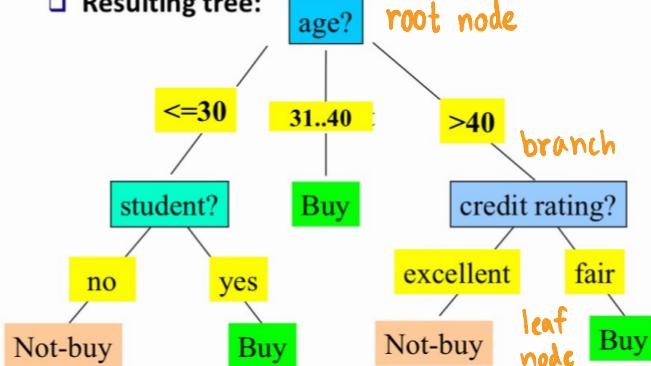


Decision Tree Induction - HW5

Resulting tree:



Training data set: Who buys computer?

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	
>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	

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

$$\text{Info}(D) = I(9,5) = -\frac{9}{14} \log_2 \left(\frac{9}{14}\right) - \frac{5}{14} \log_2 \left(\frac{5}{14}\right) = 0.904$$

$$2. \text{ Info}_A(D) = \sum_{j=1}^v \frac{|D_j|}{|D|} \times \text{Info}(D_j)$$

$$\text{Info}_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694$$

$$= \frac{5}{14} \left[-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right] + \frac{4}{14} \left[-\frac{4}{4} \log_2 \left(\frac{4}{4} \right) \right] + \frac{5}{14} \left[-\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) \right]$$

$$\begin{aligned}
 \text{Info income}^D &= \frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1) \\
 &= \frac{4}{14} \left[-\frac{2}{4} \log_2 \left(\frac{2}{4} \right) - \frac{2}{4} \log_2 \left(\frac{2}{4} \right) \right] + \\
 &\quad \frac{6}{14} \left[-\frac{4}{6} \log_2 \left(\frac{4}{6} \right) - \frac{2}{6} \log_2 \left(\frac{2}{6} \right) \right] + \\
 &\quad \frac{4}{14} \left[-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right] \\
 &= 0.911
 \end{aligned}$$

age	Pi	ni	I(Pi, ni)
<= 30	2	3	0.971
31... 40	4	0	0
>40	3	2	0.971

$$\begin{aligned}
 \text{Info student (D)} &= \frac{7}{14} I(6,1) + \frac{7}{14} I(3,4) \\
 &= \frac{7}{14} \left[-\frac{6}{7} \log_2 \left(\frac{6}{7} \right) - \frac{1}{7} \log_2 \left(\frac{1}{7} \right) \right] + \\
 &\quad \frac{7}{14} \left[-\frac{3}{7} \log_2 \left(\frac{3}{7} \right) - \frac{4}{7} \log_2 \left(\frac{4}{7} \right) \right] \\
 &= 0.789
 \end{aligned}$$

income	P_i	n_i	$I(P_i, n_i)$
high	2	2	1
median	4	2	0.918
low	3	1	0.811

$$\begin{aligned}
 \text{Info}_{\text{credit_rating}}(D) &= \frac{8}{14} I(6,2) + \frac{6}{14} I(3,3) \\
 &= \frac{8}{14} \left[-\frac{6}{8} \log_2 \left(\frac{6}{8}\right) - \frac{2}{8} \log_2 \left(\frac{2}{8}\right) \right] + \\
 &\quad \frac{6}{14} \left[-\frac{3}{6} \log_2 \left(\frac{3}{6}\right) - \frac{3}{6} \log_2 \left(\frac{3}{6}\right) \right] \\
 &= 0.892
 \end{aligned}$$

student	Pi	ni	$I(Pi, ni)$
yes	6	1	0.592
no	3	4	0.985

credit-r	Pi	ni	I(Pi, ni)
fair	6	2	0.811
excellent	3	3	1

$$3. \text{Gain}(A) = \text{Info}(D) - \text{Info}_A(D)$$

$$\text{Gain}(\text{age}) = 0.940 - 0.694 = 0.246 \quad \text{ជាអក្សរ (root node)}$$

$$\text{Gain}(\text{income}) = 0.940 - 0.911 = 0.029$$

$$\text{Gain}(\text{student}) = 0.940 - 0.789 = 0.151$$

$$\text{Gain}(\text{credit_rating}) = 0.940 - 0.892 = 0.048$$

4. ឈរកត្តុ feature ដែលត្រូវបានទាក់ទងនៅក្នុង feature age

≤ 30

$$\text{Info}(D) = I(2,3) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{3}{5} \log_2 \left(\frac{3}{5}\right) = 0.971$$

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

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

$$\text{Info}_{\text{credit_r}}(D) = \frac{3}{5} I(1,2) + \frac{2}{5} I(1,1) = 0.951$$

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
>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

$31\dots 40$

yes, no

$$\text{Info}(D) = I(4,0)$$

តាម 31...40 សមារតាម yes ឬ buys_computer ឪតើតុលិ

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
>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

>40

$$\text{Info}(D) = I(3,2) = -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) = 0.971$$

$$\text{Info}_{\text{income}}(D) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) = 0.951$$

$$\text{Info}_{\text{student}}(D) = \frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) = 0.951$$

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

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 ✓
>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

Information Gain = ?

$$\text{Gain}(\text{Income}) = 0.971 - 0.4 = 0.571$$

$$\text{Gain}(\text{student}) = 0.971 - 0 = 0.971$$

$$\text{Gain}(\text{credit_rating}) = 0.971 - 0.951 = 0.02$$

node ឱ្យ ≤ 30

$$\text{Gain}(\text{Income}) = 0.971 - 0.951 = 0.02$$

$$\text{Gain}(\text{student}) = 0.971 - 0.951 = 0.02$$

$$\text{Gain}(\text{credit_rating}) = 0.971 - 0 = 0.971$$

node ឱ្យ > 40

ឧបត្ថម្ភ Decision ចំណាំ

