



Decision Tree: The decision tree Algoedit belongs to the family of supervised machine learning algorithms. It can be used for both a classification problem as well as for regression problem.

The goal of this algorithm is to create a model that predicts the value of a target variable, for which the decision tree uses the tree structure to solve the problem in which the leaf node corresponds to a class label and attributes are represented on the internal node of the tree.

- Assumptions that we make while using the Decision Tree:
 - ↳ In the beginning, we consider the whole training set as the root.
 - ↳ Feature values are preferred to be categorical if the values continue then they are converted to discrete before building the model.
 - ↳ Based on attribute values records are distributed recursively.
 - ↳ We use a statistical method for ordering attributes as a root node or the internal node.

Entropy values range from 0 to 1", less the value of entropy more it is trusting able.

Entropy :

$$H(S) = -\text{probability of } \log_2(p+) - \text{probability of } \log_2(p-)$$

Where : $(p+)$ \rightarrow % of positive class
 $(p-)$ \rightarrow % of negative class

If we perform Entropy and get equal value like 50% yes or 50% no. Then this splitting will be going on unless and until we get a pure subset.

Pure subset: The pure subset is a situation where we will get either all yes or all no in this class.

In this case, we take other attributes to reach the leaf node and also have to take the entropy of those values and add it up to do the summation of all those entropy values for that we have the concept of info² gain.

Information Gain: Info² gain is used to decide which feature to split on at each step in building the tree. At each step, for small tree, we should choose the split that results in the perfect daughter nodes. A commonly used measure of purity is called info².

Information Gain:

$$\text{Gain}(S, A) = H(S) - \frac{|S_v|}{|S|} H(S_v)$$

The algo. calculates the info² gain for each split and split which is giving the highest value of info² gain is selected.

$S_v \rightarrow$ Total Sample after split

$S \rightarrow$ Total Sample

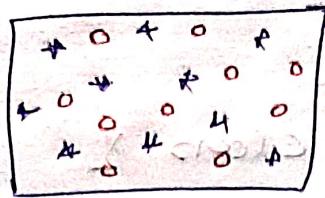
POORNIMA

The higher the value of info gain of the split the higher the chance of it getting selected for the particular split.

Gini Impurity:

Gini Impurity is a measurement used to build decision trees to determine how the features of a data set should split nodes to form the tree. Gini impurity of a data set is a number between 0 - 0.5.

→ what is decision Tree?



Total no. of student = 20

(+) Play cricket = 10 (50%)

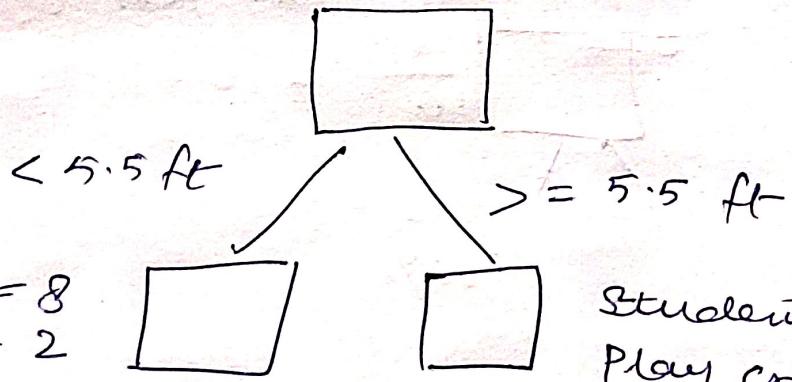
(*) Do not play cricket = 10

Features of student :-

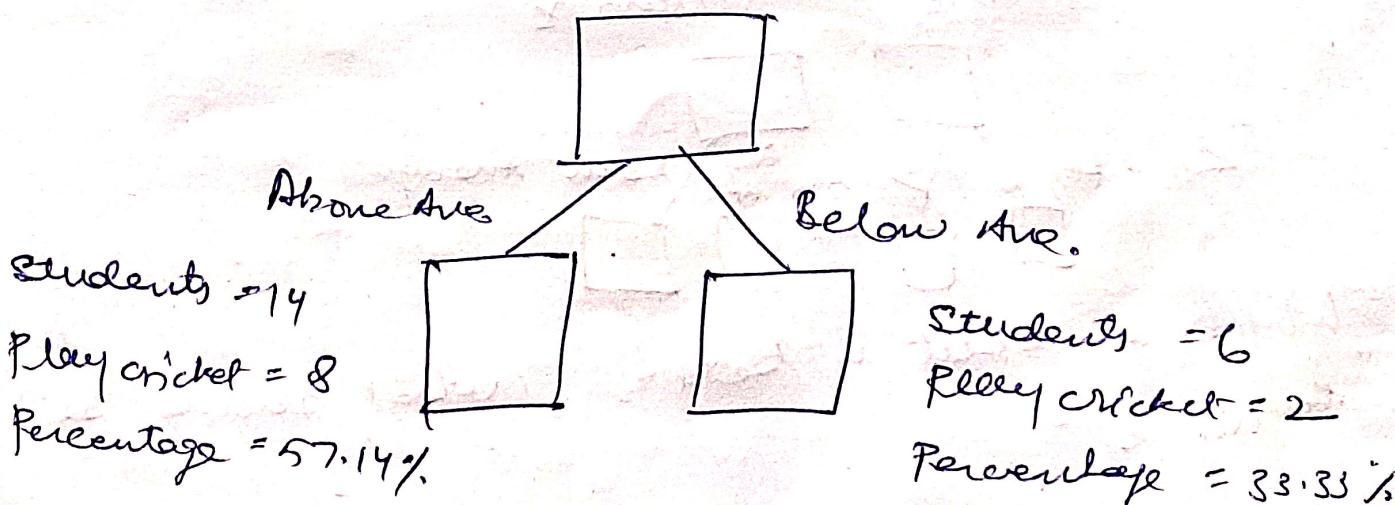
- Height
- Performance in class
- Class

Using these feature we want to train a model and predict whether they will play cricket or not.

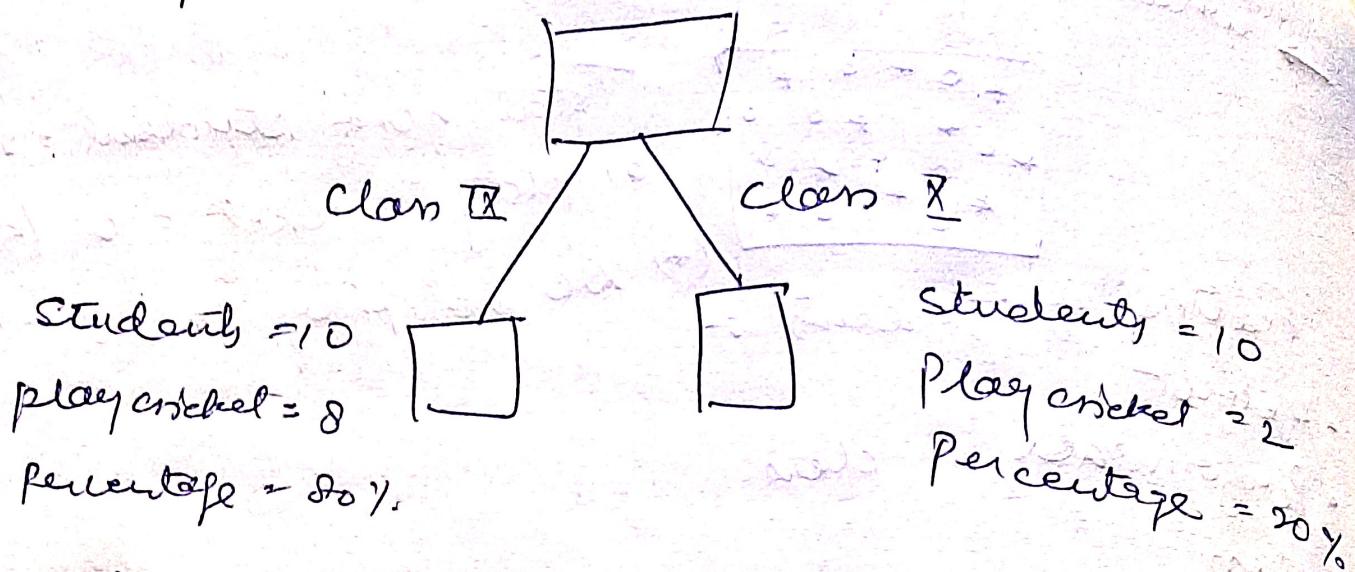
→ Split on Height



→ Split on Performance in class



Split on class



→ In decision tree we have to separate the classes.

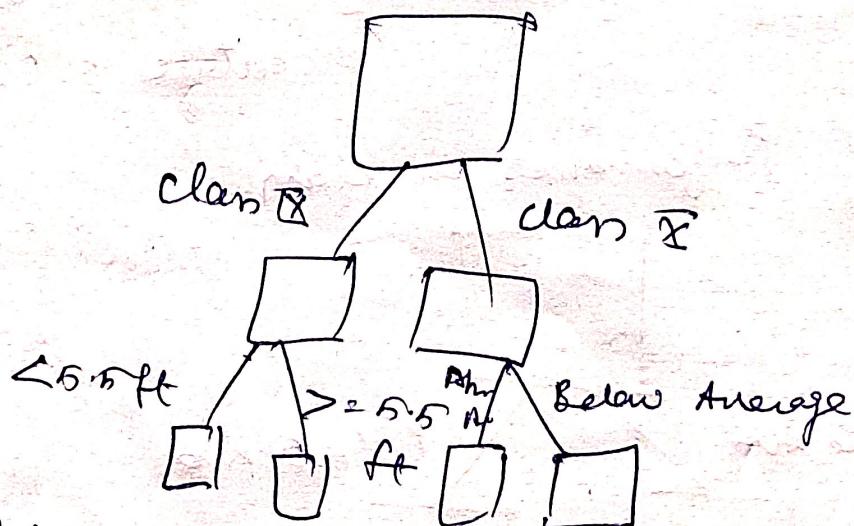
↳ Pure nodes.

↳ by which also we split the class.

Split on Played Cricket last year?



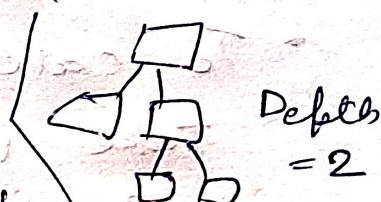
→ What is



which is split is better, what should be sequence,

Terminologies :-

- Root node
- Splitting
- Decision Nodes
- Left / Terminal Node
- Branch / Sub Tree
- Parent and child node
- Depth of tree (longest path of tree)



Select the best split point in Decision Tree:-

- Decision tree splits the nodes on all available variables
- Select the split which results in most homogeneous sub-nodes

$$\text{Gini Impurity} = 1 - \text{Gini}$$

→ Node splits are decided based on the Gini Impurity

Gini tells purity

Gini impurity tells the impurity of nodes

Selection of two random variable from the parent node. If node is pure then

Probability = 1

- Low the Gini impurity, higher the homogeneity of nodes.

~~(*)~~ Gini works on Categorical target not works on continuous targets e.g. only perform binary splits not in house predicting, etc.

Gini = Sum of Square of probabilities for each class / category

$$Gini = (p_1^2 + p_2^2 + p_3^2 + \dots + p_n^2)$$

- To calculate the gini impurity for split, take weighted gini impurity of both Sub-nodes of that split.

→ Split on Performance in class.

$$\text{Prob. play} = \frac{8}{14} = 0.57 \quad \left. \begin{array}{l} \text{Above Average} \\ \text{Average} \end{array} \right\}$$

$$\text{not play} = \frac{6}{14} = 0.43$$

$$\text{Play} \quad \frac{2}{6} = 0.33 \quad \left. \begin{array}{l} \text{Below Average} \\ \text{Average} \end{array} \right\} \quad \text{Student play cricket}$$

$$\text{not play} \quad \frac{4}{6} = 0.67 \quad \left. \begin{array}{l} \text{Student not play cricket} \\ \text{Not play} = 4 \end{array} \right\}$$

Gini impurity: Sub node Above Average:

$$1 - [(0.57) * (0.57) + (0.43) * (0.43)] = 0.49$$

Gini impurity Sub node Below Average:

$$1 - [(0.33) * (0.33) + (0.67) * (0.67)] = 0.44$$

Weight of node = $\frac{\text{No. node in child}}{\text{No. nodes in Parent}}$

for both class Above Average weight $\Rightarrow \frac{14}{20}$
Below Average = $\frac{6}{20}$

Weighted Gini Impurity: Performance in class
 $(\frac{14}{20}) * 0.49 + (\frac{6}{20}) * 0.44 = 0.475$

→ Split on class:

Split Performance in class	weighted Gini Impurity
Class	0.32