

Assignment - 1 (ML)

Name - Devarsh (roll

Batch - 191333

CS - 65

Decision tree using entropy and information gain→ Entropy (S) of entire dataset

$$S = [9^+, 5^-] = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14}$$

$$= 0.94$$

→ Find entropy of individual attribute (Age, income, student, credit_rating)

i) Age → youth, middle-aged, senior

$$\text{youth} = [2^+, 3^-] \quad \text{entropy} = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} \\ = 0.971$$

$$\text{middle-aged} = [4^+, 0^-] \quad \text{entropy} = -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} \\ = 0$$

$$\text{senior} = [3^+, 2^-] \quad \text{entropy} = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \\ = 0.971$$

$$\text{Gain}(S, \text{Age}) = \text{entropy}(S) - \sum_{i=1}^{n_s} \frac{|S_i|}{|S|} \text{entropy}(S_i)$$

$$= 0.94 - \frac{5}{14} * 0.971 - \frac{4}{14} * 0 - \frac{5}{14} * 0.971 = [0.2464]$$

ii) Income \rightarrow low, medium, high

$$\text{low} = [3^+, 1^-]$$

$$e = -\frac{3}{4} \log_{2^4} 3 - \frac{1}{4} \log_{2^4} 1 = 0.8113$$

$$\text{medium} = [4^+, 2^-]$$

$$e = -\frac{4}{6} \log_{2^6} 4 - \frac{2}{6} \log_{2^6} 2 = 0.9183$$

$$\text{high} = [2^+, 2^-]$$

$$e = -\frac{2}{4} \log_{2^4} 2 - \frac{2}{4} \log_{2^4} 2 = 1$$

$$\text{Gain}(S, \text{income}) = 0.94 - \frac{4}{14} * 0.8113 - \frac{6}{14} * 0.9183 - \frac{4}{14} * 1$$

$$\boxed{\text{Gain(income)} = 0.0289}$$

iii) Student \rightarrow NO, yes

$$\text{no} = [3^+, 4^-]$$

$$e = -\frac{3}{7} \log_{2^7} 3 - \frac{4}{7} \log_{2^7} 4 = 0.9852$$

$$\text{yes} = [6^+, 1^-]$$

$$e = -\frac{6}{7} \log_{2^7} 6 - \frac{1}{7} \log_{2^7} 1 = 0.5916$$

$$\text{Gain}(S, \text{student}) = \text{entropy}(S) - \frac{7}{14} * 0.9852 - \frac{7}{14} * 0.5916$$

$$= 0.01916$$

iv) Credit-rating \rightarrow fair, excellent

$$\text{fair} = [6^+, 2^-]$$

$$e = -\frac{6}{8} \log_{2/8} \frac{6}{8} - \frac{2}{8} \log_{2/8} \frac{2}{8} = 0.8113$$

$$\text{excellent} = [3^+, 3^-]$$

$$e = -\frac{3}{6} \log_{2/6} \frac{3}{6} - \frac{3}{6} \log_{2/6} \frac{3}{6} = 1$$

$$\text{gain}(S, \text{credit-rating}) = 0.94 - \frac{6}{14} * 1 - \frac{8}{14} * 0.8113 = 0.0478$$

$$\text{gain}(S, \text{age}) = 0.2464$$

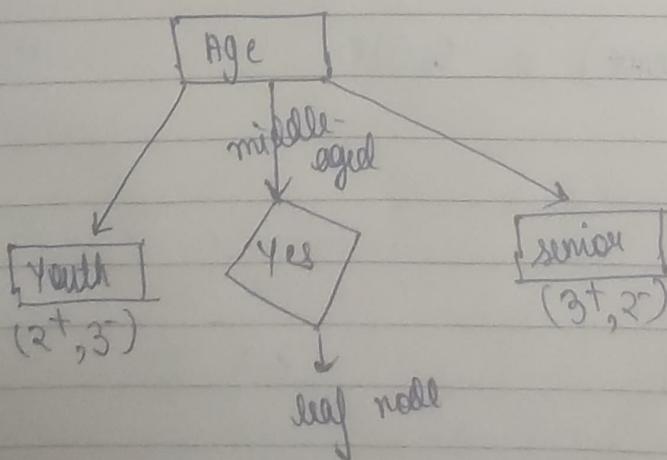
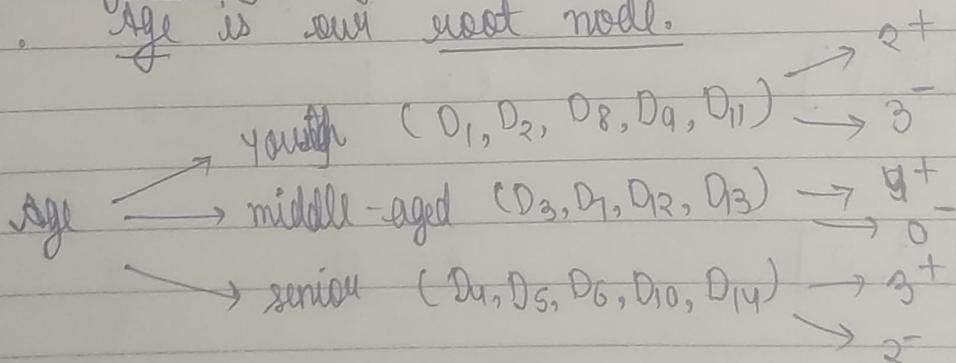
$$\text{gain}(S, \text{income}) = 0.0289$$

$$\text{gain}(S, \text{student}) = 0.01916$$

$$\text{gain}(S, \text{credit rating}) = 0.0478$$

Since, gain of age is maximum,

\therefore Age is our root node.



Now, consider only 'youth' and 'senior'.

For youth

Data point	Income	student	credit rating	buys
D ₁	high	no	fair	no
D ₂	high	no	excellent	no
D ₃	medium	no	fair	no
D ₄	low	yes	fair	yes
D ₅	medium	yes	excellent	yes

$$S_{youth} = [2^+, 3^-] \quad e = -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 0.97$$

i) income

$$\text{high} = [0^+, 2^-] \quad e = -\frac{2}{2} \log \frac{2}{2} = 0$$

$$\text{medium} = [1^+, 1^-] \quad e = 1$$

$$\text{low} = [1^+, 0^-] \quad e = 0$$

$$\text{Gain}(S_{youth}, \text{income}) = 0.510$$

ii) student

$$\text{no} = [0^+, 3^-] \quad e = 0$$

$$\text{yes} = [2^+, 0^-] \quad e = 0$$

$$\text{Gain}(S_{youth}, \text{student}) = 0.97$$

iii) credit rating

$$\text{fair} = [1^+, 2^-]$$

$$c = -\frac{1}{3} \log_{2/3} \frac{1}{2} - \frac{2}{3} \log_{2/3} \frac{2}{3} = 0.9183$$

$$\text{excellent} = [1^+, 1^-]$$

$$e = 1$$

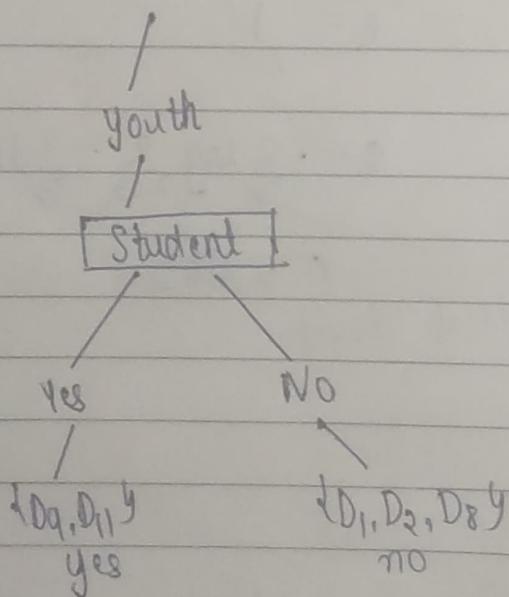
$$\text{Gain}(S_{\text{youth, credit}}) = 0.0192$$

$$\text{Gain}(S_{\text{youth, income}}) = 0.570$$

$$\text{Gain}(S_{\text{youth, student}}) = 0.97$$

\therefore Gain of student is maximum

\therefore student will be considered as node.



For senior

Data point	income	student	credit-rating	buys
D ₄	medium	no	fair	yes
D ₅	low	yes	fair	yes
D ₆	low	yes	excellent	no
D ₁₀	medium	yes	fair	yes
D ₁₄	medium	no	excellent	no

$$S_{\text{senior}} = [3^+, 2^-]$$

$$\begin{aligned} e &= -\frac{3}{5} \log \frac{3}{5} - \frac{2}{5} \log \frac{2}{3} \\ &= 0.97 \end{aligned}$$

ii) income

$$\text{low} = [1^+, 1^-]$$

$$e = 1$$

$$\text{medium} = [2^+, 1^-]$$

$$\begin{aligned} e &= -\frac{3}{5} \log \frac{2}{23} - \frac{1}{5} \log \frac{1}{23} \\ &= 0.9183 \end{aligned}$$

$$\text{Gain}(S_{\text{senior}}, \text{income}) = 0.0192$$

iii) student

$$\text{no} = [1^+, 1^-]$$

$$e = 1$$

$$\text{yes} = [2^+, 1^-]$$

$$e = 0.9183$$

$$\text{Gain}(S_{\text{senior}}, \text{student}) = 0.0192$$

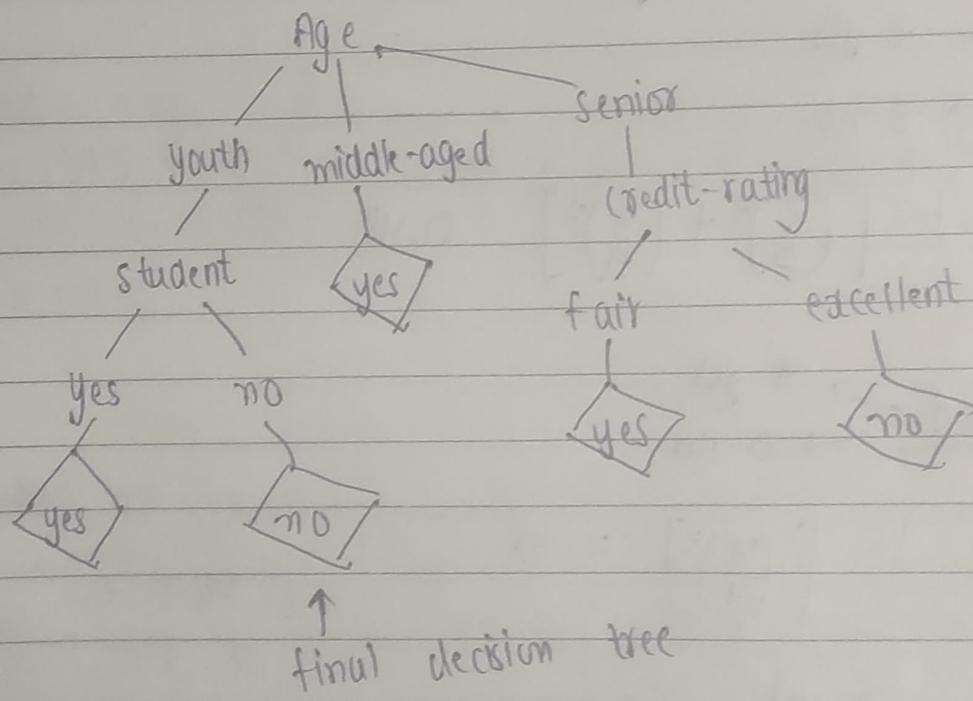
iii) credit-rating

$$\text{Jain} = [3+, 0^-] \quad e=0$$

$$\text{excellent} = [0^+, 2^-] \quad e=0$$

$$\text{Gini}(S_{\text{senior}}, \alpha) = 0.97$$

or, next node is credit-rating



Decision tree using 'gini index'

→ Gini index for overall dataset

Buy computer

$$\text{Gini}(S) = 1 - \left[\left(\frac{9}{14}\right)^2 + \left(\frac{5}{14}\right)^2 \right] = 0.459$$

→ Age → youth, middle-aged, senior

youth → 2 yes
youth → 3 no

$$\text{Gini} = 1 - \left[\left(\frac{2}{5}\right)^2 + \left(\frac{3}{5}\right)^2 \right] = 0.48$$

middle-aged → 4 yes
middle-aged → 0 no

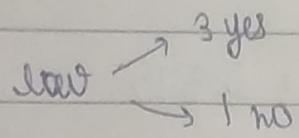
$$\text{Gini} = 1 - \left[\left(\frac{4}{4}\right)^2 \right] = 0$$

senior → 3 yes
senior → 2 no

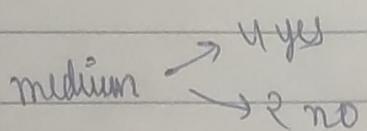
$$\text{Gini} = 1 - \left[\left(\frac{3}{5}\right)^2 + \left(\frac{2}{5}\right)^2 \right] = 0.48$$

$$\begin{aligned} \text{Weighted avg. g} &= 0.48 * \left(\frac{5}{14}\right) + 0 + 0.48 * \left(\frac{5}{14}\right) \\ &= \frac{4.8}{14} = 0.343 \end{aligned}$$

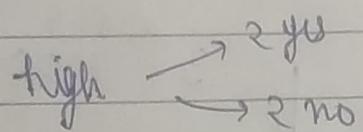
ii) Income \rightarrow low, medium, high



$$g_{(\text{minig})} = 1 - \left[\left(\frac{3}{4} \right)^2 + \left(\frac{1}{4} \right)^2 \right] = 0.375$$



$$g = 1 - \left[\left(\frac{4}{6} \right)^2 + \left(\frac{2}{6} \right)^2 \right] = 0.444$$



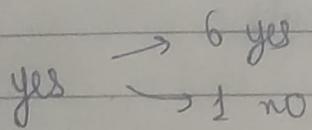
$$g = 1 - \left[2 * \left(\frac{2}{4} \right)^2 \right] = 0.5$$

Weighted avg. g

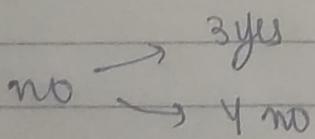
$$= 0.375 * \frac{4}{14} + 0.444 * \frac{6}{14} + \underline{0.5} * \frac{4}{14}$$

$$= 0.440$$

iii) Student \rightarrow yes, no



$$g = 1 - \left[\left(\frac{6}{7} \right)^2 + \left(\frac{1}{7} \right)^2 \right] = 0.245$$



$$g = 1 - \left[\left(\frac{3}{7} \right)^2 + \left(\frac{4}{7} \right)^2 \right] = 0.490$$

Weighted avg. g

$$= 0.245 * \frac{7}{14} + 0.49 * \frac{7}{14}$$

$$= 0.3675$$

0.735

iii) credit-rating \rightarrow fair, excellent

fair \rightarrow
2 no

$$g = 1 - \left[\left(\frac{6}{8}\right)^2 + \left(\frac{2}{8}\right)^2 \right] = 0.375$$

excellent \rightarrow
3 no

$$g = 0.5$$

Weighted avg. g

$$= 0.375 * \frac{8}{14} + 0.5 * \frac{6}{14}$$
$$= 0.429$$

$$\begin{array}{c} 6 \\ | \\ 3 \end{array} \quad \begin{array}{c} 4 \\ | \\ 3,000 \end{array}$$

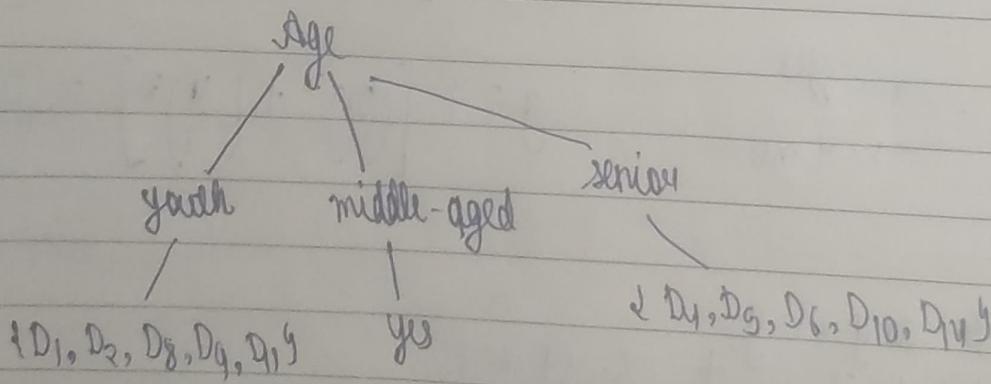
$$g(\text{age}) = 0.343$$

$$g(\text{income}) = 0.440$$

$$g(\text{student}) = 0.368$$

$$g(\text{credit-rating}) = 0.429$$

$\therefore g$ is smallest for 'age' attribute
 \therefore it will be root node



youth → 2 yes
→ 3 no

$$g = 1 - \left[\left(\frac{2}{5}\right)^2 + \left(\frac{3}{5}\right)^2 \right]$$

ii) income → low, medium, high

low → 1 yes
→ 0 no

$$g = 0$$

medium → 1 yes
→ 1 no

$$g = 0.5$$

high → 0 yes
→ 2 no

$$g = 0$$

weighted avg. $g = 0 + 0.5 * \frac{2}{5} + 0$
 $= 0.2$

iii) student → yes, no

yes → 2 yes
→ 0 no

$$g = 0$$

no → 0 yes
→ 3 no

$$g = 0$$

Weighted avg. $g = 0$

iii) credit-rating → fair, excellent

fair → 1 yes
→ 2 no

$$g = 1 - \left[\left(\frac{1}{3}\right)^2 + \left(\frac{2}{3}\right)^2 \right] = 0.444$$

excellent → 1 yes
→ 1 no

$$g = 0.5$$

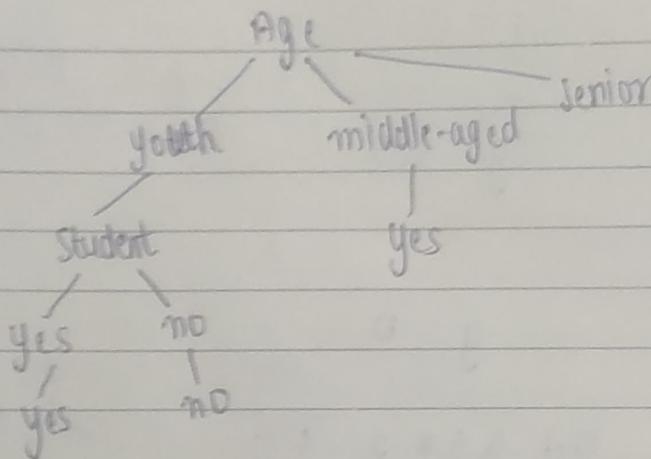
Weighted avg. $g = 0.444 * \frac{3}{5} + 0.5 * \frac{2}{5} = 0.466$

$$g(\text{income}) = 0.2$$

$$g(\text{student}) = 0$$

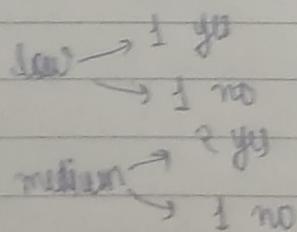
$$g(\text{u.}) = 0.466$$

$\therefore g$ of student is smallest
 \therefore next node is student node



For senior

i) income

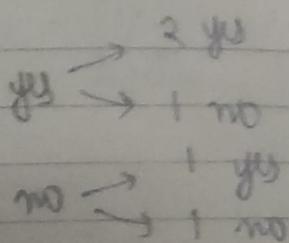


$$g = 0.5$$

$$g = \frac{1}{3} [0.5^2 + (1)^2] = 0.444$$

$$\text{Weighted avg. } g = 0.466$$

ii) student



$$g = 0.444$$

$$g = 0$$

$$\text{Weighted avg. } g = 0.3664$$

iii) credit-waiting

fair \rightarrow 3 yes
0 no

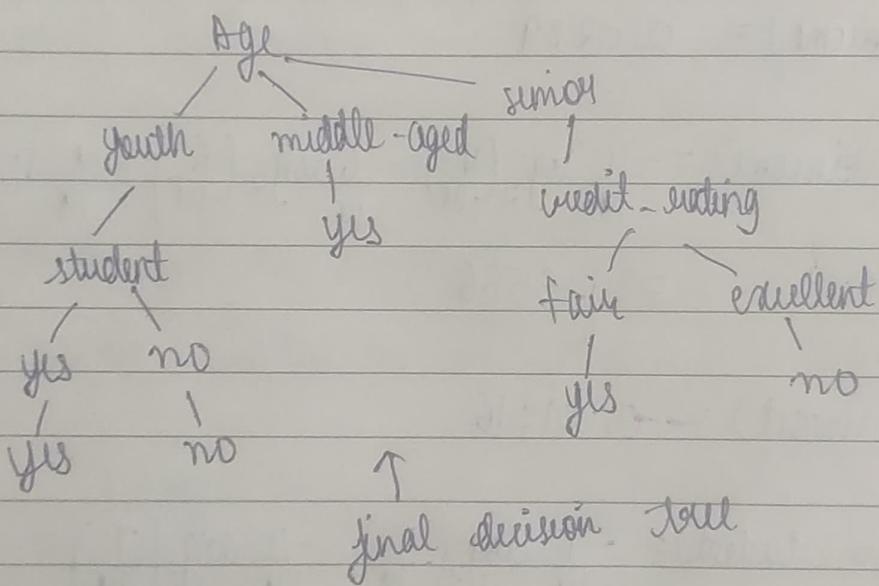
$$g = 0$$

excellent \rightarrow 0 yes
2 no

$$g = 0$$

Weighted avg. $g = 0$

or, next node is credit-waiting



Decision tree using gain ratio (information gain, split info)

$$\text{GR (gain ratio)} = \frac{\text{Gain (A)}}{\text{Split info (A)}}$$

$$\text{Split info (A)} = - \sum_{j=1}^y \frac{|D_j|}{|D|} * \log \left(\frac{|D_j|}{|D|} \right)$$

total no.
of data points

i) Age

$$\text{Gain}(s, \text{Age}) = 0.2464 \quad (\text{calculated before})$$

$$\begin{aligned} \text{split info (Age)} &= -\frac{5}{14} \log_2 \left(\frac{5}{14} \right) - \frac{4}{14} \log_2 \left(\frac{4}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) \\ &= 1.5774 \end{aligned}$$

ii) Income

$$\text{Gain}(s, \text{income}) = 0.0289$$

$$\begin{aligned} \text{split info (income)} &= -\frac{4}{14} \log_2 \left(\frac{4}{14} \right) - \frac{6}{14} \log_2 \left(\frac{6}{14} \right) - \frac{4}{14} \log_2 \left(\frac{4}{14} \right) \\ &= 1.5566 \end{aligned}$$

iii) student

$$\text{Gain}(s, \text{student}) = 0.61516$$

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

iv) credit-rating

$$\text{Gain}(s, \text{credit-rating}) = 0.0478$$

$$\begin{aligned} \text{split info (credit-rating)} &= -\frac{8}{14} \log_2 \left(\frac{8}{14} \right) - \frac{6}{14} \log_2 \left(\frac{6}{14} \right) \\ &= 0.9851 \end{aligned}$$

Gain ratio

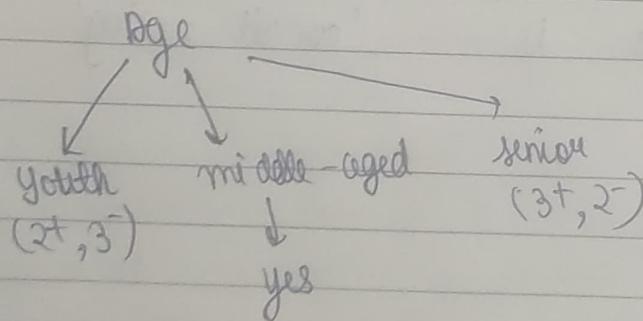
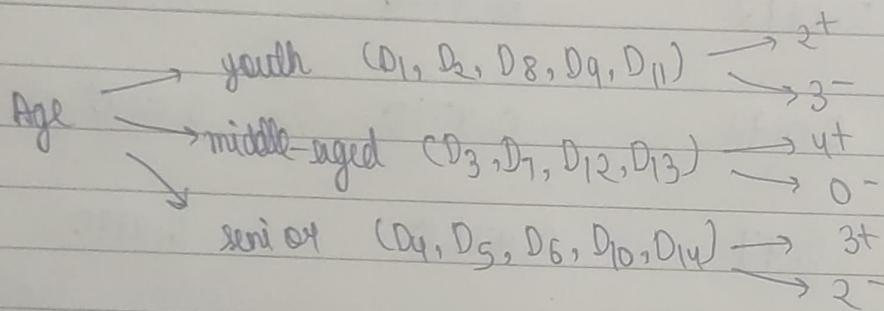
$$\text{i) age} = \frac{\text{Gain}}{\text{split info}} = \frac{0.2464}{1.5774} = 0.1562$$

ii) income = 0.0186

iii) student = 0.01816

iv) credit. = 0.0485

∴ gain ratio of age is maximum.
∴ root node is age



Now, consider only youth and senior.

For youth

i) income: Gain = 0.570
GIR = 0.3746

Split info = 1.5218

ii) student: Gain = 0.97
GIR = 0.9990

Split info = 0.9710

~~0.9994 + 3/5~~

iii) Credit-rating

$$\text{Gain} = 0.0192$$

$$\text{Split info} = 0.9710$$

$$\text{GR} = 0.0198$$

\therefore GR of student is maximum

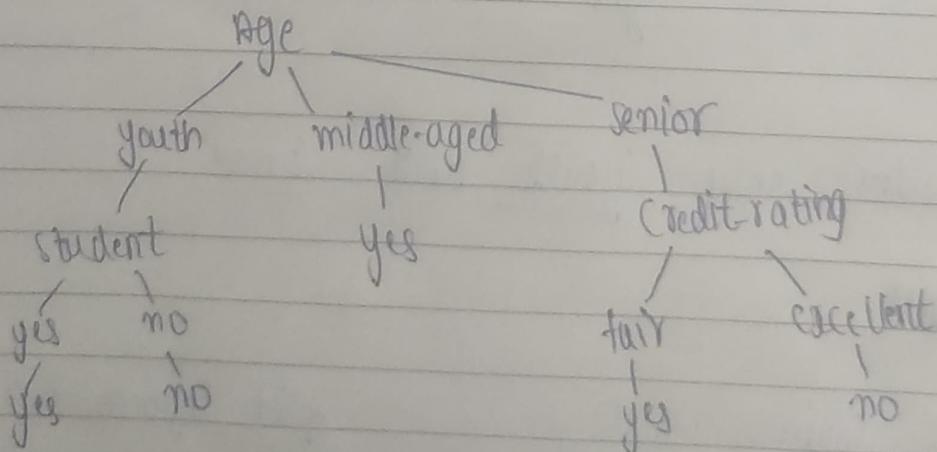
\therefore Student node will be considered.

For senior

	Gain	Split info	GR
income	0.0192	0.9710	0.0198
student	0.0192	0.9710	0.0198
credit-rating	0.97	0.9710	0.9990

\therefore next node is 'credit-rating'

Final decision tree



Shortage of ID3 (entropy + information gain)

- Attributes must be nominal values
- Dataset must not include missing data.
- Uses a greedy approach that why it does not guarantee an optimal soln; it can get stuck in local optimums.

C4.5 algo (gain, split info, gain ratio)

- C4.5 algorithm acts similar to ID3 but improves a few of ID3 behaviors:
 - A possibility to use continuous data
 - Using unknown (missing data) values

mini index

- Has difficulties when the no. of classes is large.
- Favors tests that result in equal-sized partitions with purity.