

- > The biggest challenge of decision tree algorithm is to find out which feature to split upon.
- > i.e. Identifying the feature, i.e. the data should be split, and the split should have examples / data items belonging to a single class.
- > Then only partitions are considered to be pure.

So, if we take the same example dataset first we need to understand, how many attributes?

Are they binary? Multi-valued?

So, among the  $n$  attributes, which attribute would be my root node? which would be the next level node & so on?

- > To answer this, we need to find out information gain of every attribute in the problem.
- > Once we calculate information gain, we can conclude, which attribute has the highest importance.

Now, to calculate information gain, you need to first calculate the entropy.

## Entropy

Entropy is a measure of impurity of a collection of examples. (uncertainty in a given dataset)

Entropy depends on the distribution of random variable  $P$ .

- $S$  is a collection of training examples.
- $p+$  the proportion of 'true' examples.
- $p-$  the proportion of 'false' examples.

$$\text{Entropy}(S) = -p+ \log_2 p+ - p- \log_2 p-$$

Eg:- There are 3 cases

$$\text{① If Entropy}([14+, 0-]) = -14/14 \log_2(14/14) - 0 \log_2(0) = 0$$

If you have all positive examples, like above where we have 14 true example and 0 false examples or vice versa, you don't need to apply formula as Entropy will always be 0.

(Combination of the 2-ue examples)

$$(2) \text{ Entropy } (9+, 5-) = -9/14 \log_2(9/14) - 5/14 \log_2(5/14) \\ = 0.94$$

have to apply formula.

(Equal two 4-ue examples)

$$(3) \text{ Entropy } (7+, 7-) = -7/14 \log_2(7/14) - 7/14 \log_2(7/14) \\ = 1/2 + 1/2 = 1$$

In this case entropy will always be 1.

Information Gain Information gain is calculated on the basis of decrease in entropy (s) caused by the partitioning dataset of the examples according to the attribute.

- > The attribute with 'minimum' impurity is the important node.
- > Constructing a decision tree is all about finding an attribute that returns highest information gain.

# Decision Tree - ID3 Algorithm Numerical example.

<u>Day</u>	<u>Outlook</u>	<u>Temp</u>	<u>Humidity</u>	<u>Wind</u>	<u>Play Tennis</u>
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	Rainy	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.



Attribute: Outlook

value (outlook): Sunny, Overcast, Rainy

$$S = [9+, 5-]$$

$$\begin{aligned}\text{Entropy}(S) &= -\frac{9}{14} \log_{\frac{1}{2}} \frac{9}{14} - \frac{5}{14} \log_{\frac{1}{2}} \frac{5}{14} \\ &= \underline{\underline{0.94}}\end{aligned}$$

$$\begin{aligned}S_{\text{sunny}} &\leftarrow [2+, 3-] \\ \text{Entropy}(S_{\text{sunny}}) &= -\frac{2}{5} \log_{\frac{1}{2}} \frac{2}{5} - \frac{3}{5} \log_{\frac{1}{2}} \frac{3}{5} \\ &= 0.971\end{aligned}$$

$$\begin{aligned}S_{\text{overcast}} &\leftarrow [4+, 0-] \\ \text{Entropy}(S_{\text{overcast}}) &= -\frac{4}{4} \log_{\frac{1}{2}} \frac{4}{4} - \frac{0}{4} \log_{\frac{1}{2}} \frac{0}{4} = \underline{\underline{0}}\end{aligned}$$

$$S_{\text{Rain}} \leftarrow [3+, 2-]$$

$$\begin{aligned}\text{Entropy}(S_{\text{Rain}}) &= -\frac{3}{5} \log_{\frac{1}{2}} \frac{3}{5} - \frac{2}{5} \log_{\frac{1}{2}} \frac{2}{5} \\ &= \underline{\underline{0.971}}\end{aligned}$$

$$\text{Gain}(S_{\text{outlook}}) = \text{Entropy}(S) - \sum \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$\begin{aligned}&= \text{Entropy}(S) - \frac{5}{14} \text{Entropy}(S_{\text{sunny}}) - \frac{4}{14} \text{Entropy}(S_{\text{overcast}}) \\ &\quad - \frac{5}{14} \text{Entropy}(S_{\text{Rain}})\end{aligned}$$

$$= 0.94 - \frac{5}{14} * 0.971 - \frac{4}{14} * 0 - \frac{5}{14} * 0.971$$

$$= \underline{\underline{0.2464}}$$

$$\text{Gain}(S, \text{outlook}) = \underline{\underline{0.2464}}$$

Similarly we need to calculate Information gain for other attributes.

Attributes: (temp) → Hot, mild, cool  
 Value (temp) ←

$$\text{Entropy}(S) = 0.94$$

$$S_{\text{Hot}} [2+, 2-] \quad \text{Entropy}(S_{\text{Hot}}) = 1.0$$

$$S_{\text{mild}} [4+, 2-] \quad \text{Entropy}(S_{\text{mild}}) = 0.9183$$

$$-\frac{4}{6} \log_{\frac{1}{2}} \frac{4}{6} - \frac{2}{6} \log_{\frac{1}{2}} \frac{2}{6}$$

$$S_{\text{cool}} [3+, 1-] \quad \text{Entropy}(S_{\text{cool}}) = 0.8113$$

$$-\frac{3}{4} \log_{\frac{1}{2}} \frac{3}{4} - \frac{1}{4} \log_{\frac{1}{2}} \frac{1}{4}$$

$$\text{Gain}(S, \text{temp}) = 0.94 - \frac{4}{14} * 1.0 - \frac{6}{14} * 0.9183 - \frac{4}{14} * 0.8113 = \underline{\underline{0.0289}}$$

$$= \underline{\underline{0.0289}}$$

Attribute: Humidity

value (Humidity) = High, Normal

$$S = [9+, 5-]$$

$$\text{Entropy}(S) = \underline{0.94}$$

$$S_{\text{high}} \leftarrow [3+, 4-] \quad \text{Entropy}(S_{\text{high}}) = 0.9852$$

$$S_{\text{Normal}} \leftarrow [6+, 1-] \quad \text{Entropy}(S_{\text{Normal}}) = 0.5916$$

$$\begin{aligned} \text{Gain}(S, \text{Humidity}) &= \text{Entropy}(S) - \frac{7}{14} * 0.9852 \\ &\quad - \frac{7}{14} * 0.5916 = 0.1516 \end{aligned}$$

0.1516

Attribute: Wind

$$S = [9+, 5-] \quad \text{Entropy}(S) = 0.94$$

$$S_{\text{strong}} \leftarrow [3+, 3-] \quad \text{Entropy}(S_{\text{strong}}) = 1.0$$

$$S_{\text{weak}} \leftarrow [6+, 2-] \quad \text{Entropy}(S_{\text{weak}}) = 0.8113$$

high

$$\begin{aligned}
 \text{Gain}(S, \text{wind}) &= \text{Entropy}(S) - \sum_{v \in \{\text{strong, weak}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v) \\
 &= 0.94 - \frac{6}{14} \times 1.0 - \frac{8}{14} \times 0.8113 \\
 &= 0.0478 //
 \end{aligned}$$

After calculating gain of all attributes we need to check which attribute has the highest information gain.

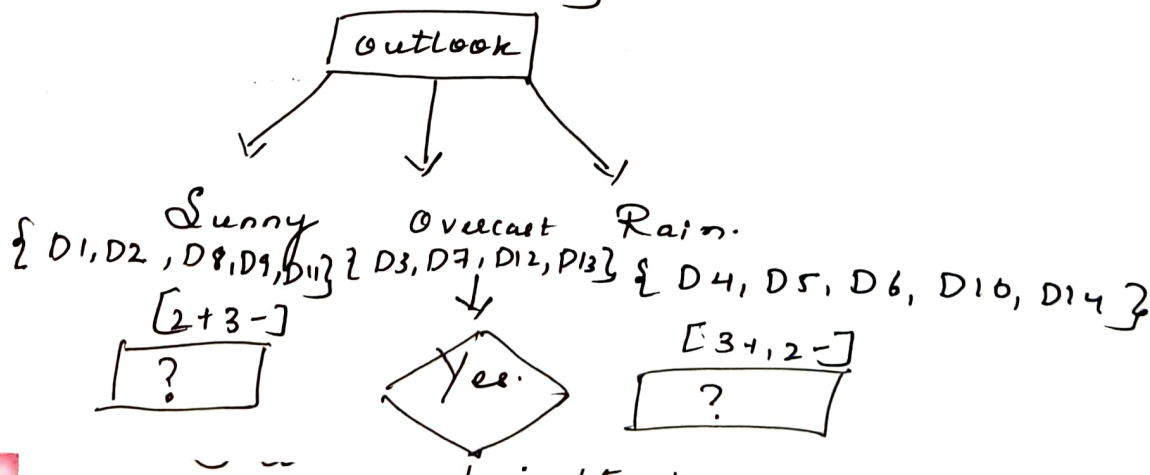
$$\text{Gain}(S, \text{outlook}) = 0.2464 \quad \checkmark$$

$$\text{Gain}(S, \text{temp}) = 0.0289$$

$$\text{Gain}(S, \text{humidity}) = 0.1516$$

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

Outlook is the root node.  
[D1.... D14]





Now, let's consider for sunny

Day	Temp	Humidity	Wind	Play Tennis
D1	Hot	High	Weak	No
D2	Hot	High	Strong	No
D8	Mild	High	Weak	No
D9	Cool	Normal	Weak	Yes
D11	Mild	Normal	Strong	Yes

Attribute: Temp  
 values(Temp) = Hot, mild, cool

$$S_{\text{sunny}} = [2+, 3-] \quad \text{Entropy}(S_{\text{sunny}}) = -\frac{2}{5} \log \frac{2}{5} - \frac{3}{5} \log \frac{3}{5} = 0.97$$

$$S_{\text{hot}} \leftarrow [0+, 2-] \quad \text{Entropy}(S_{\text{hot}}) = 0$$

$$S_{\text{mild}} \leftarrow [1+, 1-] \quad \text{Entropy}(S_{\text{mild}}) = 1.0$$

$$S_{\text{cool}} \leftarrow [1+, 0-] \quad \text{Entropy}(S_{\text{cool}}) = 0.0$$

$$\begin{aligned} \text{Gain}_{\text{temp}}(S_{\text{sunny}}, \text{temp}) &= \text{Entropy}(S_{\text{sunny}}) - \sum_{v \in \{\text{hot, mild, cool}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v) \\ &= 0.97 - \frac{2}{5} * 0 - \frac{2}{5} * 1.0 - \frac{1}{5} * 0 \\ &= 0.570 // \end{aligned}$$

Attribute: Humidity

value: High, Normal

$$S_{\text{sunny}} = 0.97$$

$$S_{\text{high}} \leftarrow [0+3-] \quad \text{Entropy}(S_{\text{high}}) = 0$$

$$S_{\text{normal}} \leftarrow [2+, 0-] \quad \text{Entropy}(S_{\text{normal}}) = 0.0$$

$$\text{Gain}(S_{\text{sunny}}, \text{humidity})$$

$$= \text{Entropy}(S_{\text{sunny}}) - \sum_{v \in \{\text{High, Normal}\}} \frac{|S_v|}{|S|} \text{Entropy}(S_v)$$

$$= 0.97 - \frac{3}{5} * 0 - \frac{2}{5} * 0$$

$$= 0.97 //$$

Attribute: Wind

value: (Strong, Weak)

$$S_{\text{sunny}} = [2+, 3-] = 0.97$$

$$S_{\text{strong}} [1+, 1-] = \text{Entropy}(S_{\text{strong}}) = 1.0$$

$$S_{\text{weak}} [1+, 2-] \quad \text{Entropy}(S_{\text{weak}}) = 0.9182$$

$$\text{Gain (Sunny, wind)} = 0.0192$$

So,  $\text{Gain (Sunny, temp)} = 0.570$

$$\text{Gain (Sunny, Humidity)} = 0.97$$

$$\text{Gain (Sunny, Wind)} = 0.0192$$

Day	temp	humidity	<del>Rain</del> Wind	Play Tennis
D <sub>4</sub>	Mild	High	Weak	Yes
D <sub>5</sub>	Cool	Normal	Weak	Yes
D <sub>6</sub>	Cool	Normal	Strong	No
D <sub>10</sub>	Mild	Normal	Weak	Yes
D <sub>14</sub>	Mild	High	Strong	No

values (temp)

$$S[\text{Rain}] = [3+, 2-] = 0.97$$

$$S_{\text{Hot}} \leftarrow [0, 0] \quad \text{Entropy}[S_{\text{Hot}}] = 0.0$$

$$S_{\text{mild}} \leftarrow [2+, 1-] \quad \text{Entropy}[S_{\text{mild}}] = 0.9183$$

$$S_{\text{cool}} \leftarrow [1+, 1-] \quad \text{Entropy}[S_{\text{cool}}] = 1.0$$

$$\text{Gain}(S_{\text{Rain}}, \text{temp}) = 0.0192$$

$$\text{Gain}(S_{\text{Rain}}, \text{humidity}) = 0.0192$$

$$\text{Gain}(S_{\text{Rain}}, \text{Wind}) = 0.97 \quad \checkmark$$

{D1, D2, D8, D9, D11}

[2+3-]

Humidity

High  
|

{D1, D2, D8}

No

Normal  
|

{D9, D11}

Yes

[D3, D7, D12, D13]

[4+0-]

Yes

[D4, D5, D6, D10, D14]

[3+2-]

?

Wind

Strong  
|

No

{D6, D14}

No

Weak  
|

Yes

{D4, D5, D10}

Yes