

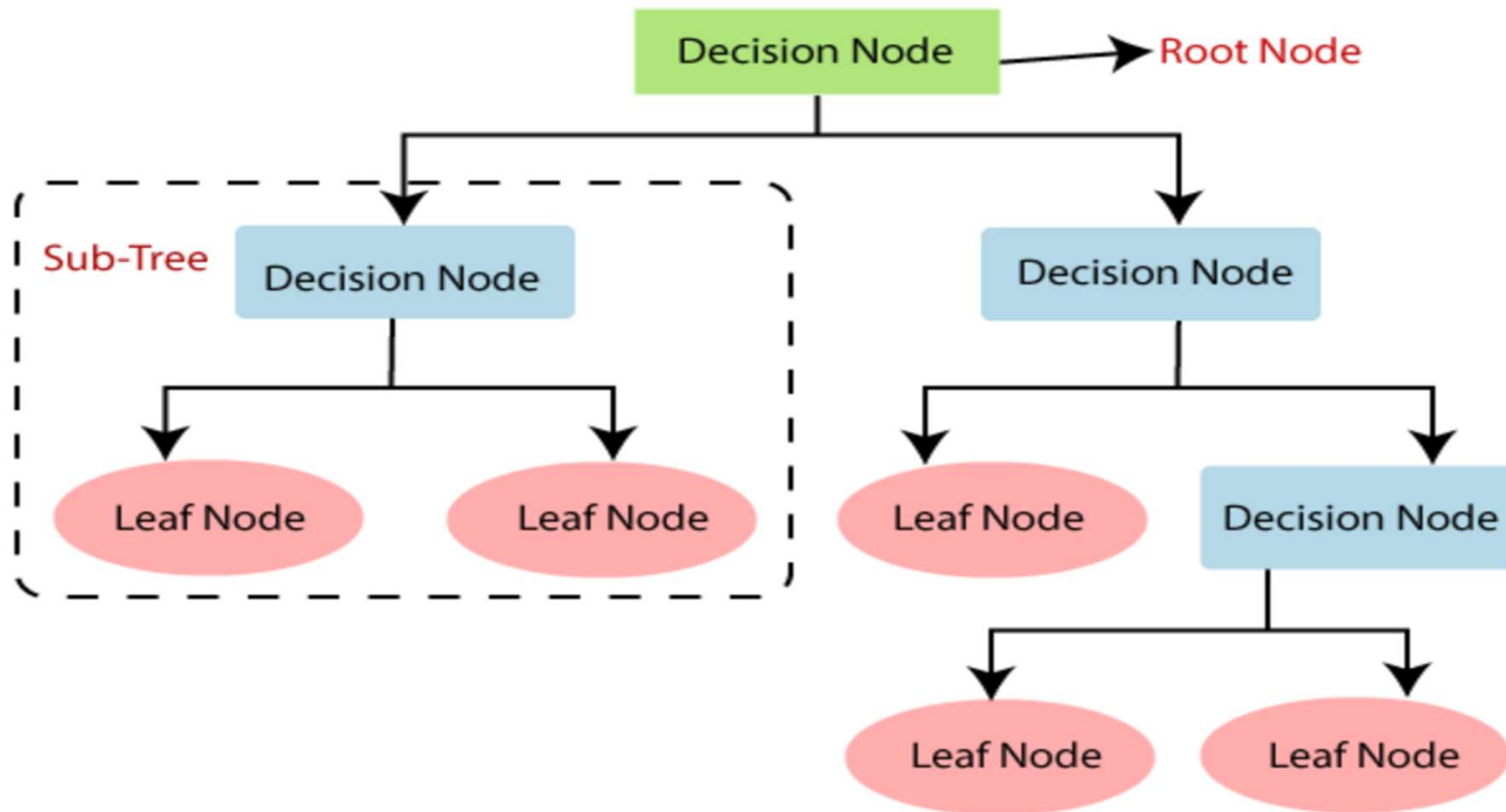
Decision Tree Machine Learning Model

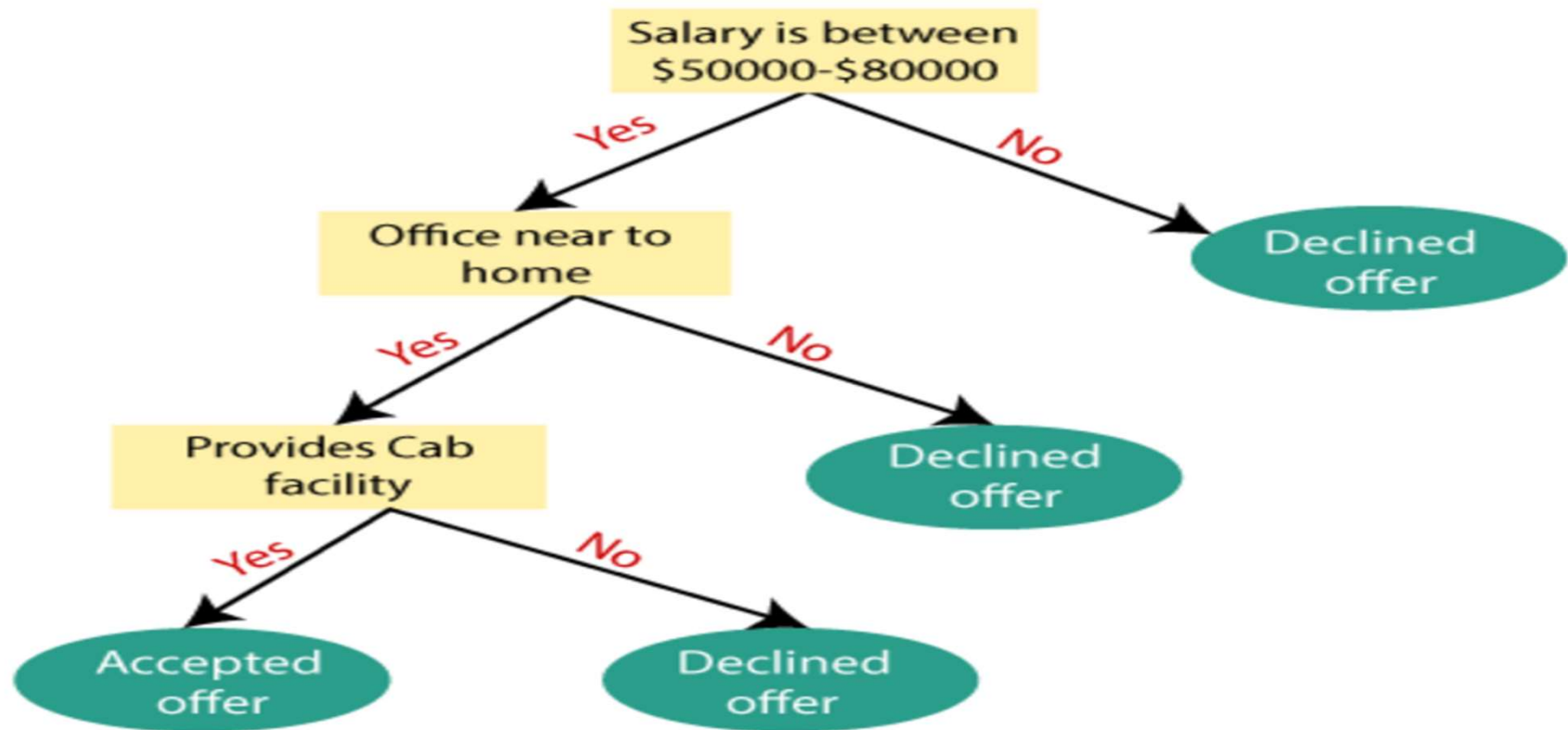
Definition

- Decision Tree is a classification and regression supervised ML model. It can be used for a classification problem and also a regression problem.
- But it is best suited for a classification problem
- In a decision tree the dataset is split into different groups based on the features.
- The aim is to use that feature for splitting which results in a purer group or in other words which results in a better classification.
- Decision Tree Model is based on 'recursive partitioning'.
- There are different algorithms based on different splitting criteria.

Basic Terminologies

- Root Node: Main node (from where the splitting starts)
- Decision Node: A point where decision about splitting needs to be taken. It is also called 'condition node'
- Parent / Child Nodes: In a split the split nodes are called child node, and the node from where they split is called parent node.
- Subnode: Same as child node
- Terminal Node/ Leaf Node: Last node of the tree that does not get divided further.
- Pruning: Once the decision tree is completed formed, it may be required to remove some of the splits and reshape the decision tree, this process is called pruning. It is completely opposite of splitting.
- In splitting the tree size increases, in pruning the tree size decreases.



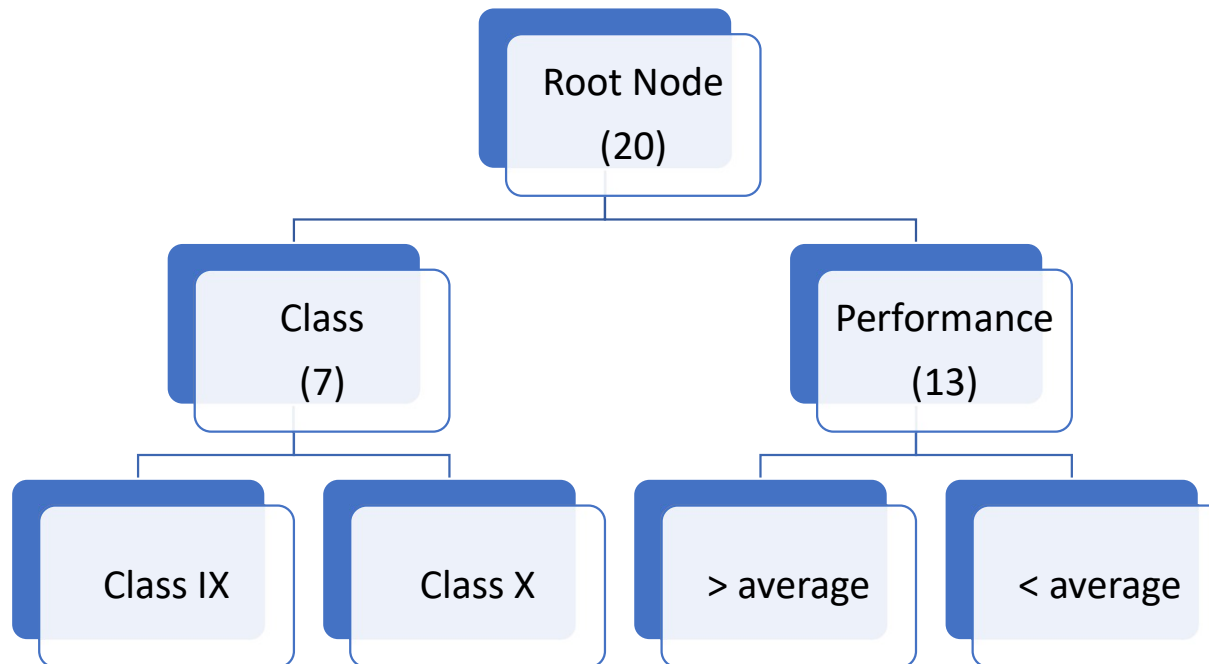


Constraints of tree splitting

- A decision tree aims to obtain pure nodes. For this it may go on splitting till the time it gets 100% pure nodes. That is the reason it is also known as a 'Greedy Algorithm'
- Since we do not want the tree to continue splitting till the last pure node, we predefine some constraints on when to stop splitting on when should it proceed to split.
- Minimum samples on a node to split
- Minimum samples for a terminal node
- Maximum depth
- Maximum number of terminal nodes
- Maximum features to be considered for a split

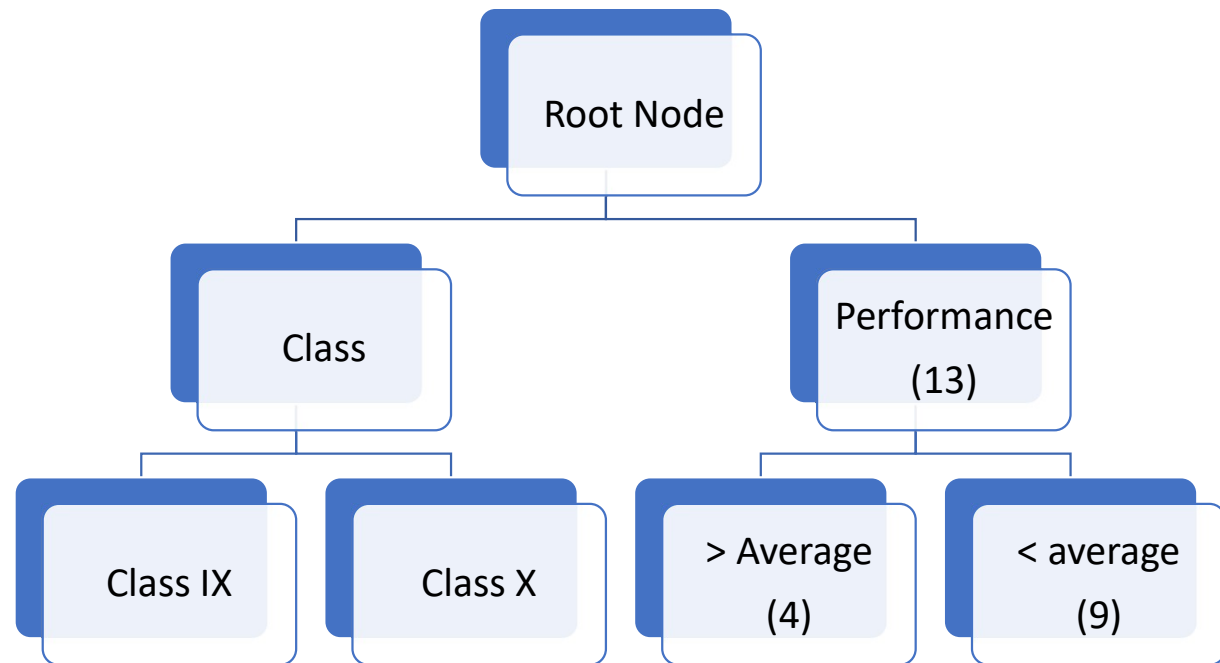
Minimum samples on a decision node to split

Lets say if the rule is that a node must have atleast 10 samples to go for further split



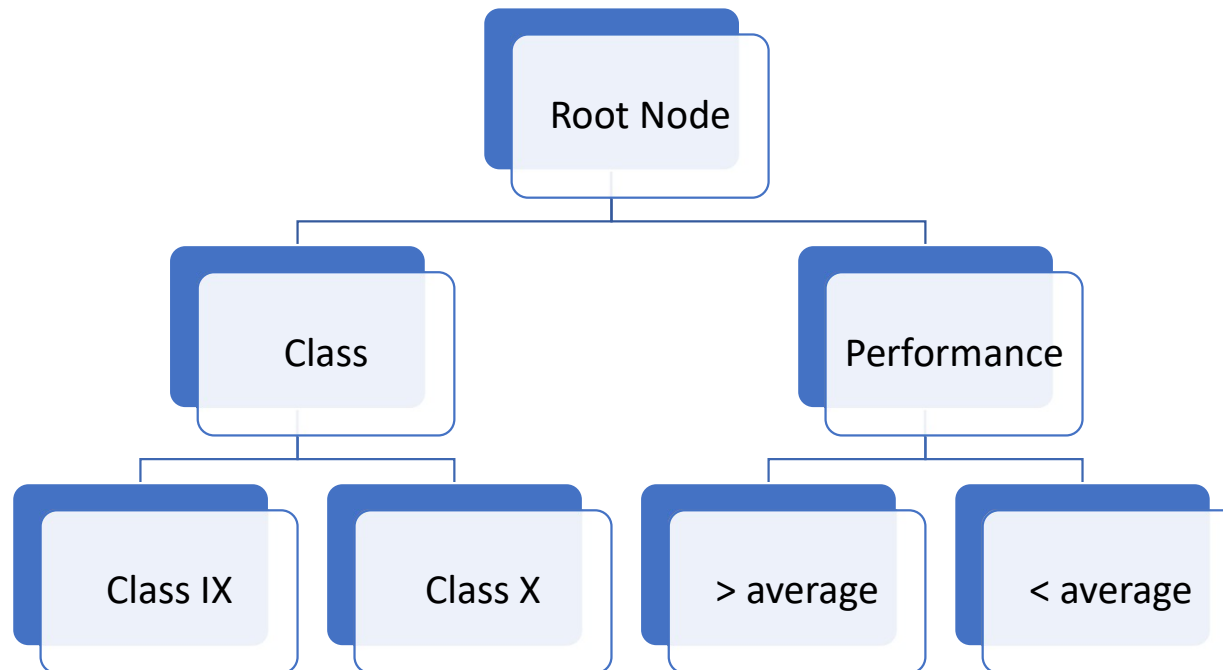
Minimum samples for a terminal node

- Suppose if the rule is that minimum samples in the terminal or leaf node must be 5

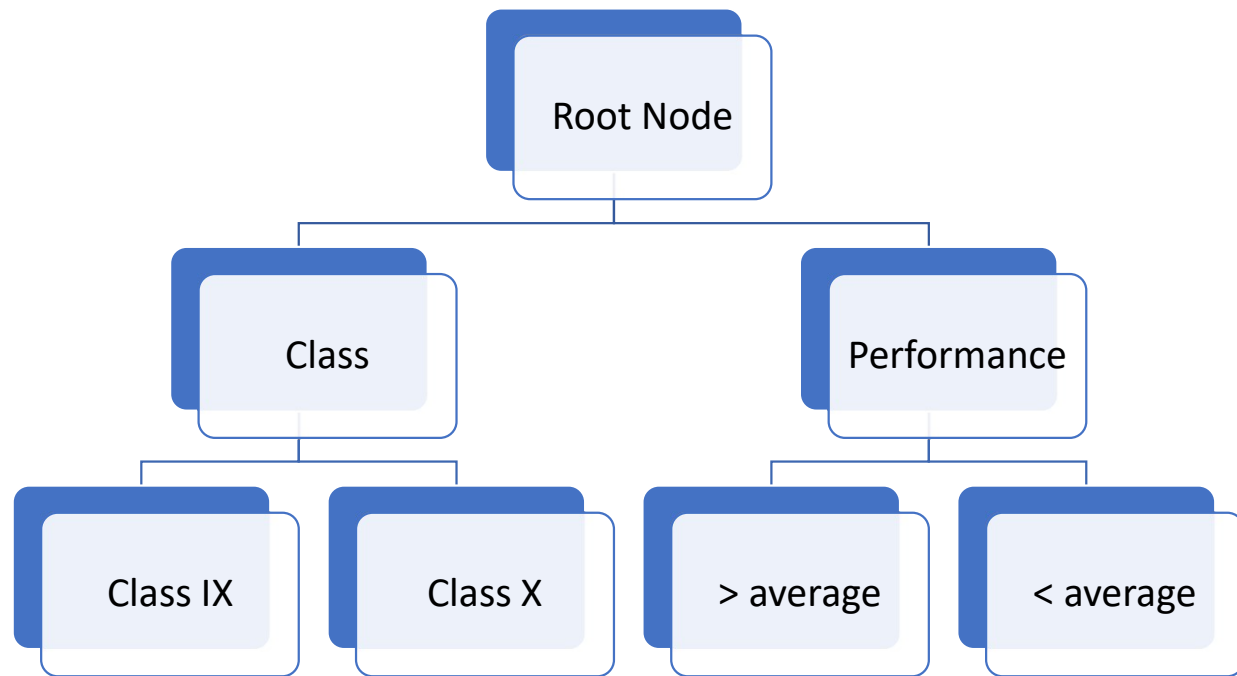


Maximum depth:

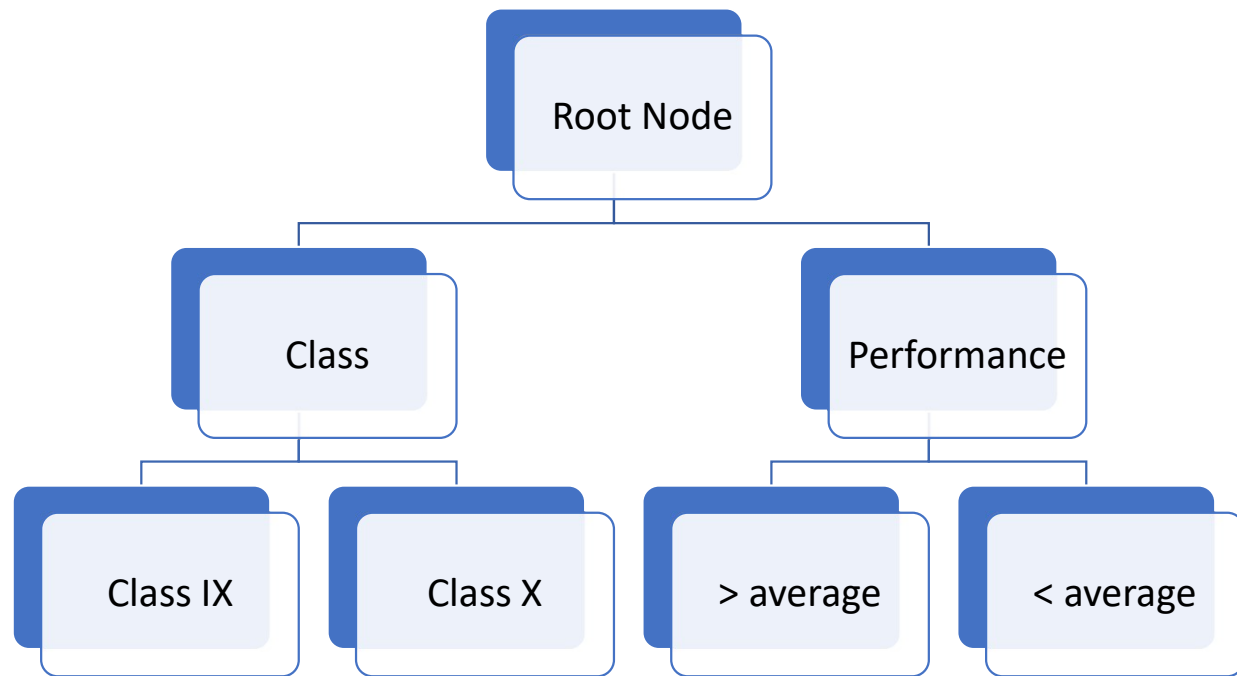
Depth is defined as the number of vertical splits the tree needs to do in order to give the final prediction. We can limit this depth by specifying the maximum depth acceptable for pruning purpose.

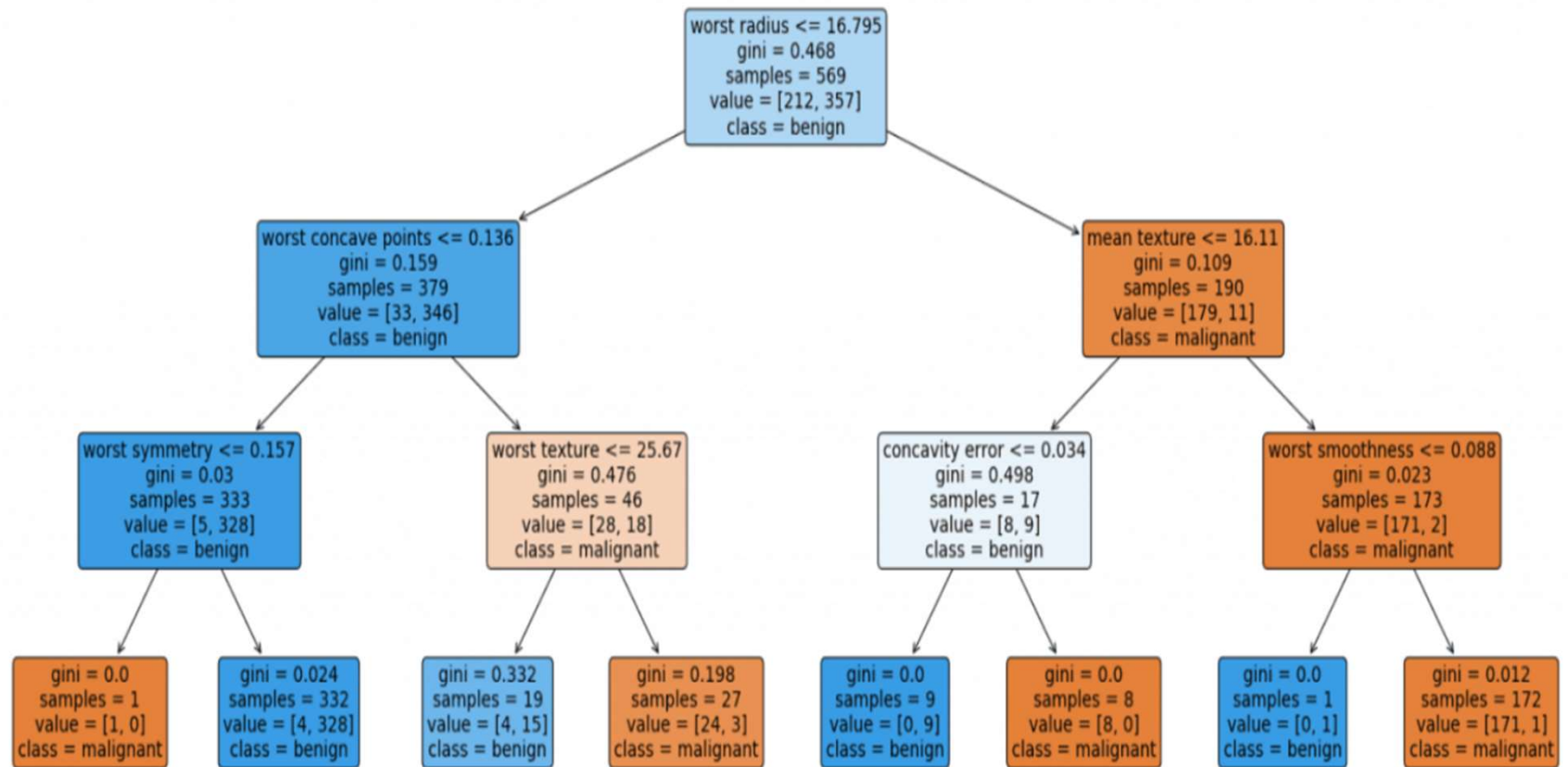


Maximum number of terminal nodes



Maximum features to be split: Generally when there are n features in a dataset, the maximum features to be split is set to \sqrt{n} or $\log(n)$





Different splitting algorithms

- CART: Classification And Regression Tree. It uses GINI as splitting criteria
- ID3: It uses Entropy as splitting criteria
- 99% of the time CART will be used in all DT models.

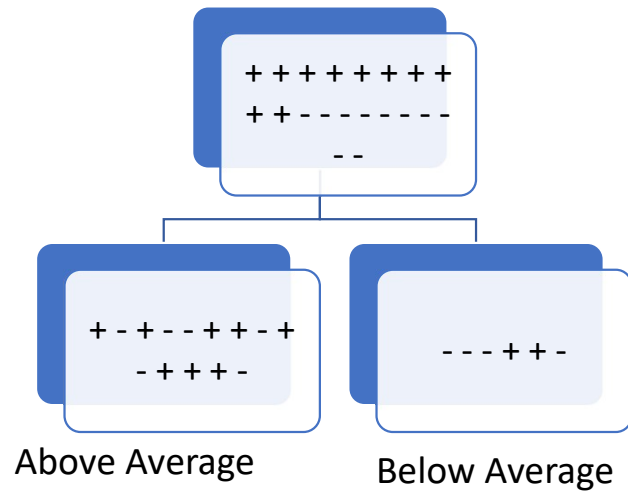
Four splitting methods are as follows:

- GINI
- ENTROPY

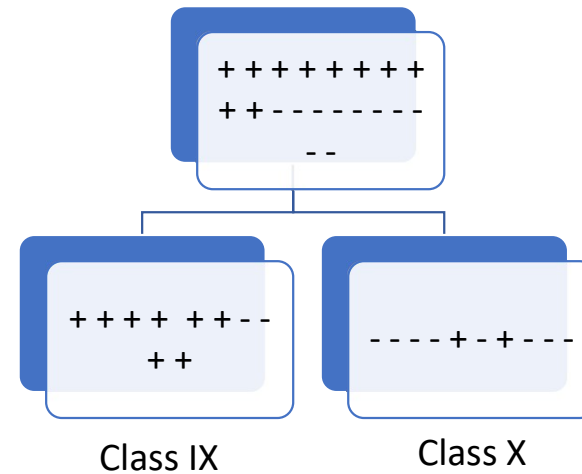
Example

Let the decision node have 20 students. + denotes those who play cricket and – denotes those who do not play cricket

Split by performance



Split by class



GINI

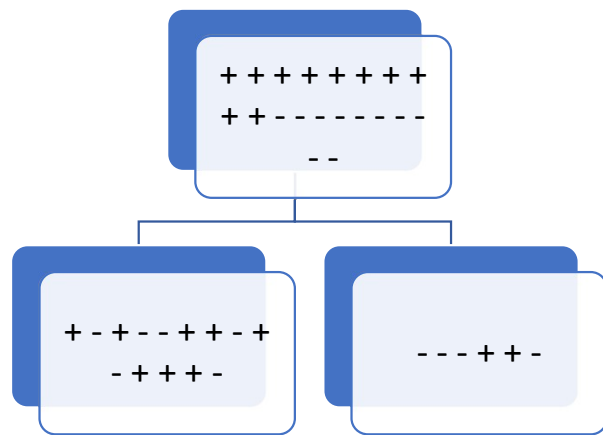
- GINI is a measure that talks about purity of a node.

$$\text{GINI} = P_1^2 + P_2^2 + P_3^2 + P_4^2 + \dots + P_n^2$$

$$\text{GI} = 1 - \text{GINI}$$

- Steps of calculation:
 1. Calculate GINI for each subnode
 2. Calculate GINI Impurity of each subnode
 3. Calculate the overall weight GINI impurity of the split.

GINI Calculation for split by performance



Above Average

Below Average

- Total = 14
 - Play cricket = 8
 - Not play cricket = 6
 - $P(+) = 8/14 = 0.57$
 - $P(-) = 6/14 = 0.43$
- Total = 6
 - Play cricket = 2
 - Not play cricket = 4
 - $P(+) = 2/6 = 0.33$
 - $P(-) = 4/6 = 0.67$

• Step I: GINI calculation of each sub node

- $GINI(AA) = (0.57)^2 + (0.43)^2 = 0.5098$
- $GINI(BA) = (0.33)^2 + (0.67)^2 = 0.5578$

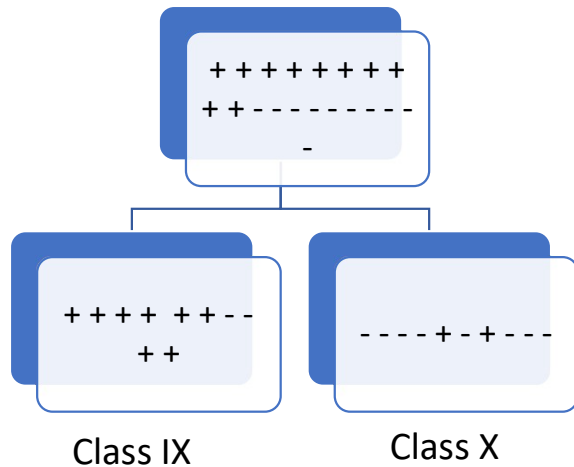
• Step II: GINI Impurity calculation for each sub node

- $GI(AA) = 1 - 0.5098 = 0.4902$
- $GI(BA) = 1 - 0.5578 = 0.4422$

• Step III: Weight GI of entire split

- $WGI = (14/20) * 0.4902 + (6/20) * 0.4422 = 0.4758$

GINI Calculation for split by class



- Total =
- Play cricket =
- Not play cricket =
- $P(+)$ =
- $P(-)$ =

- Total =
- Play cricket =
- Not play cricket =
- $P(+)$ =
- $P(-)$ =

• Step I: GINI calculation of each sub node

- $GINI(IX) =$
- $GINI(X) =$

• Step II: GINI Impurity calculation for each sub node

- $GI(IX) =$
- $GI(X) =$

• Step III: Weight GI of entire split

- $WGI = (W \text{ class IX}) * GI(IX) + (W \text{ class X}) * GI(X)$
- $WGI =$