



CS 412 Intro. to Data Mining

Chapter 8. Classification: Basic Concepts

Jiawei Han, Computer Science, Univ. Illinois at Urbana-Champaign, 2017





Information Gain: An Attribute Selection Measure

- ❑ Select the attribute with the highest information gain (used in typical decision tree induction algorithm: ID3/C4.5)
- ❑ Let p_i be the probability that an arbitrary tuple in D belongs to class C_i , estimated by $|C_{i,D}|/|D|$
- ❑ Expected information (entropy) needed to classify a tuple in D :

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

- ❑ Information needed (after using A to split D into v partitions) to classify D :

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

- ❑ Information gained by branching on attribute A

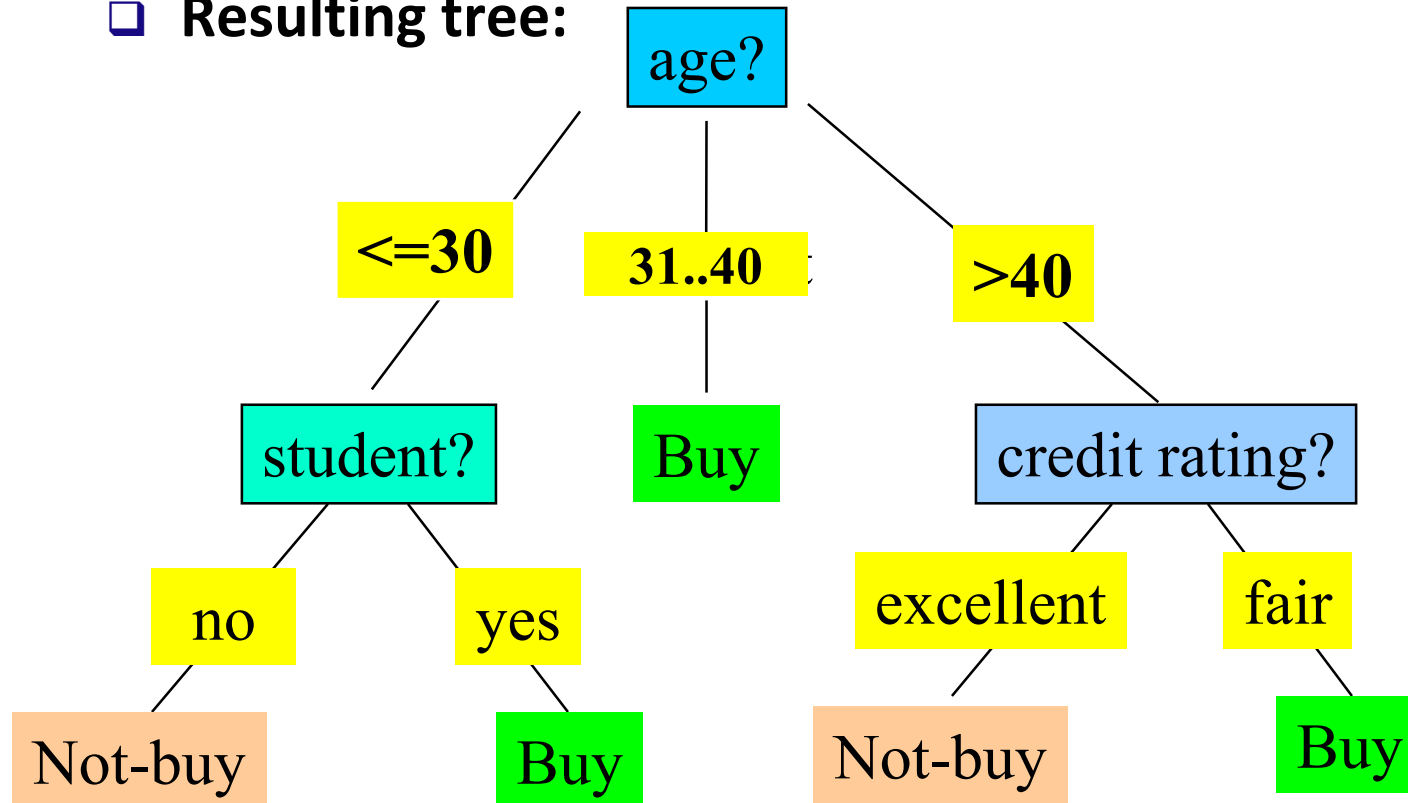
$$Gain(A) = Info(D) - Info_A(D)$$

Decision Tree Induction: An Example

Decision tree construction:

- A top-down, recursive, divide-and-conquer process

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

Note: The data set is adapted from
"Playing Tennis" example of R. Quinlan

Q7 Info (D)

$$\text{Info (D)} = 1 (0,5) = \overset{X}{-\frac{9}{14} \log_2 \left(\frac{9}{14} \right)} - \overset{N}{\frac{5}{14} \log_2 \left(\frac{5}{14} \right)}$$

$$= 0.94$$

Q7 Info age (D)

$$\text{Info age (D)} = \overset{L=30}{\frac{5}{14} I(2,3)} + \overset{37-10}{\frac{4}{14} I(4,0)} + \overset{740}{\frac{5}{14} I(3,2)}$$

$$I(2,3) = -\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) = 0.991$$

$$I(4,0) = -\frac{4}{4} \log_2 \left(\frac{4}{4} \right) - \frac{0}{4} \log_2 \left(\frac{0}{4} \right) = 0$$

$$I(3,2) = -\frac{3}{5} \log_2 \left(\frac{3}{5} \right) - \frac{2}{5} \log_2 \left(\frac{2}{5} \right) = 0.991$$

$$\text{Weighted Info age (D)} = \frac{5}{14} (0.991) + \frac{4}{14} (0) + \frac{5}{14} (0.991) = 0.694$$

Q7 Gain (age)

$$\text{Gain (age)} = 0.94 - 0.694 = 0.246$$

27 Info income (D)

$$\text{Info income (D)} = \overset{\text{high}}{\boxed{\frac{4}{14} I(2, 2)}} + \overset{\text{medium}}{\boxed{\frac{6}{14} I(4, 2)}} + \overset{\text{low}}{\boxed{\frac{4}{14} I(3, 1)}}$$

$$I(2, 2) = -\frac{2}{4} \log_2(2) \left(\frac{2}{4}\right) - \frac{2}{4} \log_2(2) \left(\frac{2}{4}\right) = 1$$

$$I(4, 2) = -\frac{4}{6} \log_2(2) \left(\frac{4}{6}\right) - \frac{2}{6} \log_2(2) \left(\frac{2}{6}\right) = 0.916$$

$$I(3, 1) = -\frac{3}{4} \log_2(2) \left(\frac{3}{4}\right) - \frac{1}{4} \log_2(2) \left(\frac{1}{4}\right) = 0.811$$

$$\text{28} \text{ } \text{Information Info income (D)} = \frac{4}{14}(1) + \frac{6}{14}(0.916) + \frac{4}{14}(0.811) = 0.911$$

29 Gain (income)

$$\text{Gain (income)} = 0.94 - 0.911 = 0.029$$

ii) Info student (D)

$$\text{Info student (D)} = \overset{\text{Yes}}{\frac{7}{10} I(6,1)} + \overset{\text{No}}{\frac{7}{14} I(3,4)}$$

$$I(6,1) = -\frac{6}{7} \log_2\left(\frac{6}{7}\right) - \frac{1}{7} \log_2\left(\frac{1}{7}\right) = 0.592$$

$$I(3,4) = -\frac{3}{7} \log_2\left(\frac{3}{7}\right) - \frac{4}{7} \log_2\left(\frac{4}{7}\right) = 0.985$$

$$\text{Weighted Info student (D)} = \frac{7}{14} (0.592) + \frac{7}{14} (0.985) = 0.749$$

iii) Gain (student)

$$\text{Gain (student)} = 0.94 - 0.749 = 0.191$$

ก) Info credit-rating (D)

$$\text{Info credit-rating (D)} = \frac{8}{16} I(\overset{\text{fair}}{b, 2}) + \frac{6}{16} I(\overset{\text{excellent}}{3, 3})$$

$$I(\overset{\text{fair}}{b, 2}) = -\frac{1}{8} \log_2\left(\frac{1}{8}\right) - \frac{1}{8} \log_2\left(\frac{2}{8}\right) = 0.41$$

$$I(\overset{\text{excellent}}{3, 3}) = -\frac{3}{6} \log_2\left(\frac{3}{6}\right) - \frac{3}{6} \log_2\left(\frac{3}{6}\right) = 1$$

$$\text{รวม Info credit-rating (D)} = \frac{8}{16} (0.41) + \frac{6}{16} (1) = 0.492$$

ข) Gain (credit-rating)

$$\text{Gain (credit-rating)} = 0.94 - 0.492 = 0.048$$

ก) Gain

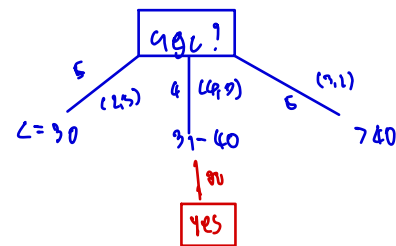
$$\text{Gain (age)} = 0.246$$

$$\text{Gain (income)} = 0.029$$

$$\text{Gain (student)} = 0.151$$

$$\text{Gain (credit-rating)} = 0.048$$

เลือก Gain ที่มากที่สุดคือ age เป็นตัวแรก จึงได้ใช่คือ Gain (age)



$$\text{age} = (L=30)$$

or $\text{info}(D)$ for $\text{age} (L=30)$

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

or $\text{info}_{\text{income}}(D)$ for age

$$\text{info}_{\text{income}}(D) \text{ for } \text{age} (L=30) = \overset{\text{high}}{\frac{2}{5} I(\vec{0}, \vec{2})} + \overset{\text{medium}}{\frac{2}{5} I(\vec{1}, \vec{1})} + \overset{\text{low}}{\frac{1}{5} I(\vec{1}, \vec{0})}$$

$$I(\vec{0}, \vec{2}) = -\frac{0}{2} \log_2\left(\frac{0}{2}\right) - \frac{2}{2} \log_2\left(\frac{1}{2}\right) = 0$$

$$I(\vec{1}, \vec{1}) = -\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right) = 1$$

$$I(\vec{1}, \vec{0}) = -\frac{1}{1} \log_2\left(\frac{1}{1}\right) - \frac{0}{1} \log_2\left(\frac{0}{1}\right) = 0$$

$$\text{which } \text{info}_{\text{income}}(D) \text{ for } \text{age} (L=30) = \frac{2}{5}(0) + \frac{2}{5}(1) + \frac{1}{5}(0) = 0.4$$

or $\text{Gain}(\text{income})$ for $\text{age} (L=30)$

$$\text{Gain}(\text{income}) \text{ for } \text{age} (L=30) = 0.971 - 0.4 = 0.571$$

၁၇ Info student (D) နှင့် age (<30)

$$\text{Info student (D) နှင့် age (<30)} = \overset{\text{yes}}{\frac{2}{5} I(2,0)} + \overset{\text{no}}{\frac{3}{5} I(0,3)}$$

ဆိုရင် Yes \rightarrow yes (buy - computer), No \rightarrow (buy - laptop)

ရေတွက်ရန် Student ကို အသက်ပေါ်မူတည်၍ အုပ်စုခွဲခြားရ

၁၈ age (>40)

၁၇ Info (D) နှင့် age (>40)

$$\text{Info (D) နှင့် age (>40)} = I(3,2) = 0.979$$

၁၉ Info income (D) နှင့် age (>40)

$$\text{Info income (D) နှင့် age (>40)} = \overset{\text{medium}}{\frac{2}{5} I(2,1)} + \overset{\text{low}}{\frac{1}{5} I(1,1)}$$

$$I(2,1) = -\frac{2}{3} \log_2 \left(\frac{2}{3}\right) - \frac{1}{3} \log_2 \left(\frac{1}{3}\right) = 0.916$$

$$I(1,1) = 1$$

$$\text{မူလ Info income (D) နှင့် age (>40)} = \frac{2}{5} (0.916) + \frac{1}{5} (1) = 0.951$$

၁၉ Gain (income) နှင့် age (>40)

$$\text{Gain (income) နှင့် age (>40)} = 0.979 - 0.951 = 0.02$$

၁၈) Info student (၀) ကို အသွယ် (>40)

$$\text{Info student (၀) ကို အသွယ် (>40)} = \overset{\text{yes}}{\frac{3}{5} I(2,1)} + \overset{\text{no}}{\frac{2}{5} I(1,1)}$$

$$I(2,1) = -\frac{2}{3} \log_2(1) \log_2(1) \log_2\left(\frac{2}{3}\right) - \frac{1}{3} \log_2(2) \log_2\left(\frac{1}{3}\right) = 0.914$$

$$I(1,1) = 1$$

$$\text{မကွာ Info student (၀) ကို အသွယ် (>40)} = \frac{3}{5} (0.914) + \frac{2}{5} (1) = 0.951$$

၁၉) Gain (student) ကို အသွယ် (>40)

$$\text{Gain (student) အသွယ် (>40)} = 0.971 - 0.951 = 0.02$$

၂၀) Info credit-rating (၀) ကို အသွယ် (>40)

$$\text{Info credit-rating (၀) ကို အသွယ် (>40)} = \overset{\text{fair}}{\frac{3}{5} I(3,0)} + \overset{\text{excellent}}{\frac{2}{5} I(0,1)}$$

အသွယ် fair → yes (buy-computer), excellent → no (buy-computer)

မလွန်သော credit-rating ကို အသွယ် အသွယ် အသွယ် အသွယ် အသွယ်

9/11

