

**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

Yes = 9, no = 5

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

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

Info<sub>A</sub>(D) =  $\sum_{j=1}^n \frac{|D_j|}{|D|} \times \text{Info}(D_j)$

1. Info<sub>age</sub>(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. Info<sub>income</sub>(D) =  $\frac{4}{14} I(2,2) + \frac{6}{14} I(4,2) + \frac{4}{14} I(3,1)$   
 $= \frac{4}{14} \left[ -\frac{2}{4} \log_2 \left( \frac{2}{4} \right) - \frac{2}{4} \log_2 \left( \frac{2}{4} \right) \right] + \frac{6}{14} \left[ -\frac{4}{6} \log_2 \left( \frac{4}{6} \right) - \frac{2}{6} \log_2 \left( \frac{2}{6} \right) \right] + \frac{4}{14} \left[ -\frac{3}{4} \log_2 \left( \frac{3}{4} \right) - \frac{1}{4} \log_2 \left( \frac{1}{4} \right) \right] = 0.911$

3. Info<sub>student</sub>(D) =  $\frac{7}{14} I(6,1) + \frac{7}{14} I(3,4)$   
 $= \frac{7}{14} \left[ -\frac{6}{7} \log_2 \left( \frac{6}{7} \right) - \frac{1}{7} \log_2 \left( \frac{1}{7} \right) \right] + \frac{7}{14} \left[ -\frac{3}{7} \log_2 \left( \frac{3}{7} \right) - \frac{4}{7} \log_2 \left( \frac{4}{7} \right) \right] = 0.989$

4. Info<sub>credit rating</sub>(D) =  $\frac{8}{14} I(6,2) + \frac{6}{14} I(3,3)$   
 $= \frac{8}{14} \left[ -\frac{6}{8} \log_2 \left( \frac{6}{8} \right) - \frac{2}{8} \log_2 \left( \frac{2}{8} \right) \right] + \frac{6}{14} \left[ -\frac{3}{6} \log_2 \left( \frac{3}{6} \right) - \frac{3}{6} \log_2 \left( \frac{3}{6} \right) \right] = 0.892$

### Information Gain

Gain(A) = Info(D) - Info<sub>A</sub>(D)

1. Gain(age) =  $0.940 - 0.694 = 0.246$  → decision attribute with highest gain

2. Gain(income) =  $0.940 - 0.911 = 0.029$

3. Gain(student) =  $0.940 - 0.989 = -0.049$

4. Gain(credit rating) =  $0.940 - 0.892 = 0.048$

• 11th feature minimum root node

1.  $1 \leq 30$

Info(D) =  $I(3,3) = 0.971$

Info<sub>income</sub>(D) =  $\frac{2}{5} I(0,2) + \frac{2}{5} I(1,1) + \frac{1}{5} I(2,0) = 0.4$

Info<sub>student</sub>(D) =  $\frac{2}{5} I(2,0) + \frac{3}{5} I(0,3) = 0$

Info<sub>credit</sub>(D) =  $\frac{2}{5} I(1,2) + \frac{3}{5} I(1,1) = 0.951$

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

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

Gain(Credit-rating) =  $0.971 - 0.951 = 0.02$

$\therefore$  When Gain(student) is node  $\leq 30$

2.  $31 \dots 40$

Yes = 4, No = 0

$\therefore$  31..40 minimum yes to no

3.  $> 40$

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

Info<sub>income</sub> =  $\frac{3}{5} I(1,1) + \frac{2}{5} I(1,1) = 0.951$

Info<sub>student</sub> =  $\frac{3}{5} I(2,1) + \frac{2}{5} I(1,1) = 0.951$

Info<sub>credit-rating</sub> =  $\frac{3}{5} I(3,0) + \frac{2}{5} I(0,2) = 0$

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

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

Gain(Credit-rating) =  $0.971 - 0 = 0.971$

$\therefore$  When Gain(Credit-rating) node  $> 40$

• Decision Tree Algorithm

