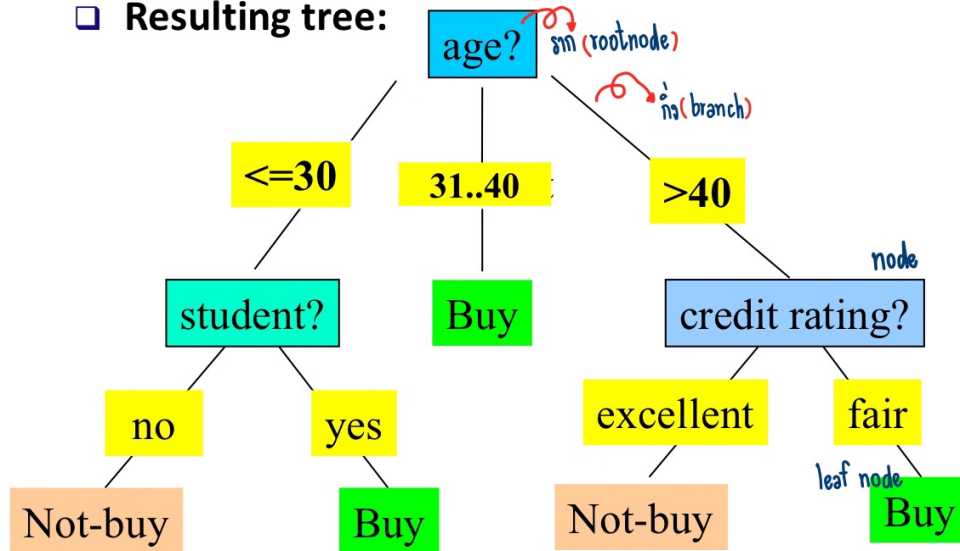


วิชาคอมพิวเตอร์ Decision - Tree

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

Class P: buys_computer = "yes"

Class N: buys_computer = "no"

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

$$\text{Info}(D) = I(9,5) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

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

$$1. \text{Info}_{\text{age}}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2)$$

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

$$= 0.694$$

$$2. \text{Info}_{\text{income}}(D) = \frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1)$$

$$= \frac{4}{14} \left[-\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) \right] + \frac{6}{14} \left[-\frac{4}{6} \log_2\left(\frac{4}{6}\right) - \frac{2}{6} \log_2\left(\frac{2}{6}\right) \right] + \frac{4}{14} \left[-\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \right]$$

$$= 0.911$$

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

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

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

3. คำนวณ Information Gain โดยที่ค่า Gain ที่มีค่าสูงที่สุดจะเป็นราก (root node)

$$3.1 \text{ Gain}(\text{age}) = 0.940 - 0.894 = 0.296$$

$$3.2 \text{ Gain}(\text{income}) = 0.940 - 0.911 = 0.029$$

$$3.3 \text{ Gain}(\text{student}) = 0.940 - 0.789 = 0.151$$

$$3.4 \text{ Gain}(\text{credit_rating}) = 0.940 - 0.892 = 0.048$$

∴ เลือก Gain (age) เป็นราก (root node) เนื่องจากมีค่าสูงที่สุดจากค่าทั้งหมด

4. แยกกลุ่มตัวอย่าง Feature ตามค่าในราก (root node)

4.1 ≤ 30

age	income	student	credit_rating	buys_computer
≤ 30	high	no	fair	no
≤ 30	high	no	excellent	no
≤ 30	medium	no	fair	no
≤ 30	low	yes	fair	yes
≤ 30	medium	yes	excellent	yes

$$\text{Info}(D) = I(2,3) = 0.971$$

$$\begin{aligned}
 \text{Info}_{\text{income}} &= \frac{2}{6} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(1,0) \\
 &= 0.4
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}_{\text{student}} &= \frac{2}{5} I(2,0) + \frac{3}{5} I(0,3) \\
 &= 0
 \end{aligned}$$

$$\begin{aligned}
 \text{Info}_{\text{credit_rating}} &= \frac{3}{5} I(1,2) + \frac{2}{6} I(1,1) \\
 &= 0.951
 \end{aligned}$$

คำนวณ Information Gain

$$\text{Gain}(\text{income}) = 0.971 - 0.400 = 0.571$$

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

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

∴ เลือก Gain (student) เป็น node ที่น้อยกว่า ≤ 30

4.2 31...40

age	income	student	credit_rating	buys_computer
31...40	high	no	fair	yes
31...40	low	yes	excellent	yes
31...40	medium	no	excellent	yes
31...40	high	yes	fair	yes

□ Class P: buys_computer = "yes" yes = 4
 □ Class N: buys_computer = "no" no = 0

เนื่องจาก 31...40 สามารถตอบ yes ใน buys_computer ได้เลย

4.3 > 40

age	income	student	credit_rating	buys_computer
>40	medium	no	fair	yes
>40	low	yes	fair	yes
>40	low	yes	excellent	no
>40	medium	yes	fair	yes
>40	medium	no	excellent	no

$$\text{Info}(0) = I\left(\frac{3}{5}, \frac{2}{5}\right) = -\frac{3}{5} \log_2 \left(\frac{3}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right)$$

$$= 0.971$$

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

$$= 0.951$$

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

$$= 0.951$$

$$\text{Info}_{(\text{credit_rating})} = \frac{3}{6} I(3,0) + \frac{2}{5} I(0,2)$$

$$= 0$$

คำนวณ Information Gain

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

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

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

∴ เลือก Gain (Credit_rating) เป็น node ในอายุ > 40

6. สามารถสร้าง Decision Tree ได้ดังนี้

