

Usama Arif Roll No 14

```
In [177]: 1 import pandas as pd
2 from sklearn.preprocessing import LabelEncoder,MinMaxScaler
3 from sklearn.tree import DecisionTreeClassifier,plot_tree
4 from sklearn.metrics import confusion_matrix,accuracy_score
5 from sklearn.model_selection import train_test_split
6 import matplotlib.pyplot as plt
7 import seaborn as sns
8 import math
```

reading a csv file

```
In [178]: 1 iris=pd.read_csv('iris.csv')
2 iris.head(10)
3
```

Out[178]:

	sepal_length	sepal_width	petal_length	petal_width	species
0	5.1	3.5	1.4	0.2	setosa
1	4.9	3.0	1.4	0.2	setosa
2	4.7	3.2	1.3	0.2	setosa
3	4.6	3.1	1.5	0.2	setosa
4	5.0	3.6	1.4	0.2	setosa
5	5.4	3.9	1.7	0.4	setosa
6	4.6	3.4	1.4	0.3	setosa
7	5.0	3.4	1.5	0.2	setosa
8	4.4	2.9	1.4	0.2	setosa
9	4.9	3.1	1.5	0.1	setosa

seprating the data

```
In [179]: 1 x=iris.drop('species',axis=1)
2 y=iris['species']
```

```
In [180]: 1 scaler=MinMaxScaler()
2 x_scaled=scaler.fit_transform(x)
3 label_encoder = LabelEncoder()
4 y_encoded = label_encoder.fit_transform(y)
```

fitting the model and training the data

```
In [186]: 1 clf=DecisionTreeClassifier()  
2 x_train,x_test,y_train,y_test=train_test_split(x_scaled,y,test_size=0.2)  
3  
4 clf.fit(x_train,y_train);
```

evaluating the model

```
In [187]: 1 accuracy_score=clf.score(x_test,y_test)  
2 print('accuracy score',accuracy_score )
```

accuracy score 0.9666666666666667

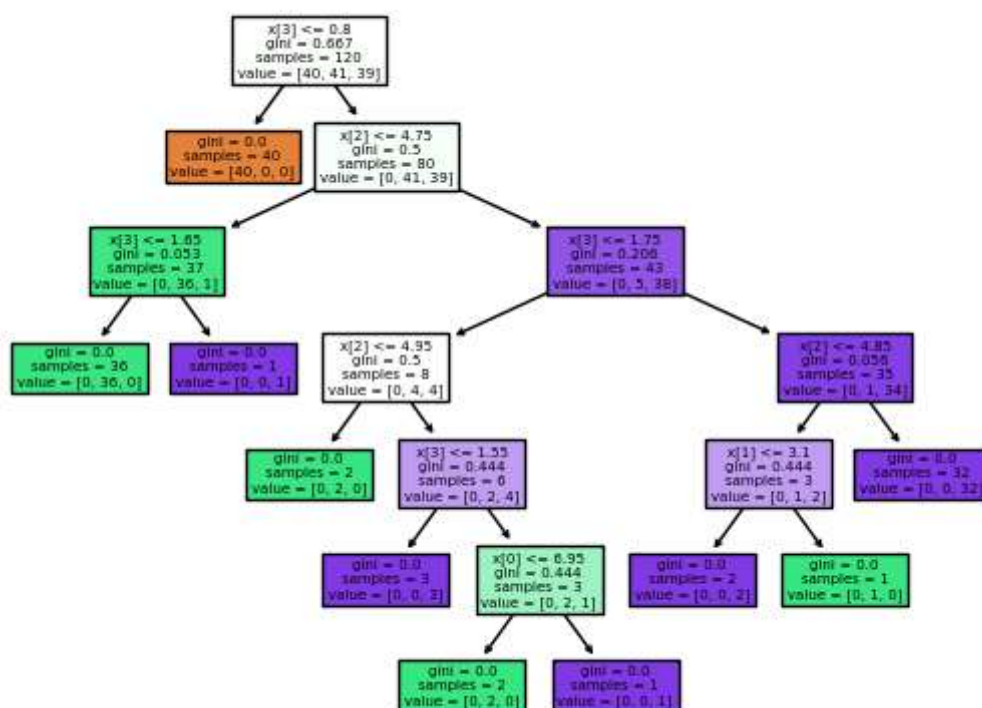
making the pridictions

```
In [188]: 1 y_preds=clf.predict(x_test)  
2 y_preds
```

```
Out[188]: array(['setosa', 'versicolor', 'setosa', 'versicolor', 'setosa', 'setosa',  
                'setosa', 'virginica', 'versicolor', 'setosa', 'versicolor',  
                'virginica', 'versicolor', 'versicolor', 'versicolor', 'setosa',  
                'setosa', 'versicolor', 'setosa', 'setosa', 'virginica', 'setosa',  
                'versicolor', 'virginica', 'setosa', 'versicolor', 'setosa',  
                'setosa', 'setosa', 'virginica'], dtype=object)
```

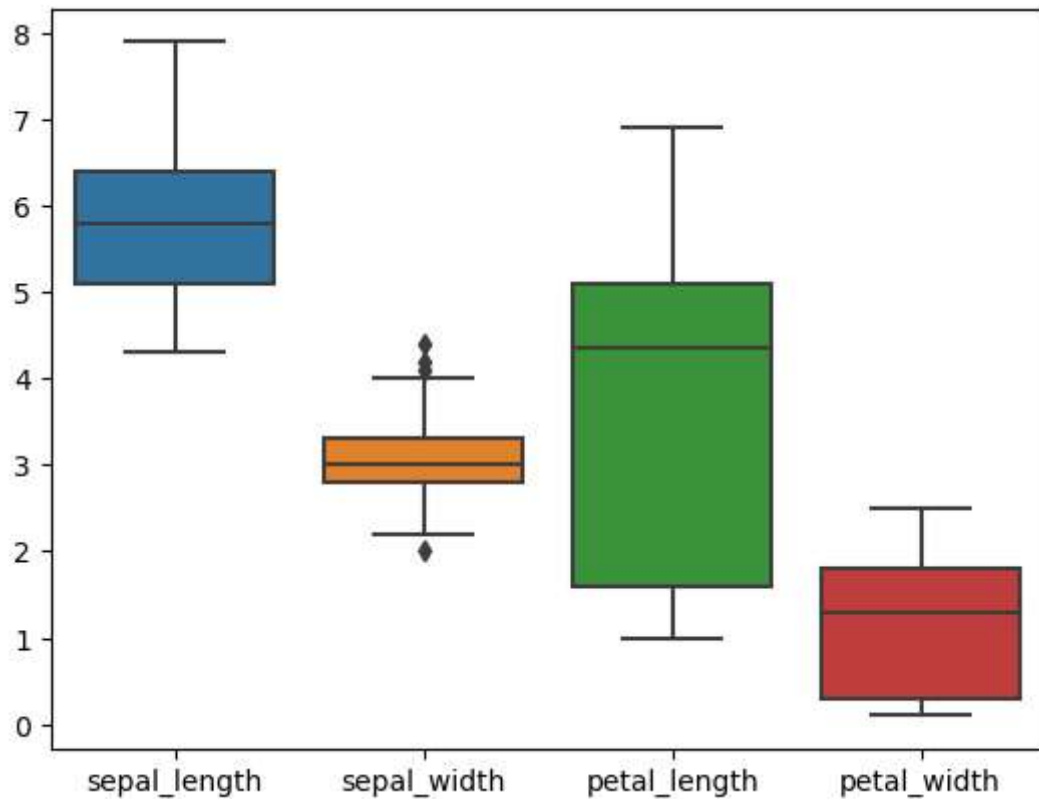
Plotting the tree at large

```
In [88]: 1 plot_tree(clf,filled=True);
          2
```



Identify outliers using a box plot

```
In [100]: 1 sns.boxplot(data=iris,saturation=0.75,width=0.8);  
2
```



Entropy

```
In [77]: 1 total_species=len(y)  
2 class_counts=y.value_counts()  
3 entropy=0  
4 for count in class_counts:  
5     probability=count/total_species  
6     entropy=entropy-probability*math.log2(probability)  
7 print(f'Entropy: {entropy}')
```

Entropy: 1.584962500721156

Gini Calculation

```
In [172]: 1 import numpy as np
          2
          3 def gini_index(labels):
          4     classes, count = np.unique(labels, return_counts=True)
          5
          6     prob = count / len(labels)
          7     gini = 1 - np.sum(prob ** 2)
          8     return gini
          9
         10 gini = gini_index(y)
         11 print("Gini Index:", gini)
```

Gini Index: 0.6666666666666667

```
In [ ]: 1
```