

# Decision Tree Classifier with Example

## ID3 Approach

outlook	temp	humidity	windy	play
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
overcast	hot	high	FALSE	yes
rainy	mild	high	FALSE	yes
rainy	cool	normal	FALSE	yes
rainy	cool	normal	TRUE	no
overcast	cool	normal	TRUE	yes
sunny	mild	high	FALSE	no
sunny	cool	normal	FALSE	yes
rainy	mild	normal	FALSE	yes
sunny	mild	normal	TRUE	yes
overcast	mild	high	TRUE	yes
overcast	hot	normal	FALSE	yes
rainy	mild	high	TRUE	no

Entropy is defined as the measure of randomness in the information being processed. It is given by the formula

$$H(C) = \text{Entropy} = -P_y \log_2 P_y - P_n \log_2 P_n$$

Information Gain is defined as

$$\text{Gain} = H(C) - \sum \frac{|S_v|}{S} H(S_v)$$

We consider that node as the root node of a Decision Tree which has the highest Information Gain.

Calculation:-

Target variable Play has 9 yes & 5 Nos.

$$\begin{aligned} \text{Entropy(Play)} &= -\frac{9}{14} \log_2 \left( \frac{9}{14} \right) - \frac{5}{14} \log_2 \left( \frac{5}{14} \right) \\ &= -\frac{9}{14} \left[ \frac{\log_{10} \left( \frac{9}{14} \right)}{\log_{10} 2} \right] - \frac{5}{14} \left[ \frac{\log_{10} \left( \frac{5}{14} \right)}{\log_{10} 2} \right] \\ &= -0.64 \left[ \frac{-0.19}{0.3010} \right] - 0.35 \left[ \frac{-0.44}{0.3010} \right] \\ &= 0.4016 + 0.51 \\ &\approx 0.91 \end{aligned}$$

Values (outlook) = sunny, overcast, rainy

$$\text{play} = [9Y, 5No]$$

$$\text{play}_{\text{sunny}} = [2Y, 3No]$$

$$\text{play}_{\text{overcast}} = [4Y, 0No]$$

$$\text{play}_{\text{rainy}} = [3Y, 2No]$$

$$\therefore IG(\text{Play}, \text{outlook}) = \text{Entropy(Play)} -$$

$$\sum_{v \in [\text{sunny}, \text{overcast}, \text{rainy}]} \frac{|S_v|}{S} \text{Entropy}(S_v)$$

$$\begin{aligned} \left[ \begin{aligned} \text{Entropy(Play}_{\text{sunny}}) &= -\frac{2}{5} \log_2 \left( \frac{2}{5} \right) - \frac{3}{5} \log_2 \left( \frac{3}{5} \right) \\ &= +0.52 + 0.44 \\ &= +0.96 \end{aligned} \right] &= 0.91 - \frac{5}{14} \text{Entropy(Play}_{\text{sunny}}) \\ &\quad - \frac{4}{14} \text{Entropy(Play}_{\text{overcast}}) \\ &\quad - \frac{5}{14} \text{Entropy(Play}_{\text{rainy}}) \end{aligned}$$

$$E(\text{Play}_{\text{overcast}}) = 0$$

outlook	temp	humidity	windy	play
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
overcast	hot	high	FALSE	yes
rainy	mild	high	FALSE	yes
rainy	cool	normal	FALSE	yes
rainy	cool	normal	TRUE	no
overcast	cool	normal	TRUE	yes
sunny	mild	high	FALSE	no
sunny	cool	normal	FALSE	yes
rainy	mild	normal	FALSE	yes
sunny	mild	normal	TRUE	yes
overcast	mild	high	TRUE	yes
overcast	hot	normal	FALSE	yes
rainy	mild	high	TRUE	no

outlook	temp	humidity	windy	play
sunny	hot	high	FALSE	no
sunny	hot	high	TRUE	no
overcast	hot	high	FALSE	yes
rainy	mild	high	FALSE	yes
rainy	cool	normal	FALSE	yes
rainy	cool	normal	TRUE	no
overcast	cool	normal	TRUE	yes
sunny	mild	high	FALSE	no
sunny	cool	normal	FALSE	yes
rainy	mild	normal	FALSE	yes
sunny	mild	normal	TRUE	yes
overcast	mild	high	TRUE	yes
overcast	hot	normal	FALSE	yes
rainy	mild	high	TRUE	no

$$E(\text{Play}, \text{rainy}) = -\frac{2}{5} \log_2 \left(\frac{2}{5}\right) - \frac{2}{5} \log_2 \left(\frac{2}{5}\right) = 0.44 + 0.44 = 0.88$$

$$\therefore \text{Information Gain}(\text{Play}, \text{outlook}) = 0.94 - \frac{5}{14} (2 \times 0.88) = 0.254$$

Similarly we compute information gain for wind feature w.r.t target play

Values (Wind) : FALSE, TRUE

play = [4Y, 5No]

play<sub>FALSE</sub> = [6Y, 2No]

play<sub>TRUE</sub> = [3Y, 3No]

$\therefore \text{Information Gain} = \text{Entropy}(\text{Play}) -$

$$\sum_{v \in \{\text{True}, \text{False}\}} \frac{|S_v|}{S} \text{Entropy } S_v$$

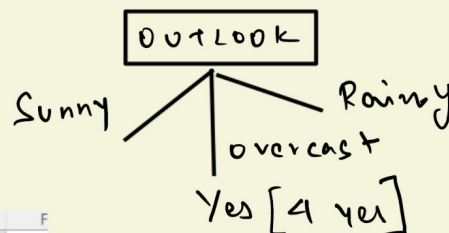
Similarly

$$\text{I.G}(\text{Play}, \text{Temp}) = 0.029$$

$$\text{I.G}(\text{Play}, \text{Humidity}) = 0.152$$

$$\begin{aligned} &= 0.94 - \frac{8}{14} \text{Entropy}(\text{Play}_{\text{FALSE}}) \\ &\quad - \frac{6}{14} \text{Entropy}(\text{Play}_{\text{TRUE}}) \\ &= 0.94 - \frac{8}{14} (0.811) - \left(\frac{6}{14}\right) \cdot \\ &= 0.048 \end{aligned}$$

Since outlook feature has the highest I.G, we choose it as root node.  
 $\therefore$  our decision tree currently looks as below



	A	B	C	D	E	F
1	outlook	temp	humidity	windy	play	
2	sunny	hot	high	FALSE	no	
3	sunny	hot	high	TRUE	no	
9	sunny	mild	high	FALSE	no	
10	sunny	cool	normal	FALSE	yes	
12	sunny	mild	normal	TRUE	yes	
16						
17						

Next we consider this table on the left

	A	B	C	D	E	F
1	outlook	temp	humidity	windy	play	
2	sunny	hot	high	FALSE	no	
3	sunny	hot	high	TRUE	no	
9	sunny	mild	high	FALSE	no	
10	sunny	cool	normal	FALSE	yes	
12	sunny	mild	normal	TRUE	yes	
16						
17						

$$\text{Entropy}(\text{Sunny}) = -P_Y \log_2 P_Y - P_N \log_2 P_N$$

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

$$= 0.971$$

Now we calculate the new I.G.

$$IG(\text{Sunny}, \text{temp}) = E(\text{Sunny}) - \sum_{v \in [\text{hot}, \text{mild}, \text{cold}]} \frac{|S_v|}{S} \text{Entropy}(S_v)$$

Values (temp) = hot, mild, cold

play = [2Y, 3No]

play<sub>hot</sub> = [0Y, 2No]

play<sub>mild</sub> = [0Y, 1No]

play<sub>cold</sub> = [2Y, 0No]

$$\begin{aligned} \therefore IG(\text{Sunny}, \text{temp}) &= 0.971 - \frac{2}{5} \text{Entropy}(\text{play}_{\text{hot}}) - \frac{1}{5} \text{Entropy}(\text{play}_{\text{mild}}) \\ &\quad - \frac{2}{5} \text{Entropy}(\text{play}_{\text{cold}}) \\ &= 0.971 - 0.4 = 0.571 \end{aligned}$$

Similarly  $IG(\text{Sunny}, \text{Humidity}) = 0.971$

$IG(\text{Sunny}, \text{Windy}) = 0.02$

Hence we choose Humidity as a node under sunny.

Similarly we can find nodes under rainy

Hence currently our D.T looks like one below.

