

# Attribute Selection with Information Gain.

age	income	student	credit_rating	buys_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
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
>40	low	yes	excellent	no

Feature

Class

(yes, no)

age	income	student	credit
$\begin{matrix} 5 \\ 2,3 \end{matrix}$ $\begin{matrix} 4 \\ 4,0 \end{matrix}$ $\begin{matrix} 3 \\ 3,2 \end{matrix}$	$\begin{matrix} 4 \\ 2,2 \end{matrix}$ $\begin{matrix} 6 \\ 4,2 \end{matrix}$ $\begin{matrix} 4 \\ 3,1 \end{matrix}$	$\begin{matrix} 7 \\ 6,1 \end{matrix}$ $\begin{matrix} 3 \\ 3,5 \end{matrix}$	$\begin{matrix} 6 \\ 3,3 \end{matrix}$ $\begin{matrix} 2 \\ 4,2 \end{matrix}$

• Overall Class on Info(D) =  $-\sum_{i=1}^m p_i \log_2(p_i)$

=  $I(9,5)$

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

= 0.940 - Expected info (entropy) within class

• Overall Feature on  $Info_A(D) = \sum_{j=1}^V \frac{|D_j|}{|D|} \times Info(D_j)$  -> minimization class

Expected info minimization Root node

$Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.6935$   
 $Info_{income}(D) = \frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1) = 0.9109$   
 $Info_{student}(D) = \frac{4}{14} I(6,1) + \frac{7}{14} I(3,4) = 0.7884$   
 $Info_{credit}(D) = \frac{2}{14} I(6,2) + \frac{6}{14} I(3,3) = 0.8929$

Gain from Root node

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

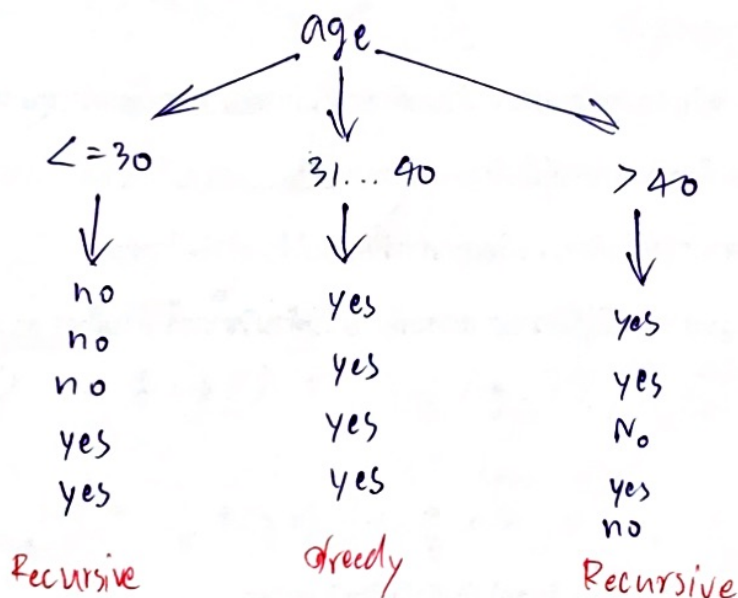
$$= 0.940 - \text{Info}_A(D)$$

$$\text{Gain}(\text{age}) = 0.940 - 0.6935 = 0.2465 \rightarrow \text{age is root node}$$

$$\text{Gain}(\text{income}) = 0.940 - 0.9109 = 0.0291$$

$$\text{Gain}(\text{student}) = 0.940 - 0.7884 = 0.1516$$

$$\text{Gain}(\text{credit\_rating}) = 0.940 - 0.8921 = 0.0479$$



Recursive age  $\leq 30$

$$\text{Info}(D) = I(2, 3) = -\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) = 0.9710$$

$$\text{Info}_{\text{income}}(D) = \frac{2}{5} I(0, 2) + \frac{2}{5} I(1, 1) + \frac{1}{5} I(1, 0) = 0.4$$

$$\text{Info}_{\text{student}}(D) = \frac{2}{5} I(1, 0) + \frac{3}{5} I(0, 3) = 0$$

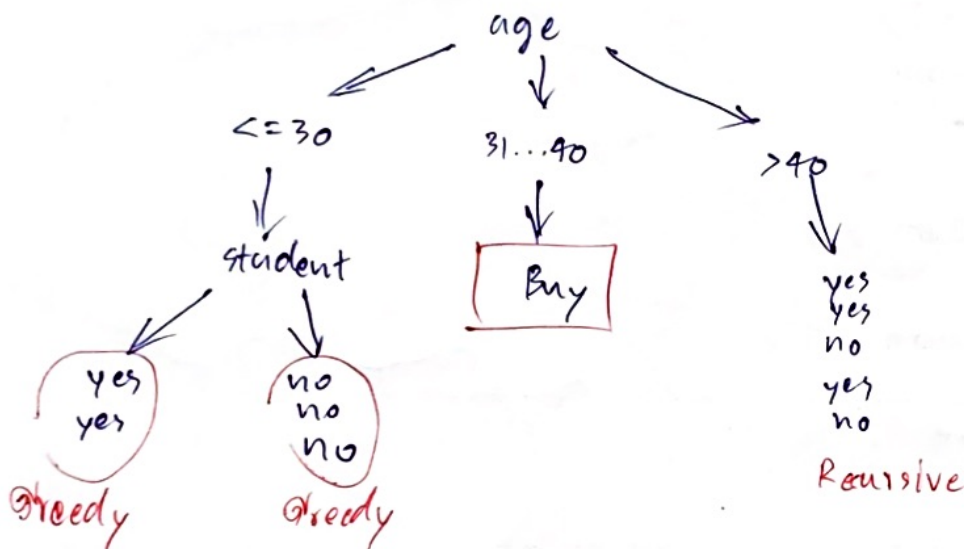
$$\text{Info}_{\text{Credit\_rating}}(D) = \frac{2}{5} I(1, 1) + \frac{3}{5} I(1, 2) = 0.9509$$

$$Gain(income) = 0.9710 - 0.4 = 0.5710$$

$$Gain(student) = 0.9710 - 0 = 0.9710 \rightarrow \text{1st on student}$$

$$Gain(credit\_rating) = 0.9710 - 0.9509 = 0.0201$$

164 root node



Recursive age > 40

$$Info(D) = I(3,2) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) = 0.9710$$

$$Info_{income}(D) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) = 0.4509$$

$$Info_{student}(D) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) = 0.4509$$

$$Info_{credit\_rating}(D) = \frac{2}{5} I(0,2) + \frac{3}{5} I(3,0) = 0$$

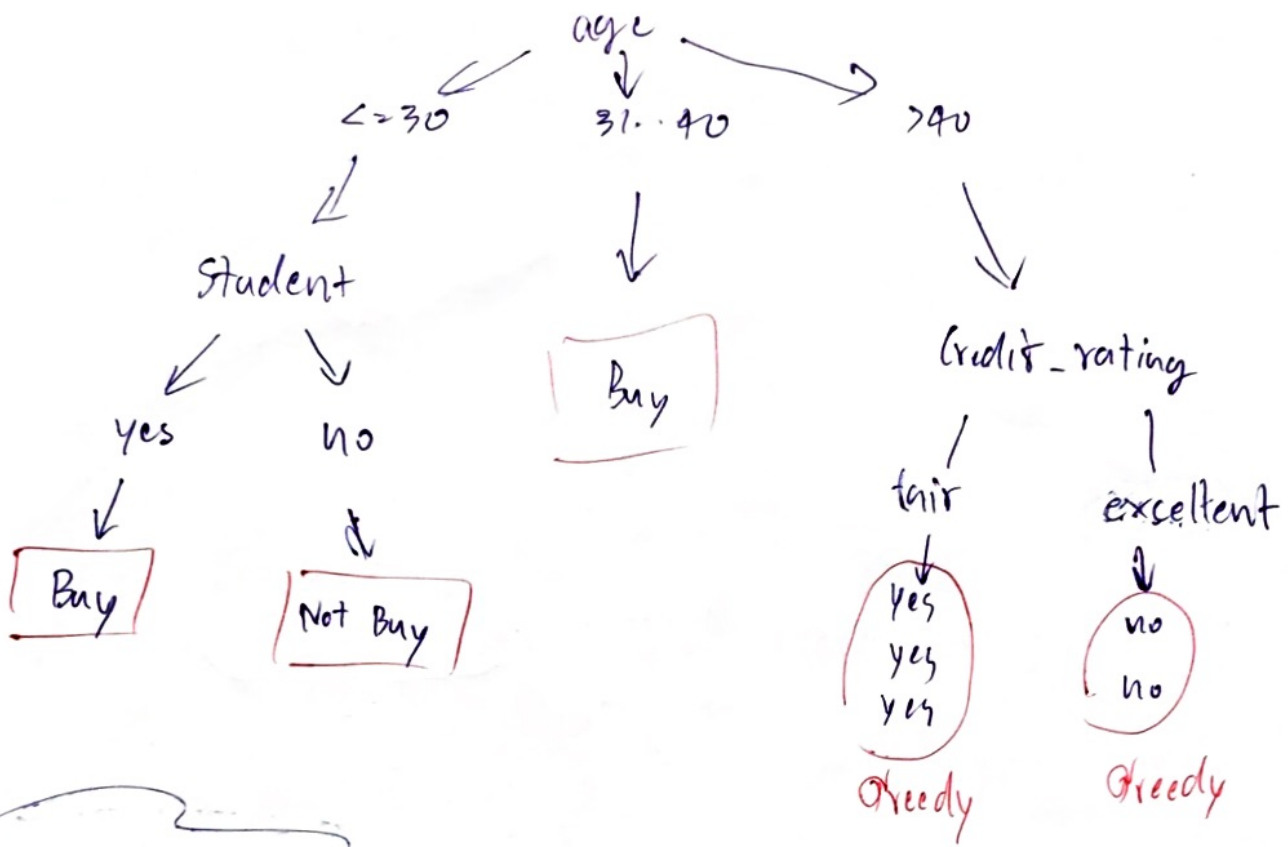
$$Gain(income) = 0.9710 - 0.4509 = 0.0201$$

$$Gain(student) = 0.9710 - 0.4509 = 0.0201$$

$$Gain(credit\_rating) = 0.9710 - 0 = 0.9710 \rightarrow \text{1st on credit\_rating}$$

166 root node





Decision Tree

