

Introduction to Decision Tree

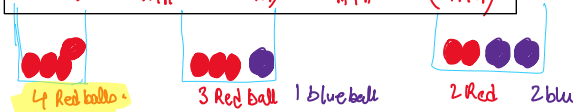
Wednesday, 17 April 2024 3:48 PM

Decision Tree

A Decision Tree takes as input an object or situation described by a set of properties, and output a 'yes/no' decision.

Entropy: Entropy is a measure of the degree of randomness or uncertainty in the dataset. It measures the randomness based on the distribution of class labels in the dataset.

$$\text{Entropy} = -\frac{m}{m+n} \log_2 \left(\frac{m}{m+n} \right) + -\frac{n}{m+n} \log_2 \left(\frac{n}{m+n} \right)$$



 $m \Rightarrow$ no. of Red balls
 $n \Rightarrow$ no. of Blue balls

Entropy for Bucket 1: $\left[\frac{4}{4+0} \right] \quad -\frac{4}{4+0} \log_2 \left(\frac{4}{4+0} \right) + -\frac{0}{0+4} \log_2 \left(\frac{0}{4+0} \right) = 0 + 0 = 0$

Entropy for Bucket 2: $\left[\frac{3}{3+1} \right] \quad -\frac{3}{3+1} \log_2 \left(\frac{3}{3+1} \right) + -\frac{1}{3+1} \log_2 \left(\frac{1}{3+1} \right) = 0.81125$

Entropy for Bucket 3: $\left[\frac{2}{2+2} \right] \quad -\frac{2}{2+2} \log_2 \left(\frac{2}{2+2} \right) + -\frac{2}{2+2} \log_2 \left(\frac{2}{2+2} \right) = \frac{1}{2} + \frac{1}{2} = 1$

Information Gain: Information Gain measures the reduction in Entropy or variance

that results from splitting a dataset based on a specific property. It is used in decision tree algorithms to determine the usefulness of a feature by partitioning the dataset into more homogeneous subsets with respect to the class labels or target variable.

$$\text{Information Gain}(H, A) = H - \sum \frac{|H_v|}{|H|} H_v$$

A is the specific attribute or class label
 $|H|$ is the entropy of dataset sample S
 $|H_v|$ is the number of instances in the subset S that have the value v for attribute A

For the Given Dataset form a decision Tree

Day	Outlook	Temp	Humidity	Wind	Play Tennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Play Tennis	
Yes	No
9	5

Outlook

$$\Rightarrow \text{Entropy}(H) = -\frac{9}{14} \log_2 \frac{9}{14} - \frac{5}{14} \log_2 \frac{5}{14} = 0.94$$

Sunny		Overcast		Rain	
Yes	No	Yes	No	Yes	No
2	3	4	0	3	2

$$\text{Entropy}(H_{\text{sunny}}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{3}{5} \log_2 \frac{3}{5} = 0.971$$

$$\text{Entropy}(H_{\text{overcast}}) = -\frac{4}{4} \log_2 \frac{4}{4} - \frac{0}{4} \log_2 \frac{0}{4} = 0$$

$$\text{Entropy}(H_{\text{rain}}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} = 0.971$$

$$\text{Gain}(H, H_v) = \text{Entropy}(H) - \sum_{v \in \{\text{sunny, overcast, rain}\}} \frac{|H_v|}{|H|} \text{Entropy}(H_v)$$

$$= 0.971 - \frac{5}{14} \text{Entropy}(H_{\text{sunny}}) - \frac{4}{14} \text{Entropy}(H_{\text{overcast}}) - \frac{5}{14} \text{Entropy}(H_{\text{rain}})$$

$$= 0.971 - \frac{5}{14}(0.971) - \frac{4}{14}(0) - \frac{5}{14}(0.971) = 0.2464$$

$$\Rightarrow 0.2464$$

$$\text{Inf. Gain}(H, \text{outlook}) = 0.2464$$

Attribute: Temperature

Hot		Mild		Cool	
Yes	No	Yes	No	Yes	No
2	2	4	2	3	1

$$\text{Entropy}(H_{\text{hot}}) = 1$$

$$\text{Entropy}(H_{\text{mild}}) = -\frac{4}{6} \log_2 \frac{4}{6} - \frac{2}{6} \log_2 \frac{2}{6} = 0.9183$$

$$\text{Entropy}(H_{\text{cool}}) = -\frac{3}{4} \log_2 \frac{3}{4} - \frac{1}{4} \log_2 \frac{1}{4} = 0.8113$$

$$\text{Inf. Gain}(H, \text{Temp}) = \text{Entropy}(H) - \frac{4}{14} \text{Entropy}(H_{\text{hot}}) - \frac{6}{14} \text{Entropy}(H_{\text{mild}}) - \frac{4}{14} \text{Entropy}(H_{\text{cool}})$$

$$\text{Inf. Gain}(H, \text{Temp}) = 0.0289$$

Attribute: Humidity = (High, Normal)

High		Normal	
Yes	No	Yes	No
3	4	6	1

$$\text{Gain}(H, \text{Humidity}) = 0.1516$$

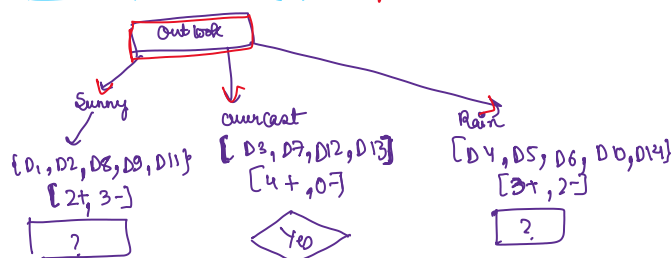
Attribute: Wind

Strong		Weak	
Yes	No	Yes	No
3	3	6	2

$$\text{Gain}(H, \text{wind}) = 0.0478$$

Attribute	I-Gain
Gain(H, outlook)	0.2464
Gain(H, Temp)	0.0289
Gain(H, Humidity)	0.1516
Gain(H, wind)	0.0478

True Root: Outlook (highest Information Gain)



Day	Temp	Humidity	Play Tennis
D1	Hot	High	No

Play Tennis	
Yes	No
3	2

$$\text{Entropy}(H) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5}$$

D ₂	Hot	High	No
D ₈	Mild	High	No
D ₉	Cool	Normal	Yes
D ₁₁	Mild	Normal	Yes

$$1 \text{ sunny } 5 \text{ (2.5)} \quad 5 \text{ '0' (5)} \\ = 0.97$$

Temp					
Hot		Mild		Cool	
Yes	No	Yes	No	Yes	No
0	2	1	1	1	0

$$\text{Entropy}(H_{\text{Hot}}) = 0$$

$$\text{Entropy}(H_{\text{Mild}}) = 1$$

$$\text{Entropy}(H_{\text{Cool}}) = 0$$

$$\text{Gain} = \text{Entropy}(H) - \sum_{v \in \{\text{Hot, Mild, Cool}\}} \frac{|H_v|}{|H|} \text{Entropy}(H_v) \\ = 0.97 - \frac{2}{5} \cdot 0.0 - \frac{2}{5} \cdot 1 - \frac{1}{5} \cdot 0.0 = 0.570$$

$$\approx 0.570$$

Attribute: Humidity

Value (Humidity) = High, Normal

$$\text{Entropy}(H_{\text{High}}) = 0$$

$$\text{Entropy}(H_{\text{Normal}}) = 0$$

Humidity			
High		Normal	
Yes	No	Yes	No
0	3	2	0

$$\text{Gain}(H, \text{Humidity}) = \text{Entropy}(H) - \sum_{v \in \{\text{High, Normal}\}} \frac{|H_v|}{|H|} \text{Entropy}(H_v) \\ = 0.97 - \frac{3}{5}(0) - \frac{2}{5}(0) \\ = 0.97$$

Attribute: Wind

Value (Wind) = Strong, Weak

Wind			
Strong		Weak	
Yes	No	Yes	No
1	1	1	2

$$\text{Entropy}(H_{\text{Strong}}) = 1.0$$

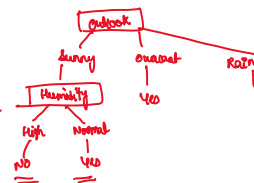
$$\text{Entropy}(H_{\text{Weak}}) = -\frac{1}{3} \log_2 \left(\frac{1}{3} \right) - \frac{2}{3} \log_2 \left(\frac{2}{3} \right) = 0.9183$$

$$\text{Gain}(H_{\text{sunny}}, \text{Wind}) = 0.97 - \frac{2}{5} \cdot 1.0 - \frac{3}{5} \cdot 0.9183 = 0.0192$$

$$\text{Gain}(H_{\text{sunny}}, \text{Temp}) \approx 0.570$$

$$\text{Gain}(H_{\text{sunny}}, \text{Humidity}) \approx 0.97$$

$$\text{Gain}(H_{\text{sunny}}, \text{Wind}) \approx 0.0192$$



Day	Temp	Humidity	Wind	Play-Tennis
D ₄	Mild	High	Weak	Yes
D ₅	Cool	Normal	Weak	Yes
D ₆	Cool	Normal	Strong	No
D ₁₀	Mild	Normal	Weak	Yes
D ₁₄	Mild	High	Strong	No

Rain	
Yes	No
3	2

$$\text{Entropy}(H_{\text{Rain}}) = -\frac{3}{5} \log_2 \frac{3}{5} - \frac{2}{5} \log_2 \frac{2}{5} \approx 0.97$$

$$\text{Entropy}(H_{\text{Hot}}) = 0.0$$

$$\text{Entropy}(H_{\text{Mild}}) = -\frac{2}{5} \log_2 \frac{2}{5} - \frac{1}{3} \log_2 \frac{1}{3} \approx 0.9183$$

$$\text{Entropy}(H_{\text{Cool}}) = 1.0$$

$$\text{Gain}(H_{\text{rain}}, \text{Temp}) = 0.97 - \frac{3}{5} \cdot 0.0 - \frac{2}{5} \cdot 0.9183 - \frac{1}{5} \cdot 1.0 = 0.0192$$

Attribute: Humidity			
High		Normal	
Yes	No	Yes	No
1	1	2	1

$$\text{Entropy}(H_{\text{High}}) = 1.0$$

$$\text{Entropy}(H_{\text{Normal}}) = -\frac{2}{3} \log_2 \frac{2}{3} - \frac{1}{3} \log_2 \frac{1}{3} \approx 0.9183$$

$$\text{Gain}(H_{\text{rain}}, \text{Humidity}) = \text{Entropy}(H) - \sum_{v \in \{\text{High}, \text{Normal}\}} \frac{|H_v|}{|H|} \text{Entropy}(H_v)$$

$$= 0.97 - \frac{2}{5} \cdot 1.0 - \frac{3}{5} \cdot 0.918 = 0.192$$

Attribute wind

Weak		Strong	
Yes	No	Yes	No
0	2	3	0

$$\text{Entropy}(H_{\text{Weak}}) = 0$$

$$\text{Entropy}(H_{\text{Strong}}) = 0$$

$$\text{Inf. Gain}(H_{\text{rain}}, \text{Wind}) = \text{Entropy}(H_{\text{rain}}) - \sum_{v \in \{\text{Strong}, \text{Weak}\}} \frac{|H_v|}{|H|} \text{Entropy}(H_v)$$

$$\text{Inf. Gain}(H_{\text{rain}}, \text{Wind}) = 0.97 - 0 = 0 = 0.97$$

