

HW2: Attribute Selection with Information Gain

$$\text{Info}(D) = I(9,5) = \frac{9}{14} \log_2 \left(\frac{9}{14} \right) - \frac{5}{14} \log_2 \left(\frac{5}{14} \right) = 0.940 - \text{Expected information (entropy) for the whole}$$

$$\text{Info}_{\text{age}}(D) = \frac{5}{14} I(2,3) + \frac{4}{14} I(4,0) + \frac{5}{14} I(3,2) = 0.694 - \text{Expected information (entropy) within part node}$$

$$\begin{aligned} \text{Info}_{\text{income}}(D) &= \frac{5}{14} I^{\text{high}}(1,2) + \frac{4}{14} I^{\text{low}}(4,1) + \frac{5}{14} I^{\text{medium}}(2,2) \\ &= \frac{5}{14} \left(-\frac{1}{5} \log_2 \left(\frac{1}{5} \right) - \frac{4}{5} \log_2 \left(\frac{4}{5} \right) \right) + \frac{4}{14} \left(-\frac{3}{4} \log_2 \left(\frac{3}{4} \right) - \frac{1}{4} \log_2 \left(\frac{1}{4} \right) \right) + \frac{5}{14} \left(-\frac{2}{5} \log_2 \left(\frac{2}{5} \right) - \frac{3}{5} \log_2 \left(\frac{3}{5} \right) \right) \\ &= 0.911 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{student}}(D) &= \frac{7}{14} I^{\text{no}}(4,5) + \frac{7}{14} I^{\text{yes}}(6,1) \\ &= \frac{7}{14} \left(-\frac{4}{7} \log_2 \left(\frac{4}{7} \right) - \frac{3}{7} \log_2 \left(\frac{3}{7} \right) \right) + \frac{7}{14} \left(-\frac{6}{7} \log_2 \left(\frac{6}{7} \right) - \frac{1}{7} \log_2 \left(\frac{1}{7} \right) \right) \\ &= 0.788 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{credit-rating}}(D) &= \frac{8}{14} \left(-\frac{6}{8} \log_2 \left(\frac{6}{8} \right) - \frac{2}{8} \log_2 \left(\frac{2}{8} \right) \right) + \frac{6}{14} \left(-\frac{3}{6} \log_2 \left(\frac{3}{6} \right) - \frac{3}{6} \log_2 \left(\frac{3}{6} \right) \right) \\ &= \frac{8}{14} I^{\text{fair}}(6,2) + \frac{6}{14} I^{\text{excellent}}(3,3) \\ &= 0.992 \end{aligned}$$

$$\text{Gain}(\text{age}) = \text{Info}(D) - \text{Info}_{\text{age}}(D) = 0.940 - 0.694 = 0.246$$

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

$$\text{Gain}(\text{student}) = \text{Info}(D) - \text{Info}_{\text{student}}(D) = 0.940 - 0.788 = 0.152$$

$$\text{Gain}(\text{Credit-rating}) = \text{Info}(D) - \text{Info}_{\text{credit-rating}}(D) = 0.940 - 0.992 = 0.048$$

Algorithm Root Node is the attribute with the highest Gain. In this case, age has the highest Gain = 0.246

in age: < 30

$$\text{Info}_{\text{age}}^{<30}(D) = I(3,3) = \frac{3}{6} \log_2 \left(\frac{3}{6} \right) - \frac{3}{6} \log_2 \left(\frac{3}{6} \right) = 0.921$$

$$\begin{aligned} \text{Info}_{\text{income}}^{<30}(D) &= \frac{1}{5} I(1,0) + \frac{2}{5} I(1,1) + \frac{2}{5} I(1,2) \\ &= \frac{1}{5} \left(-\frac{1}{1} \log_2 \left(\frac{1}{1} \right) - 0 \right) + \frac{2}{5} \left(-\frac{1}{2} \log_2 \left(\frac{1}{2} \right) - \frac{1}{2} \log_2 \left(\frac{1}{2} \right) \right) + \frac{2}{5} \left(-\frac{2}{2} \log_2 \left(\frac{2}{2} \right) - 0 \right) \end{aligned}$$

$$\text{Info}_{\text{student}}^{<30}(D) = \frac{3}{6} I(0,3) + \frac{3}{6} I(3,0) = \frac{3}{6} \left(-\frac{3}{3} \log_2 \left(\frac{3}{3} \right) - 0 \right) + \frac{3}{6} \left(-\frac{3}{3} \log_2 \left(\frac{3}{3} \right) - 0 \right)$$

= 0



$$\begin{aligned} \text{Info}_{\text{credit-rating}}(D) &= \frac{3}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{2}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) \\ &= \frac{3}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) + \frac{2}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) \\ &= 0.951 \end{aligned}$$

$$\text{Gain}_{\text{income}} = \text{Info}_{\text{age} < 30}(D) - \text{Info}_{\text{income}}(D) = 0.971 - 0.900 = 0.071$$

$$\text{Gain}_{\text{student}} = \text{Info}_{\text{age} < 30}(D) - \text{Info}_{\text{student}}(D) = 0.971 - 0 = 0.971$$

$$\text{Gain}_{\text{credit-rating}} = \text{Info}_{\text{age} < 30}(D) - \text{Info}_{\text{credit-rating}}(D) = 0.971 - 0.951 = 0.020$$

After decision made on Age student gives the gain

$$\begin{aligned} \text{Info}_{\text{age} > 30 \dots 40}(D) &= I\left(\frac{4}{8}, \frac{0}{8}\right) \\ &= -\frac{4}{8} \log_2\left(\frac{4}{8}\right) - \frac{0}{8} \log_2\left(\frac{0}{8}\right) \end{aligned}$$

After 31...40 the decision made not necessary

$$I\left(\frac{4}{8}, \frac{0}{8}\right) \text{ after } \text{age} > 31 \dots 40 \text{ not necessary}$$

$$\text{Info}_{\text{age} > 30}(D) = I\left(\frac{4}{8}, \frac{0}{8}\right) = -\frac{4}{8} \log_2\left(\frac{4}{8}\right) - \frac{0}{8} \log_2\left(\frac{0}{8}\right) = 0.971$$

$$\begin{aligned} \text{Info}_{\text{income}}(D) &= \frac{1}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{2}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) \\ &= \frac{1}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) + \frac{2}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) \\ &= 0.951 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{student}}(D) &= \frac{1}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{4}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) = \frac{1}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) + \frac{4}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) \\ &= 0.951 \end{aligned}$$

$$\begin{aligned} \text{Info}_{\text{credit-rating}}(D) &= \frac{1}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) + \frac{2}{5} I\left(\frac{1}{2}, \frac{1}{2}\right) \\ &= \frac{1}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) + \frac{2}{5} \left(-\frac{1}{2} \log_2\left(\frac{1}{2}\right) - \frac{1}{2} \log_2\left(\frac{1}{2}\right)\right) = 0 \end{aligned}$$

$$\text{Gain}_{\text{income}} = \text{Info}_{\text{age} > 30}(D) - \text{Info}_{\text{income}}(D) = 0.971 - 0.951 = 0.02$$

$$\text{Gain}_{\text{student}} = \text{Info}_{\text{age} > 30}(D) - \text{Info}_{\text{student}}(D) = 0.971 - 0.951 = 0.02$$

$$\text{Gain}_{\text{credit-rating}} = \text{Info}_{\text{age} > 30}(D) - \text{Info}_{\text{cr}}(D) = 0.971 - 0 = 0.971$$

After Credit-rating the decision made on income gives the gain

= 0 if necessary

