```python
import pandas as pd
import numpy as np
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier,plot_tree
from sklearn.metrics import accuracy_score,classification_report,confusion_matri
import matplotlib.pyplot as plt
```
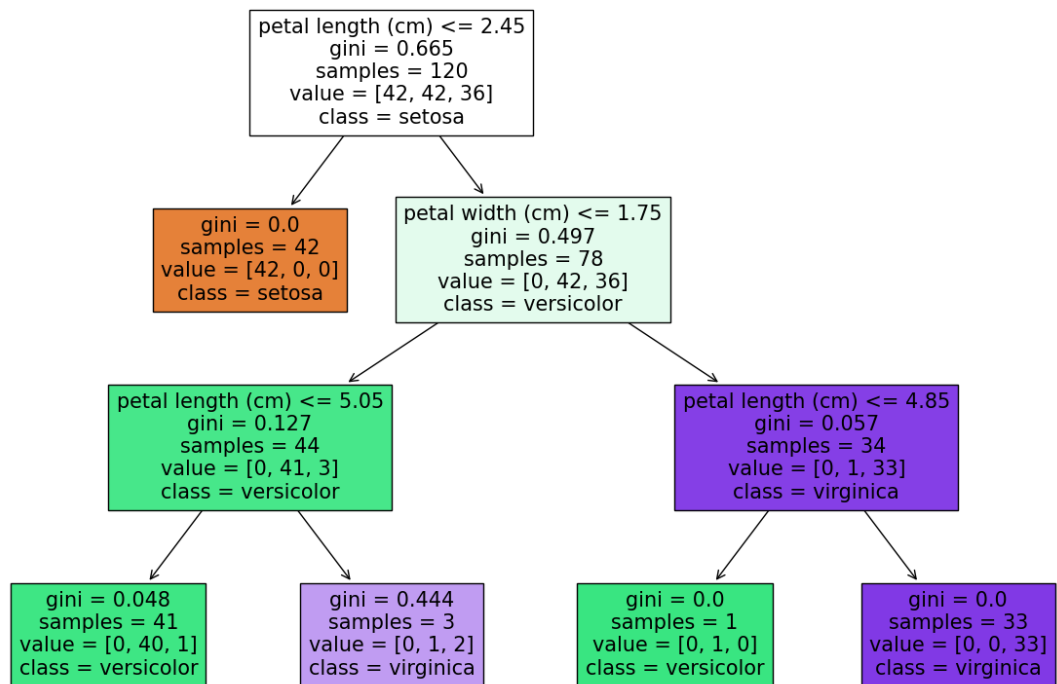
```python
X,y=load_iris(return_X_y=True)
print(X)
print(y)
```

```
[[5.1 3.5 1.4 0.2]
 [4.9 3.  1.4 0.2]
 [4.7 3.2 1.3 0.2]
 [4.6 3.1 1.5 0.2]
 [5.  3.6 1.4 0.2]
 [5.4 3.9 1.7 0.4]
 [4.6 3.4 1.4 0.3]
 [5.  3.4 1.5 0.2]
 [4.4 2.9 1.4 0.2]
 [4.9 3.1 1.5 0.1]
 [5.4 3.7 1.5 0.2]
 [4.8 3.4 1.6 0.2]
 [4.8 3.  1.4 0.1]
 [4.3 3.  1.1 0.1]
 [5.8 4.  1.2 0.2]
 [5.7 4.4 1.5 0.4]
 [5.4 3.9 1.3 0.4]
 [5.1 3.5 1.4 0.3]
 [5.7 3.8 1.7 0.3]
 [5.1 3.8 1.5 0.3]
 [5.4 3.4 1.7 0.2]
 [5.1 3.7 1.5 0.4]
 [4.6 3.6 1.  0.2]
 [5.1 3.3 1.7 0.5]
 [4.8 3.4 1.9 0.2]
 [5.  3.  1.6 0.2]
 [5.  3.4 1.6 0.4]
 [5.2 3.5 1.5 0.2]
 [5.2 3.4 1.4 0.2]
 [4.7 3.2 1.6 0.2]
 [4.8 3.1 1.6 0.2]
 [5.4 3.4 1.5 0.4]
 [5.2 4.1 1.5 0.1]
 [5.5 4.2 1.4 0.2]
 [4.9 3.1 1.5 0.2]
 [5.  3.2 1.2 0.2]
 [5.5 3.5 1.3 0.2]
 [4.9 3.6 1.4 0.1]
 [4.4 3.  1.3 0.2]
 [5.1 3.4 1.5 0.2]
 [5.  3.5 1.3 0.3]
 [4.5 2.3 1.3 0.3]
 [4.4 3.2 1.3 0.2]
 [5.  3.5 1.6 0.6]
 [5.1 3.8 1.9 0.4]
 [4.8 3.  1.4 0.3]
 [5.1 3.8 1.6 0.2]
 [4.6 3.2 1.4 0.2]
 [5.3 3.7 1.5 0.2]
 [5.  3.3 1.4 0.2]
 [7.  3.2 4.7 1.4]
 [6.4 3.2 4.5 1.5]
 [6.9 3.1 4.9 1.5]
 [5.5 2.3 4.  1.3]
 [6.5 2.8 4.6 1.5]
 [5.7 2.8 4.5 1.3]
 [6.3 3.3 4.7 1.6]
 [4.9 2.4 3.3 1. ]
 [6.6 2.9 4.6 1.3]
 [5.2 2.7 3.9 1.4]
```

```
[5.  2.  3.5 1. ]
[5.9 3.  4.2 1.5]
[6.  2.2 4.  1. ]
[6.1 2.9 4.7 1.4]
[5.6 2.9 3.6 1.3]
[6.7 3.1 4.4 1.4]
[5.6 3.  4.5 1.5]
[5.8 2.7 4.1 1. ]
[6.2 2.2 4.5 1.5]
[5.6 2.5 3.9 1.1]
[5.9 3.2 4.8 1.8]
[6.1 2.8 4.  1.3]
[6.3 2.5 4.9 1.5]
[6.1 2.8 4.7 1.2]
[6.4 2.9 4.3 1.3]
[6.6 3.  4.4 1.4]
[6.8 2.8 4.8 1.4]
[6.7 3.  5.  1.7]
[6.  2.9 4.5 1.5]
[5.7 2.6 3.5 1. ]
[5.5 2.4 3.8 1.1]
[5.5 2.4 3.7 1. ]
[5.8 2.7 3.9 1.2]
[6.  2.7 5.1 1.6]
[5.4 3.  4.5 1.5]
[6.  3.4 4.5 1.6]
[6.7 3.1 4.7 1.5]
[6.3 2.3 4.4 1.3]
[5.6 3.  4.1 1.3]
[5.5 2.5 4.  1.3]
[5.5 2.6 4.4 1.2]
[6.1 3.  4.6 1.4]
[5.8 2.6 4.  1.2]
[5.  2.3 3.3 1. ]
[5.6 2.7 4.2 1.3]
[5.7 3.  4.2 1.2]
[5.7 2.9 4.2 1.3]
[6.2 2.9 4.3 1.3]
[5.1 2.5 3.  1.1]
[5.7 2.8 4.1 1.3]
[6.3 3.3 6.  2.5]
[5.8 2.7 5.1 1.9]
[7.1 3.  5.9 2.1]
[6.3 2.9 5.6 1.8]
[6.5 3.  5.8 2.2]
[7.6 3.  6.6 2.1]
[4.9 2.5 4.5 1.7]
[7.3 2.9 6.3 1.8]
[6.7 2.5 5.8 1.8]
[7.2 3.6 6.1 2.5]
[6.5 3.2 5.1 2. ]
[6.4 2.7 5.3 1.9]
[6.8 3.  5.5 2.1]
[5.7 2.5 5.  2. ]
[5.8 2.8 5.1 2.4]
[6.4 3.2 5.3 2.3]
[6.5 3.  5.5 1.8]
[7.7 3.8 6.7 2.2]
[7.7 2.6 6.9 2.3]
[6.  2.2 5.  1.5]
```

```
            [6.9 3.2 5.7 2.3]
            [5.6 2.8 4.9 2. ]
            [7.7 2.8 6.7 2. ]
            [6.3 2.7 4.9 1.8]
            [6.7 3.3 5.7 2.1]
            [7.2 3.2 6.  1.8]
            [6.2 2.8 4.8 1.8]
            [6.1 3.  4.9 1.8]
            [6.4 2.8 5.6 2.1]
            [7.2 3.  5.8 1.6]
            [7.4 2.8 6.1 1.9]
            [7.9 3.8 6.4 2. ]
            [6.4 2.8 5.6 2.2]
            [6.3 2.8 5.1 1.5]
            [6.1 2.6 5.6 1.4]
            [7.7 3.  6.1 2.3]
            [6.3 3.4 5.6 2.4]
            [6.4 3.1 5.5 1.8]
            [6.  3.  4.8 1.8]
            [6.9 3.1 5.4 2.1]
            [6.7 3.1 5.6 2.4]
            [6.9 3.1 5.1 2.3]
            [5.8 2.7 5.1 1.9]
            [6.8 3.2 5.9 2.3]
            [6.7 3.3 5.7 2.5]
            [6.7 3.  5.2 2.3]
            [6.3 2.5 5.  1.9]
            [6.5 3.  5.2 2. ]
            [6.2 3.4 5.4 2.3]
            [5.9 3.  5.1 1.8]]
            [0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0 0
             0 0 0 0 0 0 0 0 0 0 0 0 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
             1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 2 2 2
             2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2
             2 2]
```

```python
x_train,x_test,y_train,y_test=train_test_split(X,y,test_size=0.2,random_state=33
```

```python
cl=DecisionTreeClassifier()
cl.fit(x_train,y_train)
print("Accuracy",cl.score(x_test,y_test))
y_pred=cl.predict(x_test)
print(confusion_matrix(y_test,y_pred))
```

```
Accuracy 0.8666666666666667
[[ 8  0  0]
 [ 0  8  0]
 [ 0  4 10]]
```

```python
iris=load_iris()
fig = plt.figure(figsize=(15,10))
_ = plot_tree(cl,
                feature_names=iris.feature_names,
                class_names=iris.target_names,
                filled=True)
```
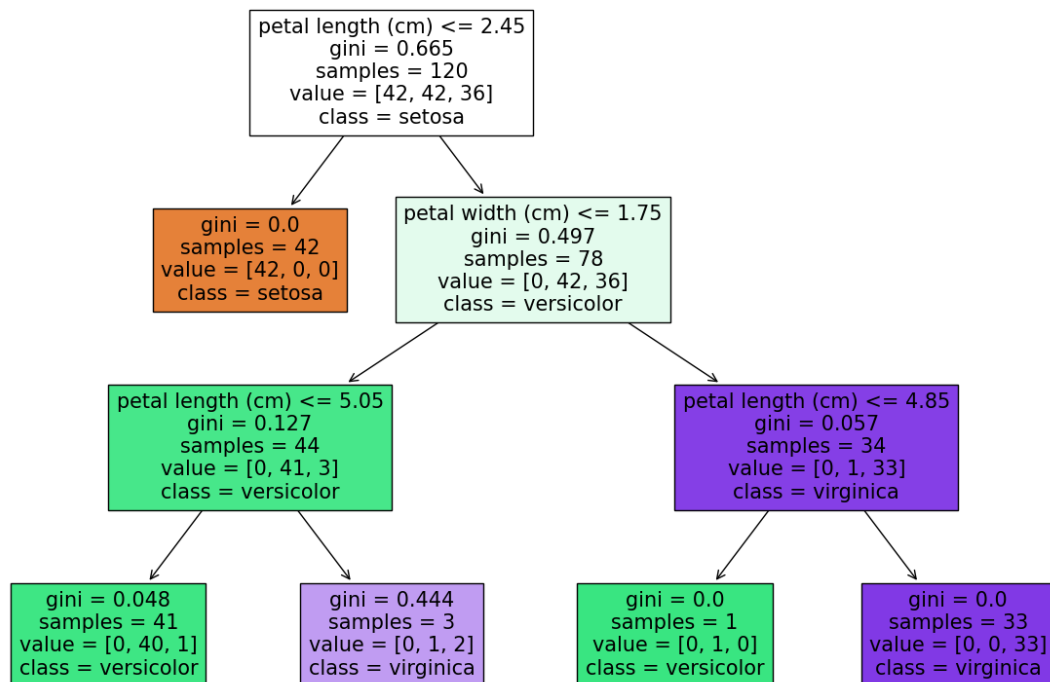
```
                        petal length (cm) <= 2.45
                               gini = 0.665
                              samples = 120
                            value = [42, 42, 36]
                              class = setosa
```

```
         gini = 0.0                    petal width (cm) <= 1.75
       samples = 42                          gini = 0.497
     value = [42, 0, 0]                     samples = 78
       class = setosa                     value = [0, 42, 36]
                                          class = versicolor
```

```
petal length (cm) <= 5.05                        petal length (cm) <= 4.85
       gini = 0.127                                      gini = 0.057
      samples = 44                                      samples = 34
    value = [0, 41, 3]                                value = [0, 1, 33]
    class = versicolor                                class = virginica
```

```
gini = 0.048      gini = 0.444          gini = 0.0          gini = 0.0
samples = 41      samples = 3          samples = 1         samples = 33
value = [0, 40, 1] value = [0, 1, 2]   value = [0, 1, 0]   value = [0, 0, 33]
class = versicolor class = virginica   class = versicolor  class = virginica
```

# Pre-Pruning

```
In [ ]: cl_pre=DecisionTreeClassifier(max_depth=3,max_leaf_nodes=6)
        cl_pre.fit(x_train,y_train)
        print("Accuracy",cl_pre.score(x_test,y_test))
        y_pred=cl_pre.predict(x_test)
        print(confusion_matrix(y_test,y_pred))
```

```
Accuracy 0.9
[[ 8  0  0]
 [ 0  8  0]
 [ 0  3 11]]
```

```
In [ ]: iris=load_iris()
        fig = plt.figure(figsize=(15,10))
        _ = plot_tree(cl_pre,
                        feature_names=iris.feature_names,
                        class_names=iris.target_names,
                        filled=True)
```

```
                         petal length (cm) <= 2.45
                               gini = 0.665
                              samples = 120
                            value = [42, 42, 36]
                               class = setosa
```

```
              gini = 0.0                    petal width (cm) <= 1.75
            samples = 42                          gini = 0.497
          value = [42, 0, 0]                      samples = 78
            class = setosa                      value = [0, 42, 36]
                                                class = versicolor
```

```
    petal length (cm) <= 5.05                              petal length (cm) <= 4.85
          gini = 0.127                                          gini = 0.057
         samples = 44                                          samples = 34
        value = [0, 41, 3]                                   value = [0, 1, 33]
       class = versicolor                                    class = virginica
```

```
  gini = 0.048       gini = 0.444              gini = 0.0          gini = 0.0
 samples = 41       samples = 3             samples = 1          samples = 33
value = [0, 40, 1]  value = [0, 1, 2]      value = [0, 1, 0]    value = [0, 0, 33]
class = versicolor  class = virginica      class = versicolor   class = virginica
```

# Post Pruning

```
In [ ]:  path=cl.cost_complexity_pruning_path(x_train,y_train)
         ccp_alphas,impurities=path.ccp_alphas,path.impurities
         print("ccp alpha wil give list of values :\n",ccp_alphas)
         print("*********************************************************")
         print("Impurities in Decision Tree :\n",impurities)
```

```
         ccp alpha wil give list of values :
          [0.         0.00813008 0.01111111 0.01617647 0.01921964 0.26030954
          0.34192308]
         *********************************************************
         Impurities in Decision Tree :
          [0.         0.01626016 0.02737127 0.04354774 0.06276738 0.32307692
          0.665      ]
```

```
In [ ]:  clfs=[]
         for ccp_alpha in ccp_alphas:
             clf=DecisionTreeClassifier(random_state=23,ccp_alpha=ccp_alpha)
             clf.fit(x_train,y_train)
             clfs.append(clf)
         print("Last node in Decision Tree is {} and ccp_alpha for last node is {}".forma
```
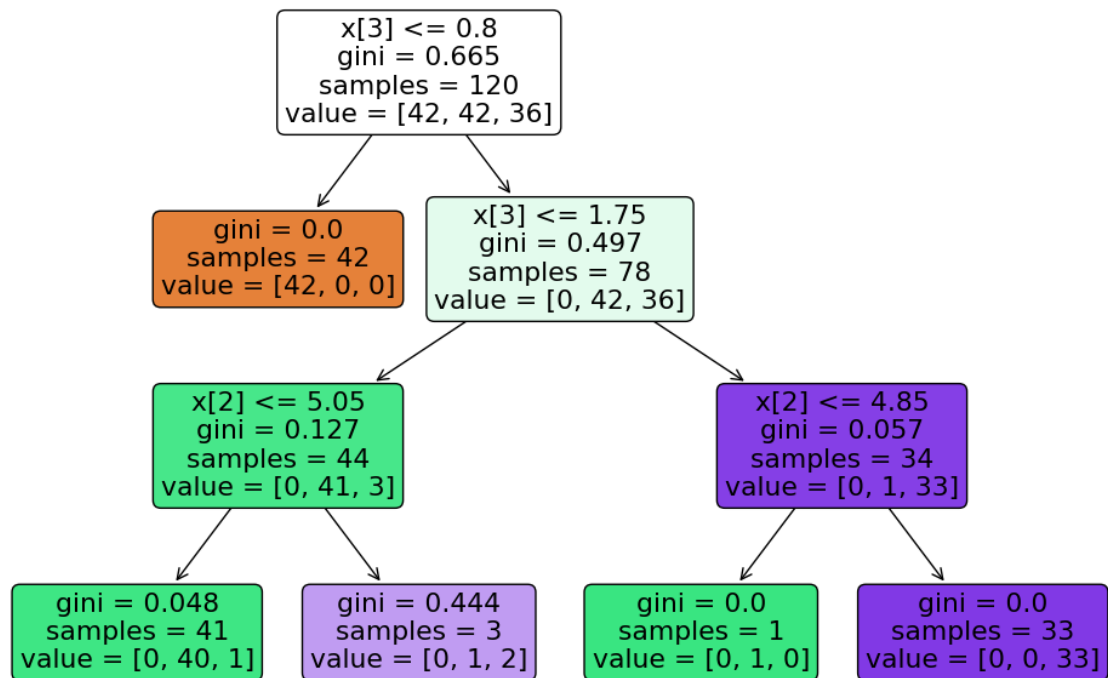
```
         Last node in Decision Tree is DecisionTreeClassifier(ccp_alpha=0.34192307692307
         71, random_state=23) and ccp_alpha for last node is 0.3419230769230771
```

```
In [ ]:  train_scores = [clf.score(x_train, y_train) for clf in clfs]
         test_scores = [clf.score(x_test, y_test) for clf in clfs]
         fig, ax = plt.subplots()
         ax.set_xlabel("alpha")
         ax.set_ylabel("accuracy")
         ax.set_title("Accuracy vs alpha for training and testing sets")
         ax.plot(ccp_alphas, train_scores, marker='o', label="train",drawstyle="steps-pos
         ax.plot(ccp_alphas, test_scores, marker='o', label="test",drawstyle="steps-post"
```

```
ax.legend()
plt.show()
```



Accuracy vs alpha for training and testing sets

```
from sklearn import tree
clf=DecisionTreeClassifier(random_state=0,ccp_alpha=0.015)
clf.fit(x_train,y_train)
plt.figure(figsize=(12,8))
tree.plot_tree(clf,rounded=True,filled=True)
plt.show()
```

```
                        x[3] <= 0.8
                        gini = 0.665
                      samples = 120
                    value = [42, 42, 36]
                    /                  \
          gini = 0.0              x[3] <= 1.75
        samples = 42             gini = 0.497
      value = [42, 0, 0]         samples = 78
                              value = [0, 42, 36]
                              /                  \
                    x[2] <= 5.05            x[2] <= 4.85
                    gini = 0.127            gini = 0.057
                    samples = 44            samples = 34
                  value = [0, 41, 3]      value = [0, 1, 33]
                  /            \          /            \
        gini = 0.048    gini = 0.444    gini = 0.0    gini = 0.0
       samples = 41     samples = 3    samples = 1    samples = 33
     value = [0, 40, 1] value = [0, 1, 2] value = [0, 1, 0] value = [0, 0, 33]
```

```python
import numpy as np
from collections import Counter
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix
```

```python
class Node:
    '''
    Helper class which implements a single tree node.
    '''
    def __init__(self, feature=None, threshold=None, data_left=None, data_right=
        self.feature = feature
        self.threshold = threshold
        self.data_left = data_left
        self.data_right = data_right
        self.gain = gain
        self.value = value


class DecisionTree:
    '''
    Class which implements a decision tree classifier algorithm.
    '''
    def __init__(self, min_samples_split=2, max_depth=5):
        self.min_samples_split = min_samples_split
        self.max_depth = max_depth
        self.root = None

    @staticmethod
    def _entropy(s):
        '''
        Helper function, calculates entropy from an array of integer values.

        :param s: list
        :return: float, entropy value
        '''
        # Convert to integers to avoid runtime errors
        counts = np.bincount(np.array(s, dtype=np.int64))
        # Probabilities of each class label
        percentages = counts / len(s)

        # Caclulate entropy
        entropy = 0
        for pct in percentages:
            if pct > 0:
                entropy += pct * np.log2(pct)
        return -entropy

    def _information_gain(self, parent, left_child, right_child):
        '''
        Helper function, calculates information gain from a parent and two child

        :param parent: list, the parent node
        :param left_child: list, left child of a parent
        :param right_child: list, right child of a parent
        :return: float, information gain
        '''
        num_left = len(left_child) / len(parent)
        num_right = len(right_child) / len(parent)
```

```python
        # One-Liner which implements the previously discussed formula
        return self._entropy(parent) - (num_left * self._entropy(left_child) + r

    def _best_split(self, X, y):
        '''
        Helper function, calculates the best split for given features and target

        :param X: np.array, features
        :param y: np.array or list, target
        :return: dict
        '''
        best_split = {}
        best_info_gain = -1
        n_rows, n_cols = X.shape

        # For every dataset feature
        for f_idx in range(n_cols):
            X_curr = X[:, f_idx]
            # For every unique value of that feature
            for threshold in np.unique(X_curr):
                # Construct a dataset and split it to the left and right parts
                # Left part includes records lower or equal to the threshold
                # Right part includes records higher than the threshold
                df = np.concatenate((X, y.reshape(1, -1).T), axis=1)
                df_left = np.array([row for row in df if row[f_idx] <= threshold
                df_right = np.array([row for row in df if row[f_idx] > threshold

                # Do the calculation only if there's data in both subsets
                if len(df_left) > 0 and len(df_right) > 0:
                    # Obtain the value of the target variable for subsets
                    y = df[:, -1]
                    y_left = df_left[:, -1]
                    y_right = df_right[:, -1]

                    # Caclulate the information gain and save the split paramete
                    # if the current split if better then the previous best
                    gain = self._information_gain(y, y_left, y_right)
                    if gain > best_info_gain:
                        best_split = {
                            'feature_index': f_idx,
                            'threshold': threshold,
                            'df_left': df_left,
                            'df_right': df_right,
                            'gain': gain
                        }
                        best_info_gain = gain
        return best_split

    def _build(self, X, y, depth=0):
        '''
        Helper recursive function, used to build a decision tree from the input

        :param X: np.array, features
        :param y: np.array or list, target
        :param depth: current depth of a tree, used as a stopping criteria
        :return: Node
        '''

        n_rows, n_cols = X.shape
```

```python
            # Check to see if a node should be leaf node
            if n_rows >= self.min_samples_split and depth <= self.max_depth:
                # Get the best split
                best = self._best_split(X, y)
                # If the split isn't pure
                if best['gain'] > 0:
                    # Build a tree on the left
                    left = self._build(
                        X=best['df_left'][:, :-1],
                        y=best['df_left'][:, -1],
                        depth=depth + 1
                    )
                    right = self._build(
                        X=best['df_right'][:, :-1],
                        y=best['df_right'][:, -1],
                        depth=depth + 1
                    )
                    return Node(
                        feature=best['feature_index'],
                        threshold=best['threshold'],
                        data_left=left,
                        data_right=right,
                        gain=best['gain']
                    )
            # Leaf node - value is the most common target value
            return Node(
                value=Counter(y).most_common(1)[0][0]
            )

    def fit(self, X, y):
        '''
        Function used to train a decision tree classifier model.

        :param X: np.array, features
        :param y: np.array or list, target
        :return: None
        '''
        # Call a recursive function to build the tree
        self.root = self._build(X, y)

    def _predict(self, x, tree):
        '''
        Helper recursive function, used to predict a single instance (tree trave

        :param x: single observation
        :param tree: built tree
        :return: float, predicted class
        '''
        # Leaf node
        if tree.value != None:
            return tree.value
        feature_value = x[tree.feature]

        # Go to the left
        if feature_value <= tree.threshold:
            return self._predict(x=x, tree=tree.data_left)

        # Go to the right
        if feature_value > tree.threshold:
            return self._predict(x=x, tree=tree.data_right)
```

```python
    def predict(self, X):
        '''
        Function used to classify new instances.

        :param X: np.array, features
        :return: np.array, predicted classes
        '''
        # Call the _predict() function for every observation
        return [self._predict(x, self.root) for x in X]
```

In [ ]:
```python
iris = load_iris()

X = iris['data']
y = iris['target']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_
model=DecisionTree()
model.fit(X_train,y_train)
y_pred=model.predict(X_test)

print("Test Results:",y_test)
print("Predicted Results: ",y_pred)
# print(np.concatenate((np.array(y_test).reshape(len(y_test),1),np.array(y_pred)

print("Confusion Matrix\n",confusion_matrix(y_test,y_pred))
```

```
Test Results: [1 0 2 1 1 0 1 2 1 1 2 0 0 0 0 1 2 1 1 2 0 2 0 2 2 2 2 2 0 0]
Predicted Results:  [1.0, 0.0, 2.0, 1.0, 1.0, 0.0, 1.0, 2.0, 1.0, 1.0, 2.0, 0.
0, 0.0, 0.0, 0.0, 1.0, 2.0, 1.0, 1.0, 2.0, 0.0, 2.0, 0.0, 2.0, 2.0, 2.0, 2.0,
2.0, 0.0, 0.0]
Confusion Matrix
 [[10  0  0]
 [ 0  9  0]
 [ 0  0 11]]
```