

Probability

Temp - Cool, Hot, Mild

Cricket - Yes, No

Total Rows = 20

Temp --- Yes --- No

Cool --- 5 --- 2 --- 7/20

Hot --- 3 --- 4 --- 7/20

Mild --- 3 --- 3 --- 6/20

| |

Total - 11/20 - 9/20

Decision Tree

Entropy - To check the purity of column values
(Categories)

Information Gain - To check how much
information a column is gaining

Gini index - Same as Entropy, will check the purity of column values (Categories)

In Dataset -

Independent column, input column, input features are the same - Can be many

Similarly, Target Column, output column, dependent feature are the same. - Only one

Dataset - Play Tennis -

Input Column - Outlook, Temperature, Humidity, Wind

Output Column - Play Tennis

Prob(Yes) - 9/14

Prob(No) - 5/14

For complete Dataset

$\text{Prob}(\text{Sunny}) - 5/14 - \text{Prob}(\text{Sunny} | \text{Yes}) - 2/5 -$
 $\text{Prob}(\text{Sunny} | \text{No}) - 3/5$

$\text{Prob}(\text{Overcast}) - 4/14 - \text{Prob}(\text{Over} | \text{Yes}) - 4/4 -$
 $\text{Prob}(\text{Over} | \text{No}) - 0/4$

$\text{Prob}(\text{Rain}) - 5/14 - \text{Prob}(\text{Rain} | \text{Yes}) - 3/5 -$
 $\text{Prob}(\text{Rain} | \text{No}) - 2/5$

We will calculate the entropy and information gain for every column

To create the tree we have to start from root node,

To decide which column should be root node we will consider the max information gain of column

Entropy range from 0 to 1

0 means pure and 1 means impure

Entropy = $-\sum p \times \log_2(p)$

$$\begin{aligned}\text{Entropy} &= -P_{\text{yes}} \log_2(P_{\text{yes}}) - P_{\text{no}} \log_2(P_{\text{no}}) \\ &= -(9/14) * \log_2(9/14) - (5/14) * \log_2(5/14) \\ &= 0.94\end{aligned}$$

$$\begin{aligned}\text{Entropy(sunny)} &= -P_{\text{yes}} \log_2(P_{\text{yes}}) - P_{\text{no}} \log_2(P_{\text{no}}) \\ &= -(2/5) * \log_2(2/5) - (3/5) * \log_2(3/5) \\ &= 0.97\end{aligned}$$

$$\begin{aligned}\text{Entropy(over.)} &= -P_{\text{yes}} \log_2(P_{\text{yes}}) - P_{\text{no}} \log_2(P_{\text{no}}) \\ &= -(4/4) * \log_2(4/4) - (0/4) * \log_2(0/4) \\ &= 0\end{aligned}$$

$$\begin{aligned}\text{Entropy(Rain)} &= -P_{\text{yes}} \log_2(P_{\text{yes}}) - P_{\text{no}} \log_2(P_{\text{no}}) \\ &= -(3/5) * \log_2(3/5) - (2/5) * \log_2(2/5) \\ &= 0.97\end{aligned}$$

Information Gain =

Overall Entropy – Weighted Entropy of column values

$$= E(\text{table}) - [(\text{sunny}/\text{total}) * E(\text{sunny}) + (\text{over}/\text{total}) * E(\text{over}) + (\text{rain}/\text{total}) * E(\text{rain})]$$

$$= 0.94 - [(5/14) * 0.97 + (4/14) * 0 + (5/14) * 0.97]$$

$$= 0.94 - 0.69$$

$$IG(\text{Outlook}) = 0.25$$