.9305 - 0.9592 \\
 \cancel{\text{Info}} \text{Gain}(\text{credit_ratio}) &= 0.0113 \text{ bits}
 \end{aligned}$$

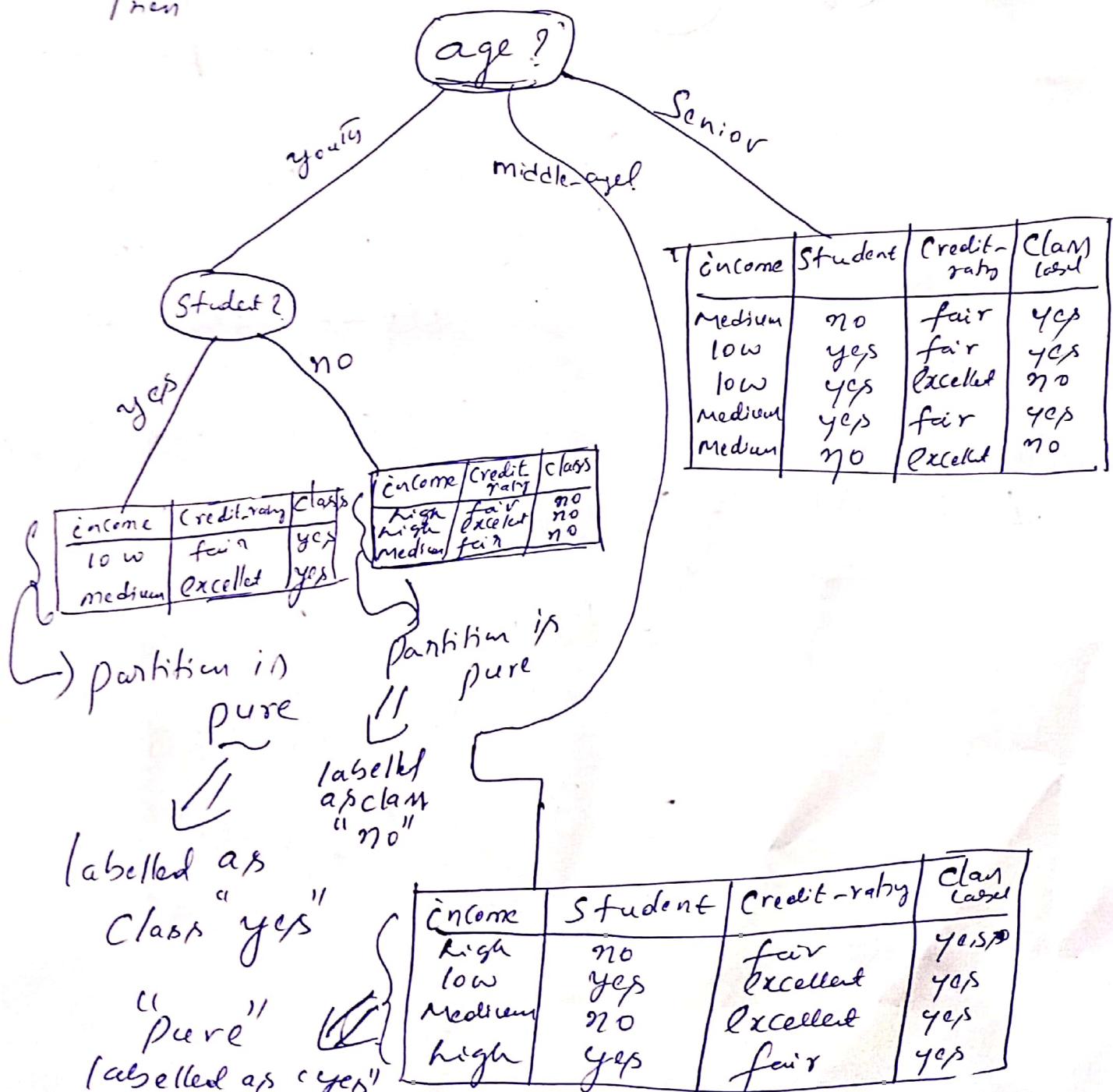
Since, ~~Gain~~ Gain (income) = 0.5705 bits

~~Gain~~ Gain (student) = 0.9705 bits

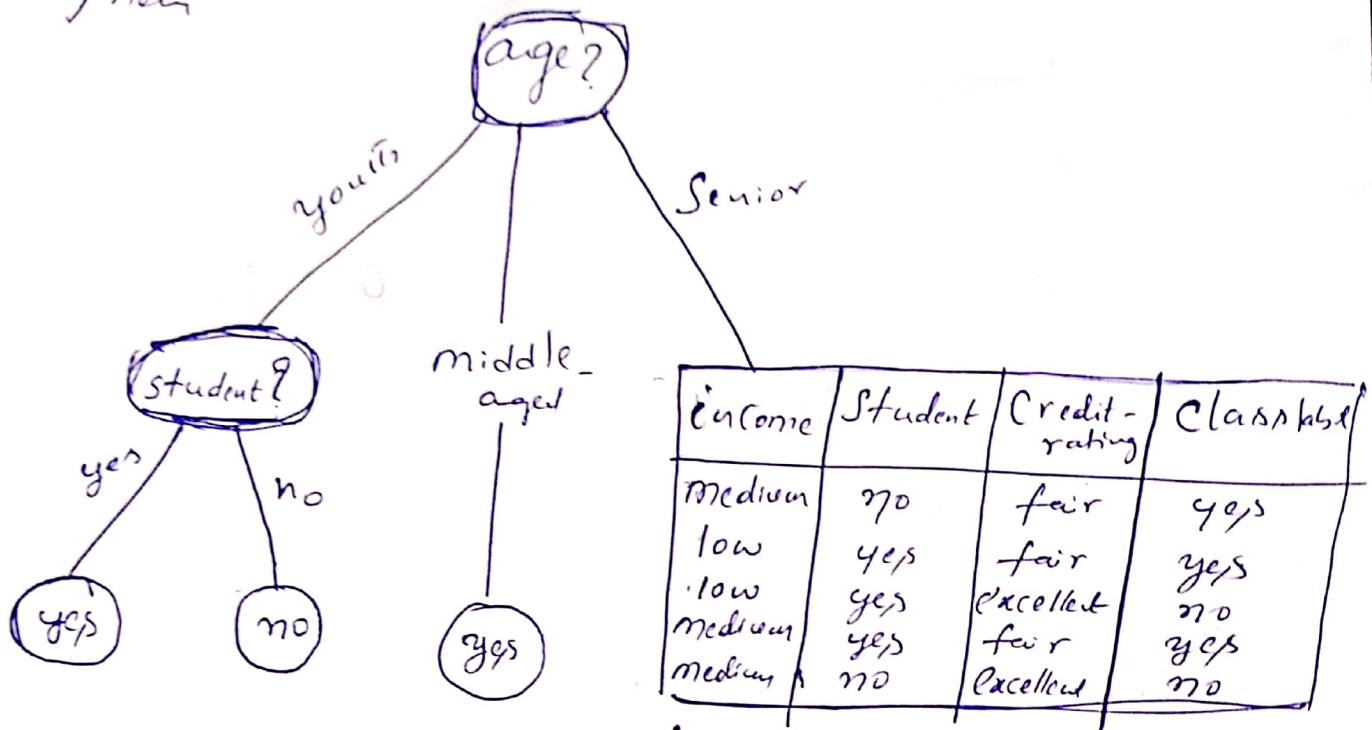
~~Gain~~ Gain (Credit-rating) = 0.0113 bits

Therefore student attribute will be selected for splitting partition.

Then



Then



Since impure parti Split
the partition using best attribute
from remaining attribute

list

$$= \{ \text{income, Credit-rating} \}$$

$$\begin{aligned}
 \text{Info}(D) &= - \sum_{i=1}^2 p_i \log_2 p_i \\
 &= - [P_1 \log_2 P_1 + P_2 \log_2 P_2] \\
 &= - [\frac{3}{5} \log_2 \frac{3}{5} + \frac{2}{5} \log_2 \frac{2}{5}] \\
 &= - [0.6 \log_2 0.6 + 0.4 \log_2 0.4] \\
 &= \underline{0.9705 \text{ bits}}
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}(D)_{\text{income}} &= \sum_{j=1}^{|D|} \frac{|D_j|}{|D|} \text{Info}(D_j) \\
 &= \frac{|D_{low}|}{|D|} \text{Info}(D_{low}) \\
 &\quad + \frac{|D_{median}|}{|D|} \text{Info}(D_{med}) \\
 &= \frac{2}{5} \left[-\frac{1}{2} \log_2(\frac{1}{2}) - \frac{1}{2} \log_2(\frac{1}{2}) \right] \\
 &\quad + \frac{3}{5} \left[-\frac{2}{3} \log_2(\frac{2}{3}) - \frac{1}{3} \log_2(\frac{1}{3}) \right] \\
 &= \frac{2}{5} [1] \\
 &\quad + \frac{3}{5} [0.9321] \\
 &= 0.9592
 \end{aligned}$$

$$\begin{aligned}
 \text{Gain(income)} &= \text{Info}(D) - \text{Info}(D)_{\text{income}} \\
 &= 0.9705 - 0.9592 \\
 \boxed{\text{Gain(income)} = 0.0113 \text{ bits}}
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}(D)_{\text{Credit-rat}} &= \sum_{j=1}^9 \frac{|D_j|}{|D|} \text{Info}(D_j) \\
 &\equiv \frac{|D_{\text{fair}}|}{|D|} \text{Info}(D_{\text{fair}}) \\
 &\quad + \frac{|D_{\text{excellent}}|}{|D|} \text{Info}(D_{\text{excellent}}) \\
 &= \frac{3}{5} \left[-\frac{3}{3} \log_2 \left(\frac{3}{3} \right) - \frac{0}{3} \log_2 \left(\frac{0}{3} \right) \right] \\
 &\quad + \frac{2}{5} \left[-\frac{0}{2} \log_2 \left(\frac{0}{2} \right) - \frac{2}{2} \log_2 \left(\frac{2}{2} \right) \right] \\
 &= 0 \text{ bits.}
 \end{aligned}$$

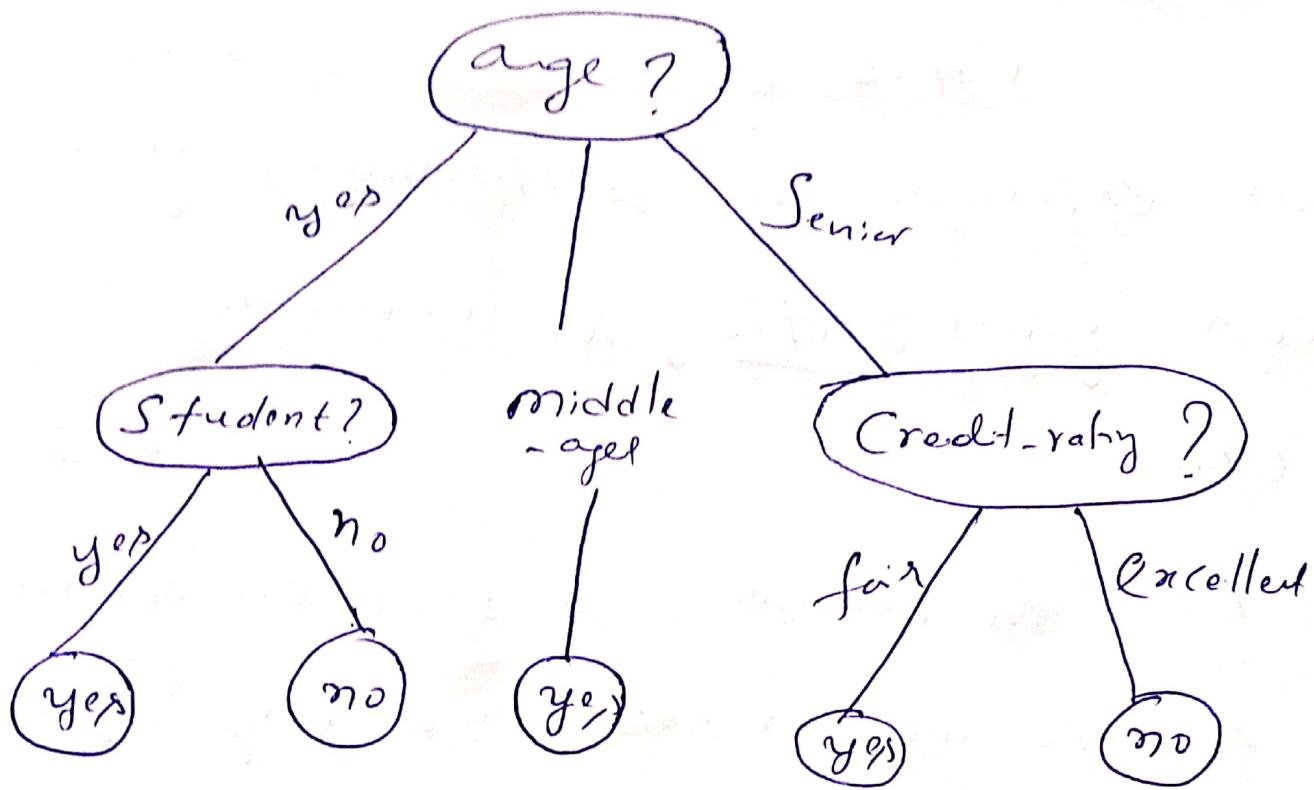
~~Info~~ Gain (Credit-rat)

$$= 0.9705 - 0$$

$$= \underline{\underline{0.9705}}$$

Therefore, Attribute Credit-rat
is selected for splitting
a partition.

Then



This value represents the potential information generated by splitting the training Data set D into v partitions, corresponding to the v outcomes of a test on attribute A .

Note:- The attribute with the maximum gain ratio is selected as the splitting attribute.

Ex:- Consider attribute income,

$$\text{GainRatio}(\text{income}) = \frac{\text{Gain}(\text{income})}{\text{Split}(D)}$$

$$\text{WIST, } \text{Gain}(\text{income}) = 0.029$$

$$\text{Split}(D) = -\frac{4}{14} \times \log_2\left(\frac{4}{14}\right) - \frac{6}{14} \times \log_2\left(\frac{6}{14}\right)$$

$$-\left(\frac{4}{14}\right) \times \log_2\left(\frac{4}{14}\right)$$

$$= 1.557$$

$$\text{Gain Ratio}(\text{income}) = \frac{0.029}{1.557} = \underline{\underline{0.019}}$$

Next consider attribute age

$$\text{Gainratio(Age)} = \frac{\text{Gain(Age)}}{\text{Splitinfo}(D)}$$

L. R. $\text{Gain(Age)} = 0.246$

$$\text{Splitinfo}(D)$$

$$= - \sum_{j=1}^{\text{age}} \frac{|D_j|}{|D|} \times \log_2 \left(\frac{|D_j|}{|D|} \right)$$

∴

$$= - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) - \frac{4}{14} \log_2 \left(\frac{4}{14} \right) \\ - \frac{5}{14} \log_2 \left(\frac{5}{14} \right)$$

$$= -0.35 \log_2(0.35) - 0.28 \log_2(0.28) \\ - 0.35 \log_2(0.35)$$

$$= -[-0.5308 - 0.5082 - 0.5308]$$

$$= -[-1.5748]$$

$$= 1.5748$$

$$\text{Gainratio(Age)} = \frac{0.246}{1.5748} = 0.1562$$

Next Consider an attribute Gender of Student

$$\text{Gainratio}(\text{student}) = \frac{\text{Gain}(\text{student})}{\text{Splitinfo}(D)}$$

LHS

$$\text{Gain}(\text{student}) = \underline{\underline{0.151}}$$

$$\begin{aligned}\text{Splitinfo}(D)_{\text{student}} &= -\frac{7}{14} \log_2\left(\frac{7}{14}\right) - \frac{7}{14} \log_2\left(\frac{7}{14}\right) \\ &= -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\end{aligned}$$

~~Explain~~

$$= -Y_2(-1) - Y_2(-1)$$

$$= Y_2 + Y_2$$

$$= 1$$

$$\text{Gainratio}(\text{Student}) = \frac{0.151}{1}$$

$$= \underline{\underline{0.151}}$$

Next Consider attribute Credit-rating

$$\text{Gain ratio}(\text{Credit-rating}) = \frac{\text{Gain}(\text{Credit-rating})}{\text{SplitInfo}(D)}$$

SplitInfo(D)

Credit-rating

$$\text{LKR } \text{Gain}(\text{Credit-rating}) = \underline{\underline{0.048}}$$

SplitInfo(D)

$$\text{Credit-rating} = -\frac{8}{14} \log_2 \frac{8}{14} - \frac{6}{14} \log_2$$

$$= -0.57 \log_2(0.57) - 0.42 \log_2(0.42)$$

$$= -[-0.4622 - 0.5256]$$

$$= 0.9878$$

Gain ratio(Credit-rating)

$$= \frac{0.048}{0.9878}$$

$$\approx \underline{\underline{0.04859}}$$

$$\begin{aligned} \therefore \left\{ \begin{array}{l} \text{Gain ratio}(\text{age}) = \text{Gain ratio}(\text{age}) = 0.1562 \\ \text{Gain ratio}(\text{Student}) = 0.151 \\ \text{Gain ratio}(\text{Income}) = 0.019 \\ \text{Gain ratio}(\text{Credit-rating}) = 0.0485 \end{array} \right. \end{aligned}$$