<div align="center">**Experiment No.: 10**</div>

**Title:** Implementation of Decision Tree.

**Objectives:** To learn Decision Tree.

**Theory:**

A decision tree is a simple representation for classifying examples. Assume that all of the input features have finite discrete domains, and there is a single target feature called the "classification". Each element of the domain of the classification is called a *class*. A decision tree or a classification tree is a tree in which each internal (non-leaf) node is labeled with an input feature. The arcs coming from a node labeled with an input feature are labeled with each of the possible values of the target or output feature or the arc leads to a subordinate decision node on a different input feature. Each leaf of the tree is labeled with a class or a probability distribution over the classes, signifying that the data set has been classified by the tree into either a specific class, or into a particular probability distribution (which, if the decision tree is well-constructed, is skewed towards certain subsets of classes).

A tree is built by splitting the source set, constituting the root node of the tree, into subsets - which constitute the successor children. The splitting is based on a set of splitting rules based on classification features. This process is repeated on each derived subset in a recursive manner called recursive partitioning. The recursion is completed when the subset at a node has all the same values of the target variable, or when splitting no longer adds value to the predictions.

**Assumptions we make while using Decision tree:**

- At the beginning, we consider the whole training set as the root.
- Attributes are assumed to be categorical for information gain and for gini index, attributes are assumed to be continuous.
- On the basis of attribute values records are distributed recursively.
- We use statistical methods for ordering attributes as root or internal node.

**Pseudo code:**

1. Find the best attribute and place it on the root node of the tree.

2. Now, split the training set of the dataset into subsets. While making the subset make sure that each subset of training dataset should have the same value for an attribute.

3. Find leaf nodes in all branches by repeating 1 and 2 on each subset.

**Implementation procedure**:

1. Building Phase

- Preprocess the dataset.
- Split the dataset from train and test.
- Train the classifier.

2. Operational Phase

- Make predictions.
- Calculate the accuracy.