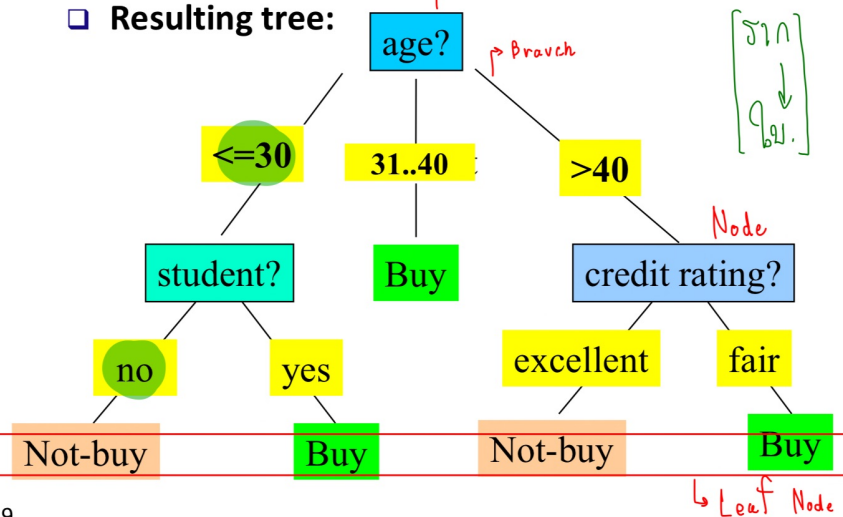


conquer process

Resulting tree:



Class P: buys_computer = "yes"

Class N: buys_computer = "no"

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

$$1.) Info(D) = -\sum_{i=1}^m p_i \log_2(p_i) \rightarrow Info(D) = -\frac{9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

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

$$2.1 Info_{age}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = \frac{5}{14} \left[-\frac{2}{5} \log_2\left(\frac{2}{5}\right) - \frac{3}{5} \log_2\left(\frac{3}{5}\right) \right] + \frac{4}{14} \left[-\frac{4}{4} \log_2\left(\frac{4}{4}\right) \right] + \frac{5}{14} \left[-\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) \right] = 0.694$$

$$2.2 Info_{income}(D) = \frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1) = \frac{4}{14} \left[-\frac{2}{4} \log_2\left(\frac{2}{4}\right) - \frac{2}{4} \log_2\left(\frac{2}{4}\right) \right] + \frac{6}{14} \left[-\frac{1}{6} \log_2\left(\frac{1}{6}\right) - \frac{5}{6} \log_2\left(\frac{5}{6}\right) \right] + \frac{4}{14} \left[-\frac{3}{4} \log_2\left(\frac{3}{4}\right) - \frac{1}{4} \log_2\left(\frac{1}{4}\right) \right] = 0.911$$

$$2.3 Info_{student}(D) = \frac{7}{14} I(1,1) + \frac{7}{14} I(3,4) = \frac{7}{14} \left[-\frac{1}{7} \log_2\left(\frac{1}{7}\right) - \frac{6}{7} \log_2\left(\frac{6}{7}\right) \right] + \frac{7}{14} \left[-\frac{1}{7} \log_2\left(\frac{1}{7}\right) - \frac{6}{7} \log_2\left(\frac{6}{7}\right) \right] = 0.789$$

$$2.4 Info_{credit_rating}(D) = \frac{8}{14} I(4,2) + \frac{6}{14} I(3,3) = \frac{8}{14} \left[-\frac{1}{4} \log_2\left(\frac{1}{4}\right) - \frac{3}{4} \log_2\left(\frac{3}{4}\right) \right] + \frac{6}{14} \left[-\frac{2}{6} \log_2\left(\frac{2}{6}\right) - \frac{4}{6} \log_2\left(\frac{4}{6}\right) \right] = 0.892$$

3. Gain(A) = Info(D) - Info_A(D) Gain Information Gain Total Gain Highest Gain (root node)

$$3.1 Gain(age) = 0.940 - 0.694 = 0.246 \rightarrow \text{Gain}(age) \text{ is the highest gain (root node) selected}$$

$$3.2 Gain(income) = 0.940 - 0.911 = 0.029$$

$$3.3 Gain(student) = 0.940 - 0.789 = 0.151$$

$$3.4 Gain(credit_rating) = 0.940 - 0.892 = 0.048$$

4. Feature Selection (root node)

$$4.1 <=30 \rightarrow Info(D) = I(9,5) = 0.971$$

$$\rightarrow Info_{income}(D) = \frac{2}{9} I(0,2) + \frac{2}{9} I(1,1) + \frac{1}{9} I(1,0) = 0.4$$

$$\rightarrow Info_{student}(D) = \frac{2}{9} I(2,0) + \frac{3}{9} I(0,3) = 0$$

$$\rightarrow Info_{credit}(D) = \frac{3}{9} I(2,2) + \frac{2}{9} I(1,1) = 0.951$$

4.2 31...40

$$4.3 >40 \rightarrow Info(D) = I(3,2) = -\frac{3}{5} \log_2\left(\frac{3}{5}\right) - \frac{2}{5} \log_2\left(\frac{2}{5}\right) = 0.971$$

$$\rightarrow Info_{income}(D) = \frac{3}{5} I(1,1) + \frac{2}{5} I(1,1) = 0.951$$

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

Information Gain

$$Gain(Income) = 0.971 - 0.4 = 0.571$$

$$Gain(student) = 0.971 - 0 = 0.971$$

$$Gain(credit_rating) = 0.971 - 0.951 = 0.02$$

* Gain(student) is the highest gain node for <=30

yes = 4, no = 0 selected 31...40

Information Gain

Information Gain

$$Gain(Income) = 0.971 - 0.951 = 0.02$$

$$Gain(student) = 0.971 - 0.951 = 0.02$$

$$Gain(credit_rating) = 0.971 - 0 = 0.971$$

* Gain(credit_rating) is the highest gain node for >40

5. វិធានការសម្រាប់ Decision Tree ក្នុងការសម្រេច

លេខកូដ ឈ្មោះ
633021019-9

