

ISOM3360 Assignment: Decision Tree

Zhang Yichen

March 24, 2020

1 Answer

The feature ***Snow*** should be the first split condition for root node as it will create a highest information gain after split.

2 Procedures

According to ID3 algorithm, the feature to split the root should be able to split it into subsets for which the information gain is maximum. Hence, I will compare the information gain by using different feature as the root node splitting condition. I denote the root node set as ***S*** and the two subsets as ***T*₁** and ***T*₂**.

Select ***Snow*** as the root node

$$\begin{aligned} IG(feature = \mathbf{Snow}, \mathbf{S}) &= H(\mathbf{S}) - p(\mathbf{T}_1)H(\mathbf{T}_1) - p(\mathbf{T}_2)H(\mathbf{T}_2) \\ &= -\left(\frac{6}{10} \log_2\left(\frac{6}{10}\right) + \frac{4}{10} \log_2\left(\frac{4}{10}\right)\right) + \frac{6}{10} * \left(\frac{1}{6} \log_2\left(\frac{1}{6}\right) + \frac{5}{6} \log_2\left(\frac{5}{6}\right)\right) + \frac{4}{10} * \left(\frac{1}{4} \log_2\left(\frac{1}{4}\right) + \frac{3}{4} \log_2\left(\frac{3}{4}\right)\right) \\ &= 0.971 - 0.390 - 0.325 \\ &= \mathbf{0.256} \end{aligned}$$

Select ***Season*** as the root node

$$\begin{aligned} IG(feature = \mathbf{Season}, \mathbf{S}) &= H(\mathbf{S}) - p(\mathbf{T}_1)H(\mathbf{T}_1) - p(\mathbf{T}_2)H(\mathbf{T}_2) \\ &= -\left(\frac{6}{10} \log_2\left(\frac{6}{10}\right) + \frac{4}{10} \log_2\left(\frac{4}{10}\right)\right) + \frac{5}{10} * \left(\frac{1}{5} \log_2\left(\frac{1}{5}\right) + \frac{4}{5} \log_2\left(\frac{4}{5}\right)\right) + \frac{5}{10} * \left(\frac{2}{5} \log_2\left(\frac{2}{5}\right) + \frac{3}{5} \log_2\left(\frac{3}{5}\right)\right) \\ &= 0.971 - 0.361 - 0.485 \\ &= \mathbf{0.125} \end{aligned}$$

Select *Weather* as the root node

$$\begin{aligned} IG(feature = \mathbf{Weather}, \mathbf{S}) &= H(\mathbf{S}) - p(\mathbf{T}_1)H(\mathbf{T}_1) - p(\mathbf{T}_2)H(\mathbf{T}_2) \\ &= -(\frac{6}{10} \log_2(\frac{6}{10}) + \frac{4}{10} \log_2(\frac{4}{10})) + \frac{4}{10} * (\frac{2}{4} \log_2(\frac{2}{4}) \\ &\quad + \frac{2}{4} \log_2(\frac{2}{4})) + \frac{6}{10} * (\frac{2}{6} \log_2(\frac{2}{6}) + \frac{4}{6} \log_2(\frac{4}{6})) \\ &= 0.971 - 0.4 - 0.551 \\ &= \mathbf{0.020} \end{aligned}$$

Therefore,

$$IG(feature = \mathbf{Snow}, \mathbf{S}) > IG(feature = \mathbf{Season}, \mathbf{S}) > IG(feature = \mathbf{Weather}, \mathbf{S})$$

The feature *Snow* should be the first split condition for root node as it will creat a highest information gain after split.