

Decision Tree

Decision Tree is a supervised learning method used for classification and regression. It is a tree which helps us by assisting us in decision-making!

Decision tree builds classification or regression models in the form of a tree structure. It breaks down a data set into smaller and smaller subsets and simultaneously decision tree is incrementally developed. The final tree is a tree with decision nodes and leaf nodes. A decision node has two or more branches. Leaf node represents a classification or decision. We cannot do more split on leaf nodes.

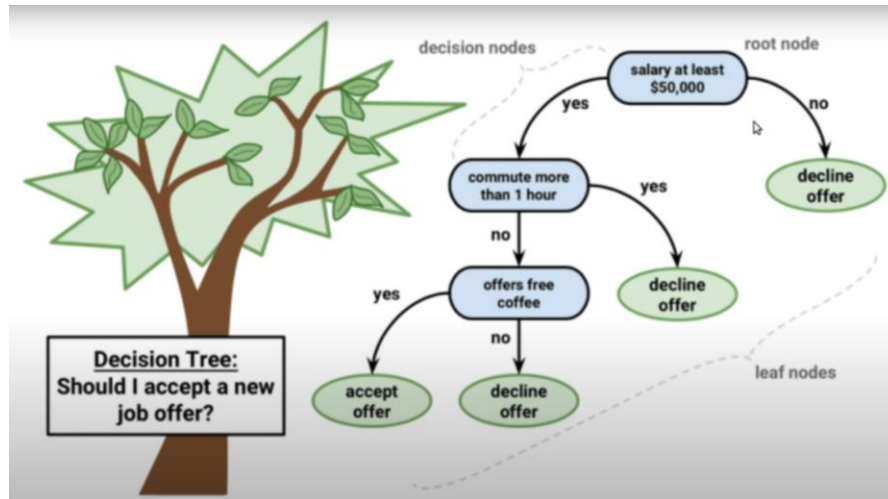
The topmost decision node in a tree which corresponds to the best predictor called root node. Decision trees can handle both categorical and numerical data.

Common terms used with Decision trees:

- **Root Node:** It represents entire population or sample and this further gets divided into two or more homogeneous sets.
- **Splitting:** It is a process of dividing a node into two or more sub-nodes.
- **Decision Node:** When a sub-node splits into further sub-nodes, then it is called decision node.
- **Leaf/ Terminal Node:** Nodes do not split is called Leaf or Terminal node.
- **Pruning:** When we remove sub-nodes of a decision node, this process is called pruning. You can say opposite process of splitting.
- **Branch / Sub-Tree:** A sub section of entire tree is called branch or sub-tree.
- **Parent and Child Node:** A node, which is divided into sub-nodes is called parent node of sub-nodes whereas sub-nodes are the child of parent node.

How does Decision Tree works ?

Decision tree is a type of supervised learning algorithm (having a pre-defined target variable) that is mostly used in classification problems. It works for both categorical and continuous input and output variables. In this technique, we split the population or sample into two or more homogeneous sets (or sub-populations) based on most significant splitter / differentiator in input variables.



There are two algorithms to solve a problem using Decision Tree.

1. CART: Gini Index
2. ID3: Entropy, Gain

Let us look at the example below:

S. No.	Outlook	Temperature	Humidity	Windy	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rainy	Mild	High	Weak	Yes
5	Rainy	Cool	Normal	Weak	Yes
6	Rainy	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rainy	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rainy	Mild	High	Strong	No

Now, we need to make a decision tree that predicts whether Tennis will be played or not based on the above data.

Here the class attribute is: “Play Tennis”

ID3 (Iterative Dichotomiser 3) algorithm has two important concepts.

- i. Entropy
- ii. Information Gain

Entropy states the uncertainty in the dataset, i.e. the number positive and number of negative samples. E.g. if we have equal number of positive and equal number of negative samples, then the Entropy will be 1. If we have only positive samples, or only negatives, then Entropy will be 0.

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

Information Gain is the difference in Entropy before and after splitting dataset on attribute.

$$Gain = Entropy(S) - I(Attribute)$$

At first, we need to calculate the Entropy for entire dataset (shown above) then we will calculate information gain for each attribute separately. See the table below.

S. No.	Outlook	Temperature	Humidity	Windy	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rainy	Mild	High	Weak	Yes
5	Rainy	Cool	Normal	Weak	Yes
6	Rainy	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rainy	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rainy	Mild	High	Strong	No

After calculating separately, we will sum the individual information gain. This is known as Average Information ($I_{\text{attribute}}$)

$$I(Attribute) = \sum \frac{p_i + n_i}{p+n} Entropy(A)$$

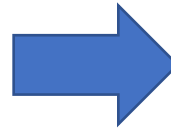
Pseudocode:

1. Computer Entropy for entire dataset, Entropy (S)
2. For Every Attribute:
 - a. Calculate Entropy for all other values (Entropy (A))
 - b. Take Average Information Entropy for the current Attribute
 - c. Calculate Gain for the current Attribute
3. Pick the Highest Gain Attribute
4. Repeat until we get the desired Tree constructed.

Let us solve the problem:

- **Step 1: We will find out number of positive samples and negative samples from the dataset.**

S. No.	Outlook	Temperature	Humidity	Windy	PlayTennis
1	Sunny	Hot	High	Weak	No
2	Sunny	Hot	High	Strong	No
3	Overcast	Hot	High	Weak	Yes
4	Rainy	Mild	High	Weak	Yes
5	Rainy	Cool	Normal	Weak	Yes
6	Rainy	Cool	Normal	Strong	No
7	Overcast	Cool	Normal	Strong	Yes
8	Sunny	Mild	High	Weak	No
9	Sunny	Cool	Normal	Weak	Yes
10	Rainy	Mild	Normal	Weak	Yes
11	Sunny	Mild	Normal	Strong	Yes
12	Overcast	Mild	High	Strong	Yes
13	Overcast	Hot	Normal	Weak	Yes
14	Rainy	Mild	High	Strong	No



Yes= 9
No = 5
Total = 14

- **Step 2: Calculating the Entropy:**

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

$$Entropy(S) = \frac{-9}{9+5} \log_2\left(\frac{9}{9+5}\right) - \frac{5}{9+5} \log_2\left(\frac{5}{9+5}\right)$$

$$Entropy(S) = \frac{-9}{14} \log_2\left(\frac{9}{14}\right) - \frac{5}{14} \log_2\left(\frac{5}{14}\right) = 0.940$$

- For each Attribute: (let say **Outlook**)
 - Calculate Entropy for each Values, i.e for 'Sunny', 'Rainy', 'Overcast'

Outlook	PlayTennis
Sunny	No
Sunny	No
Sunny	No
Sunny	Yes
Sunny	Yes

Outlook	PlayTennis
Rainy	Yes
Rainy	Yes
Rainy	No
Rainy	Yes
Rainy	No

Outlook	PlayTennis
Overcast	Yes
Overcast	Yes
Overcast	Yes
Overcast	Yes

Outlook	p	n	Entropy
Sunny	2	3	0.971
Rainy	3	2	0.971
Overcast	4	0	0

Calculate **Entropy(Outlook='Value')**:

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

$$E(\text{Outlook=sunny}) = -\frac{2}{5} \log \left(\frac{2}{5} \right) - \frac{3}{5} \log \left(\frac{3}{5} \right) = 0.971$$

$$E(\text{Outlook=overcast}) = -1 \log(1) - 0 \log(0) = 0$$

$$E(\text{Outlook=rainy}) = -\frac{3}{5} \log \left(\frac{3}{5} \right) - \frac{2}{5} \log \left(\frac{2}{5} \right) = 0.971$$

- Calculate **Average Information Entropy**:

$$I(Outlook) = \frac{p_{sunny} + n_{sunny}}{p + n} Entropy(Outlook = Sunny) +$$

$$\frac{p_{rainy} + n_{rainy}}{p + n} Entropy(Outlook = Rainy) +$$

$$\frac{p_{Overcast} + n_{Overcast}}{p + n} Entropy(Outlook = Overcast)$$

$$I(Outlook) = \frac{3 + 2}{9 + 5} * 0.971 + \frac{2 + 3}{9 + 5} * 0.971 + \frac{4 + 0}{9 + 5} * 0 = 0.693$$

- Calculate **Gain**: attribute is Outlook

$$Gain = Entropy(S) - I(Attribute)$$

$$Entropy(S) = 0.940$$

$$Gain(Outlook) = 0.940 - 0.693 = 0.247$$

- For each Attribute: (let say **Temperature**)
 - Calculate Entropy for each Temp, i.e for 'Hot', 'Mild' and 'Cool'

Temperature	PlayTennis
Hot	No
Hot	No
Hot	Yes
Hot	Yes

Temperature	PlayTennis
Mild	Yes
Mild	No
Mild	Yes
Mild	Yes
Mild	Yes
Mild	No

Temperature	PlayTennis
Cool	Yes
Cool	No
Cool	Yes
Cool	Yes

Temperature	p	n	Entropy
Hot	2	2	1
Mild	4	2	0.918
Cool	3	1	0.811

- Calculate **Average Information Entropy**:

$$I(\text{Temperature}) = \frac{p_{\text{hot}} + n_{\text{hot}}}{p + n} \text{Entropy}(\text{Temperature} = \text{Hot}) +$$

$$\frac{p_{\text{mild}} + n_{\text{mild}}}{p + n} \text{Entropy}(\text{Temperature} = \text{Mild}) +$$

$$\frac{p_{\text{Cool}} + n_{\text{Cool}}}{p + n} \text{Entropy}(\text{Temperature} = \text{Cool})$$

$$I(\text{Temperature}) = \frac{2 + 2}{9 + 5} * 1 + \frac{4 + 2}{9 + 5} * 0.918 + \frac{3 + 1}{9 + 5} * 0.811 \Rightarrow 0.911$$

- Calculate **Gain**: attribute is Temperature

$$Gain = Entropy(S) - I(Attribute)$$

$$Entropy(S) = 0.940$$

$$Gain(Temperature) = 0.940 - 0.911 = 0.029$$

For each Attribute: (let say **Humidity**)

- Calculate Entropy for each Humidity, i.e for 'High', 'Normal'

Humidity	PlayTennis
Normal	Yes
Normal	No
Normal	Yes
Normal	Yes
Normal	Yes
Normal	Yes
Normal	Yes

→

Humidity	PlayTennis
High	No
High	No
High	Yes
High	Yes
High	No
High	Yes
High	No

Humidity	p	n	Entropy
High	3	4	0.985
Normal	6	1	0.591

- Calculate **Average Information Entropy**:

$$I(Humidity) = \frac{p_{High} + n_{High}}{p + n} Entropy(Humidity = High) + \frac{p_{Normal} + n_{Normal}}{p + n} Entropy(Humidity = Normal)$$

$$I(Humidity) = \frac{3 + 4}{9 + 5} * 0.985 + \frac{6 + 1}{9 + 5} * 0.591 => 0.788$$

- Calculate **Gain**: attribute is Humidity

$$Gain = Entropy(S) - I(Attribute)$$

$$Entropy(S) = 0.940$$

$$Gain(Humidity) = 0.940 - 0.788 = 0.152$$

- For each Attribute: (let say **Windy**)
 - Calculate Entropy for each Windy, i.e for 'Strong' and 'Weak'

Windy	PlayTennis
Weak	No
Weak	Yes
Weak	Yes
Weak	Yes
Weak	No
Weak	Yes
Weak	Yes
Weak	Yes

Windy	PlayTennis
Strong	No
Strong	No
Strong	Yes
Strong	Yes
Strong	Yes
Strong	No

Windy	p	n	Entropy
Strong	3	3	1
Weak	6	2	0.811

- Calculate **Average Information Entropy**:

$$I(Windy) = \frac{p_{Strong} + n_{Strong}}{p + n} Entropy(Windy = Strong) + \frac{p_{Weak} + n_{Weak}}{p + n} Entropy(Windy = Weak)$$

$$I(Windy) = \frac{3 + 3}{9 + 5} * 1 + \frac{6 + 2}{9 + 5} * 0.811 \Rightarrow 0.892$$

$$Gain(Windy) = 0.940 - 0.892 = 0.048$$

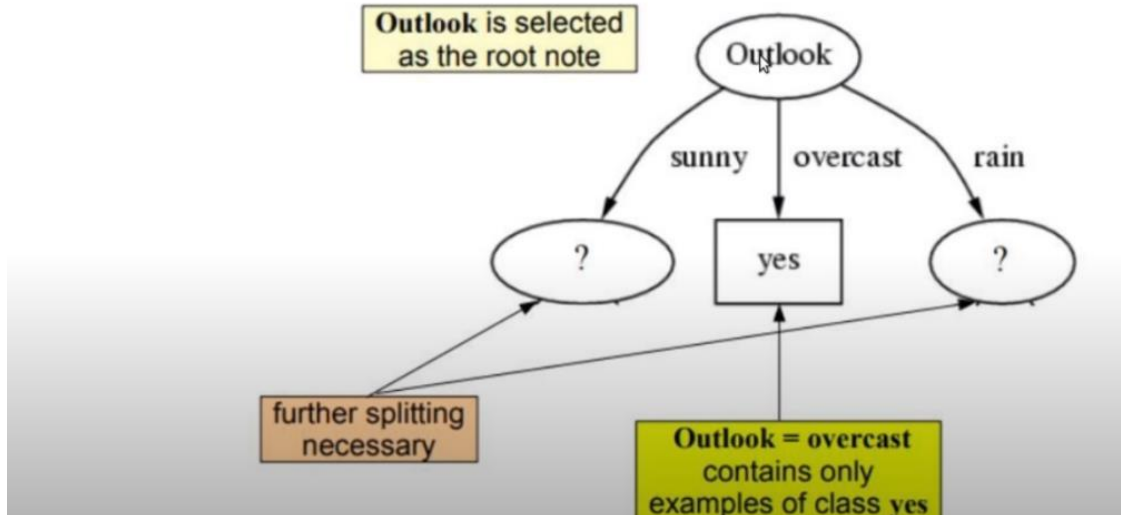
- PICK THE HIGHEST GAIN ATTRIBUTE.

Attributes	Gain
Outlook	0.247
Temperature	0.029
Humidity	0.152
Windy	0.048

ROOT NODE:
OUTLOOK

Outlook	Temperature	Humidity	Windy	PlayTennis
Overcast	Hot	High	Weak	Yes
Overcast	Cool	Normal	Strong	Yes
Overcast	Mild	High	Strong	Yes
Overcast	Hot	Normal	Weak	Yes

Outlook is selected
as the root node



Outlook	Temperature	Humidity	Windy	PlayTennis
Sunny	Hot	High	Weak	No
Sunny	Hot	High	Strong	No
Sunny	Mild	High	Weak	No
Sunny	Cool	Normal	Weak	Yes
Sunny	Mild	Normal	Strong	Yes

$$P = \frac{2}{5} \quad N = \frac{3}{5}$$

• ENTROPY:

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

$$Entropy(S_{sunny}) = \frac{-2}{2+3} \log_2\left(\frac{2}{2+3}\right) - \frac{3}{2+3} \log_2\left(\frac{3}{2+3}\right)$$

$$\Rightarrow 0.971$$

• For each Attribute: (let say **Humidity**):

- Calculate Entropy for each Humidity, i.e for 'High' and 'Normal'

Outlook	Humidity	PlayTennis
Sunny	High	No
Sunny	High	No
Sunny	High	No
Sunny	Normal	Yes
Sunny	Normal	Yes

Humidity	p	n	Entropy
high	0	3	0
normal	2	0	0

- Calculate **Average Information Entropy**:

$$I(\text{Humidity}) = 0$$

- Calculate **Gain**:

$$\text{Gain} = 0.971$$

- For each Attribute: (let say **Windy**):
 - Calculate Entropy for each Windy, i.e for 'Strong' and 'Weak'

Outlook	Windy	PlayTennis
Sunny	Strong	No
Sunny	Strong	Yes
Sunny	Weak	No
Sunny	Weak	No
Sunny	Weak	Yes

Windy	p	n	Entropy
Strong	1	1	1
Weak	1	2	0.918

- Calculate **Average Information Entropy**: $I(\text{Windy}) = 0.951$
- Calculate **Gain**: $\text{Gain} = 0.020$

- For each Attribute: (let say **Temperature**):
 - Calculate Entropy for each Windy, i.e for 'Cool', 'Hot' and 'Mild'

Outlook	Temperature	PlayTennis
Sunny	Cool	Yes
Sunny	Hot	No
Sunny	Hot	No
Sunny	Mild	No
Sunny	Mild	Yes

Temperature	p	n	Entropy
Cool	1	0	0
Hot	0	2	0
Mild	1	1	1

- Calculate **Average Information Entropy**: $I(\text{Temp}) = 0.4$
- Calculate **Gain**: $\text{Gain} = 0.571$

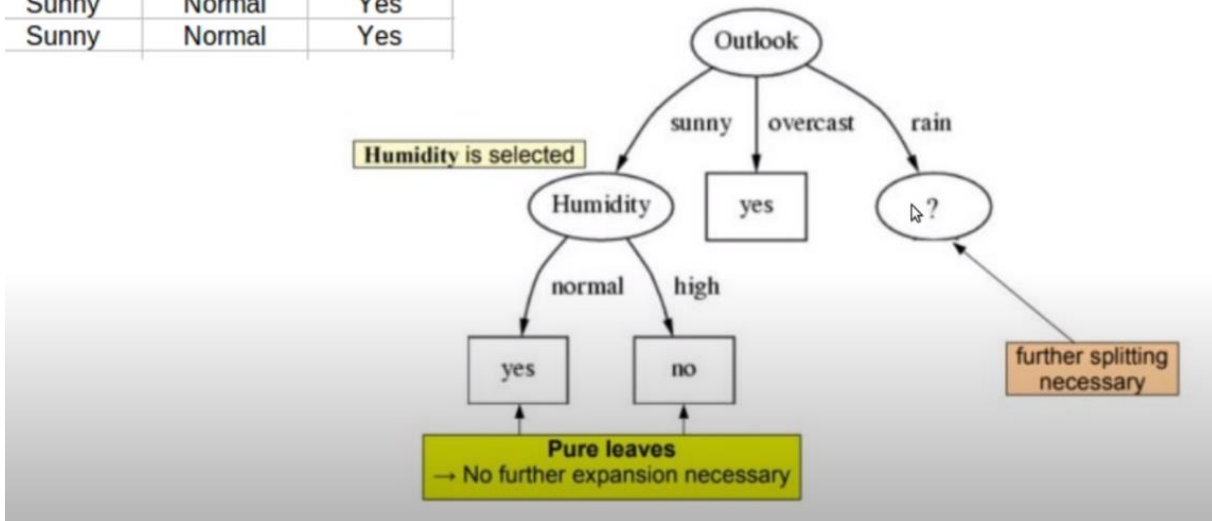
- PICK THE HIGHEST GAIN ATTRIBUTE.

Attributes	Gain
Temperature	0.571
Humidity	0.971
Windy	0.02

NEXT NODE IN SUNNY:

HUMIDITY

Outlook	Humidity	PlayTennis
Sunny	High	No
Sunny	High	No
Sunny	High	No
Sunny	Normal	Yes
Sunny	Normal	Yes



Repeat the above steps to find the node after outlook (shown in ? mark above) and you will find “windy”. The final tree will look like:

