# DECISION TREE ALGORITHM

# Introduction to Decision Tree algorithm

- A Decision Tree algorithm is one of the most popular machine learning algorithms. It uses a tree like structure and their possible combinations to solve a particular problem. It belongs to the class of supervised learning algorithms where it can be used for both classification and regression purposes.

- A decision tree is a structure that includes a root node, branches, and leaf nodes. Each internal node denotes a test on an attribute, each branch denotes the outcome of a test, and each leaf node holds a class label. The topmost node in the tree is the root node.

- We make some assumptions while implementing the Decision-Tree algorithm. These are listed below:-

  1) At the beginning, the whole training set is considered as the root.

  2) Feature values need to be categorical. If the values are continuous then they are discretized prior to building the model.

  3) Records are distributed recursively on the basis of attribute values.

  4) Order to placing attributes as root or internal node of the tree is done by using some statistical approach.

# Decision Tree algorithm terminology

In a Decision Tree algorithm, there is a tree like structure in which each internal node represents a test on an attribute, each branch represents the outcome of the test, and each leaf node represents a class label. The paths from the root node to leaf node represent classification rules.

We can see that there is some terminology involved in Decision Tree algorithm. The terms involved in Decision Tree algorithm are as follows:-

1) Root Node

  It represents the entire population or sample. This further gets divided into two or more homogeneous sets.

2) Splitting

   It is a process of dividing a node into two or more sub-nodes.

3) Decision Node

   When a sub-node splits into further sub-nodes, then it is called a decision node.

4) Leaf/Terminal Node

   Nodes that do not split are called Leaf or Terminal nodes.

5) Pruning

   When we remove sub-nodes of a decision node, this process is called pruning. It is the opposite process of splitting.
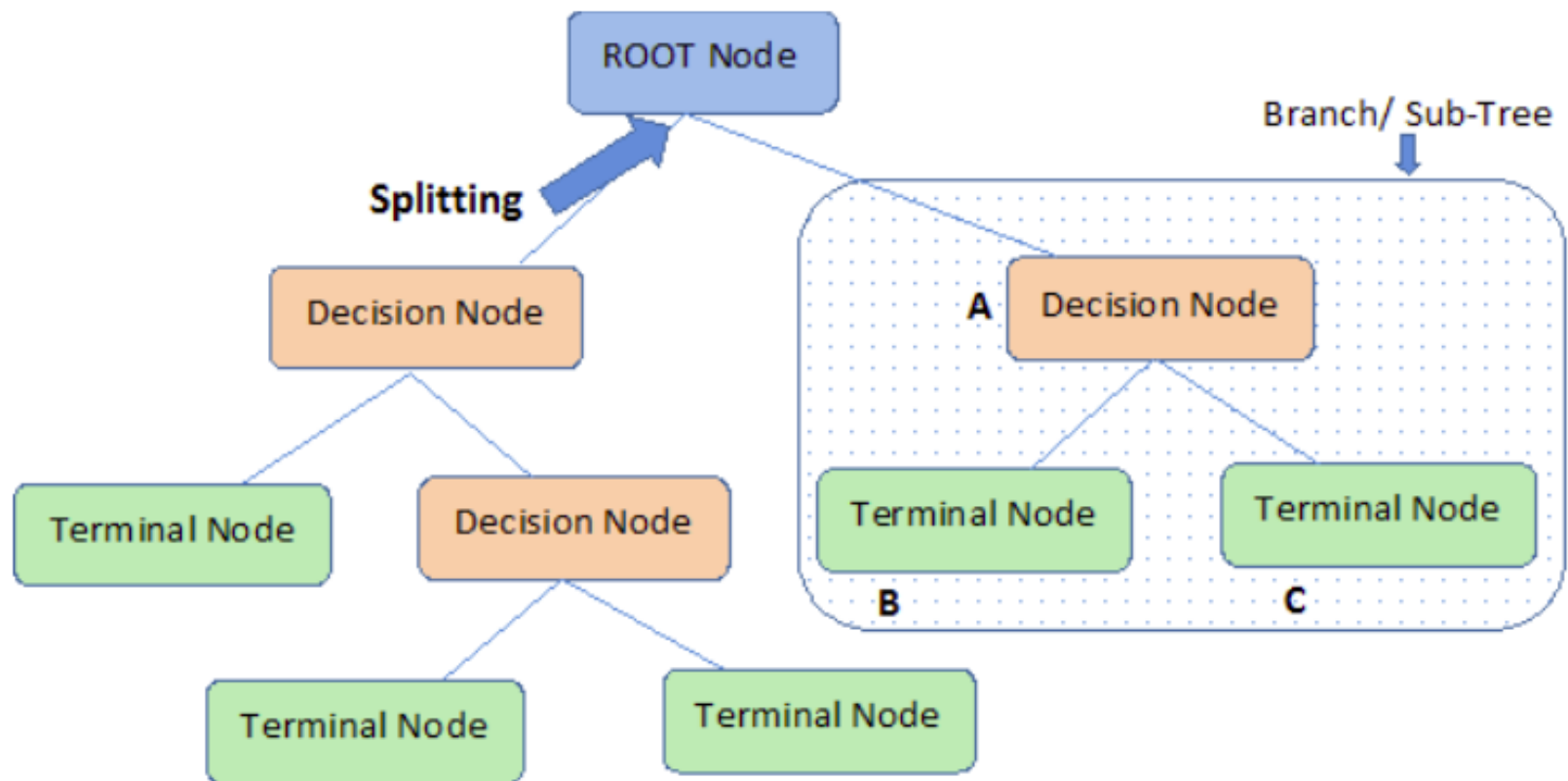
6) Branch/Sub-Tree

   A sub-section of an entire tree is called a branch or sub-tree.

7) Parent and Child Node

   A node, which is divided into sub-nodes is called the parent node of sub-nodes where sub-nodes are the children of a parent        node.

## Decision-Tree terminology

# Decision Tree algorithm intuition

The Decision-Tree algorithm is one of the most frequently and widely used supervised machine learning algorithms that can be used for both classification and regression tasks. The intuition behind the Decision-Tree algorithm is very simple to understand.

The Decision Tree algorithm intuition is as follows:-

For each attribute in the dataset, the Decision-Tree algorithm forms a node. The most important attribute is placed at the root node.

For evaluating the task in hand, we start at the root node and we work our way down the tree by following the corresponding node that meets our condition or decision.

This process continues until a leaf node is reached. It contains the prediction or the outcome of the Decision Tree.

# Attribute selection measures

The primary challenge in the Decision Tree implementation is to identify the attributes which we consider as the root node and each level. This process is known as the attributes selection. There are different attributes selection measure to identify the attribute which can be considered as the root node at each level.

There are 2 popular attribute selection measures. They are as follows:-

Information gain

Gini index

While using Information gain as a criterion, we assume attributes to be categorical and for Gini index attributes are assumed to be continuous. These attribute selection measures are described below.

# Information gain

By using information gain as a criterion, we try to estimate the information contained by each attribute. To understand the concept of Information Gain, we need to know another concept called Entropy.

Entropy

Entropy measures the impurity in the given dataset. In Physics and Mathematics, entropy is referred to as the randomness or uncertainty of a random variable X. In information theory, it refers to the impurity in a group of examples. Information gain is the decrease in entropy. Information gain computes the difference between entropy before split and average entropy after split of the dataset based on given attribute values.

Entropy is represented by the following formula:-

$$Entropy = \sum_{i=1}^{C} -p_i * \log_2(p_i)$$

Here, **c** is the number of classes and **pi** is the probability associated with the ith class.

Here, again c is the number of classes and pi is the probability associated with the ith class.

**Gini index** says, if we randomly select two items from a population, they must be of the same class and probability for this is 1 if the population is pure.

It works with the categorical target variable **"Success" or "Failure"**. It performs only binary splits. The higher the value of Gini, higher the homogeneity. CART (Classification and Regression Tree) uses the Gini method to create binary splits.

**Steps to Calculate Gini for a split:-**

Calculate Gini for sub-nodes, using formula sum of the square of probability for success and failure (p^2+q^2).

Calculate Gini for split using weighted Gini score of each node of that split.

In case of a discrete-valued attribute, the subset that gives the minimum gini index for that chosen is selected as a splitting attribute. In the case of continuous-valued attributes, the strategy is to select each pair of adjacent values as a possible split-point and point with smaller gini index chosen as the splitting point. The attribute with minimum Gini index is chosen as the splitting attribute.

# Overfitting in Decision Tree algorithm

Overfitting is a practical problem while building a Decision-Tree model. The problem of overfitting is considered when the algorithm continues to go deeper and deeper to reduce the training-set error but results with an increased test-set error. So, accuracy of prediction for our model goes down. It generally happens when we build many branches due to outliers and irregularities in data.

Two approaches which can be used to avoid overfitting are as follows:-

**Pre-Pruning**

In pre-pruning, we stop the tree construction a bit early. We prefer not to split a node if its goodness measure is below a threshold value. But it is difficult to choose an appropriate stopping point.

**Post-Pruning**

In post-pruning, we go deeper and deeper in the tree to build a complete tree. If the tree shows the overfitting problem then pruning is done as a post-pruning step. We use the cross-validation data to check the effect of our pruning. Using cross-validation data, we test whether expanding a node will result in improve or not. If it shows an improvement, then we can continue by expanding that node. But if it shows a reduction in accuracy then it should not be expanded. So, the node should be converted to a leaf node.