# Task 6- Exploratory Data Analysis - Prediction using Decision tree algorithm

Exploratory Data Analysis - Prediction using Decision tree algorithm

## Author:Aniket Suresh Hendre

Importing the required Libraries

```
In [1]: %matplotlib inline

import matplotlib.pyplot as plt
%matplotlib inline
import pandas as pd

from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier

from sklearn import tree
```

```
In [2]: #Importing the dataset
df=pd.read_csv('Iris.csv')
df
```

Out[2]:

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|---|---|---|---|---|---|
| 0 | 1 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 2 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 3 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |
| ... | ... | ... | ... | ... | ... | ... |
| 145 | 146 | 6.7 | 3.0 | 5.2 | 2.3 | Iris-virginica |
| 146 | 147 | 6.3 | 2.5 | 5.0 | 1.9 | Iris-virginica |
| 147 | 148 | 6.5 | 3.0 | 5.2 | 2.0 | Iris-virginica |
| 148 | 149 | 6.2 | 3.4 | 5.4 | 2.3 | Iris-virginica |
| 149 | 150 | 5.9 | 3.0 | 5.1 | 1.8 | Iris-virginica |

150 rows × 6 columns

```
In [3]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 150 entries, 0 to 149
Data columns (total 6 columns):
 #   Column         Non-Null Count  Dtype
---  ------         --------------  -----
 0   Id             150 non-null    int64
 1   SepalLengthCm  150 non-null    float64
 2   SepalWidthCm   150 non-null    float64
 3   PetalLengthCm  150 non-null    float64
 4   PetalWidthCm   150 non-null    float64
 5   Species        150 non-null    object
dtypes: float64(4), int64(1), object(1)
memory usage: 7.2+ KB
```

```
In [4]: df.head()
```

Out[4]:

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|---|---|---|---|---|---|
| 0 | 1 | 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 1 | 2 | 4.9 | 3.0 | 1.4 | 0.2 | Iris-setosa |
| 2 | 3 | 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 3 | 4 | 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 4 | 5 | 5.0 | 3.6 | 1.4 | 0.2 | Iris-setosa |

```
In [5]: df.tail()
```

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm | Species |
|---|---|---|---|---|---|---|
| 145 | 146 | 6.7 | 3.0 | 5.2 | 2.3 | Iris-virginica |
| 146 | 147 | 6.3 | 2.5 | 5.0 | 1.9 | Iris-virginica |
| 147 | 148 | 6.5 | 3.0 | 5.2 | 2.0 | Iris-virginica |
| 148 | 149 | 6.2 | 3.4 | 5.4 | 2.3 | Iris-virginica |
| 149 | 150 | 5.9 | 3.0 | 5.1 | 1.8 | Iris-virginica |

In [6]:
```python
df.describe()
```

Out[6]:

| | Id | SepalLengthCm | SepalWidthCm | PetalLengthCm | PetalWidthCm |
|---|---|---|---|---|---|
| count | 150.000000 | 150.000000 | 150.000000 | 150.000000 | 150.000000 |
| mean | 75.500000 | 5.843333 | 3.054000 | 3.758667 | 1.198667 |
| std | 43.445368 | 0.828066 | 0.433594 | 1.764420 | 0.763161 |
| min | 1.000000 | 4.300000 | 2.000000 | 1.000000 | 0.100000 |
| 25% | 38.250000 | 5.100000 | 2.800000 | 1.600000 | 0.300000 |
| 50% | 75.500000 | 5.800000 | 3.000000 | 4.350000 | 1.300000 |
| 75% | 112.750000 | 6.400000 | 3.300000 | 5.100000 | 1.800000 |
| max | 150.000000 | 7.900000 | 4.400000 | 6.900000 | 2.500000 |

## Checking out the null Values in the dataset

In [7]:
```python
df.isnull().sum()
```

Out[7]:
```
Id               0
SepalLengthCm    0
SepalWidthCm     0
PetalLengthCm    0
PetalWidthCm     0
Species          0
dtype: int64
```

## Training the Model

In [30]:
```python
features = ['SepalLengthCm','SepalWidthCm','PetalLengthCm','PetalWidthCm']

# Create features matrix
x = df.loc[:, features].values
```

In [31]:
```python
y=df.Species
x_train,x_test,y_train,y_test=train_test_split(x, y, random_state=0)
```

In [32]:
```python
clf = DecisionTreeClassifier(max_depth = 2,
                             random_state = 0)
```

In [33]:
```python
clf.fit(x_train, y_train)
```

Out[33]:
```
DecisionTreeClassifier(max_depth=2, random_state=0)
```

In [34]:
```python
best_clf = clf
```

In [35]:
```python
clf.predict(x_test[0:1])
```

Out[35]:
```
array(['Iris-virginica'], dtype=object)
```

In [36]:
```python
from sklearn import metrics
import seaborn as sns
```

In [37]:
```python
score = clf.score(x_test, y_test)
print(score)
```
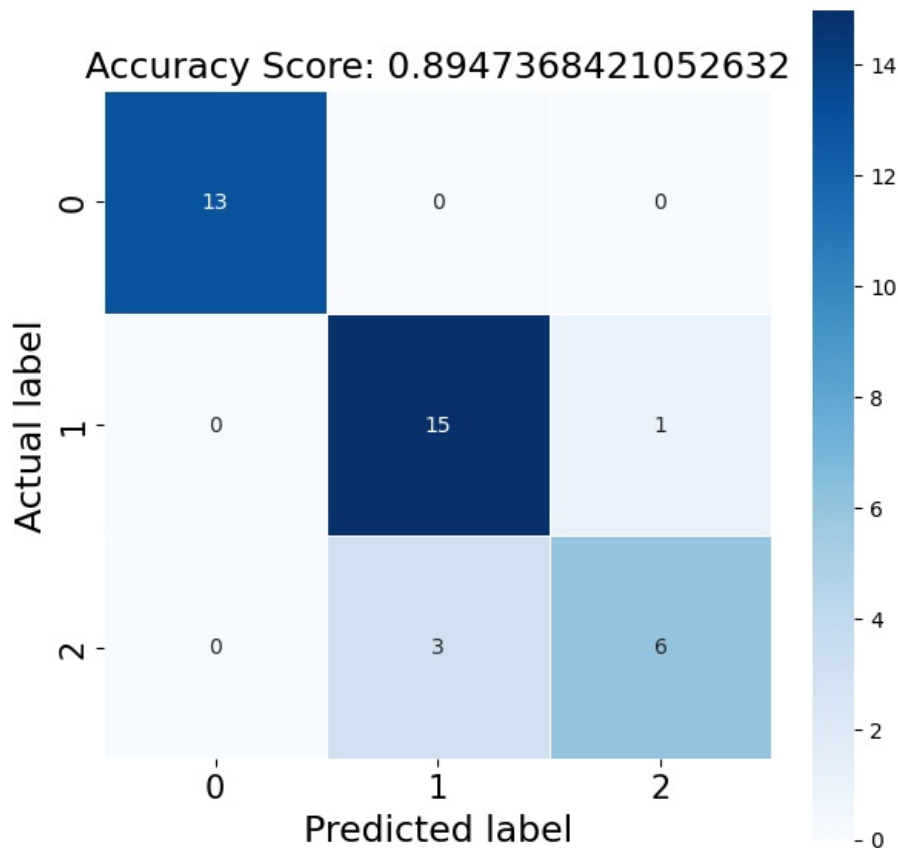```
0.8947368421052632
```

In [38]:
```python
print(metrics.classification_report(y_test,clf.predict(x_test)))
```

|              | precision | recall | f1-score | support |
|--------------|-----------|--------|----------|---------|
| Iris-setosa     | 1.00 | 1.00 | 1.00 | 13 |
| Iris-versicolor | 0.83 | 0.94 | 0.88 | 16 |
| Iris-virginica  | 0.86 | 0.67 | 0.75 | 9 |
|              |      |      |      |    |
| accuracy     |      |      | 0.89 | 38 |
| macro avg    | 0.90 | 0.87 | 0.88 | 38 |
| weighted avg | 0.90 | 0.89 | 0.89 | 38 |

## Visualing the data

```
In [39]:  cm = metrics.confusion_matrix(y_test, clf.predict(x_test))

          plt.figure(figsize=(7,7))
          sns.heatmap(cm, annot=True,
                      fmt=".0f",
                      linewidths=.5,
                      square = True,
                      cmap = 'Blues');
          plt.ylabel('Actual label', fontsize = 17);
          plt.xlabel('Predicted label', fontsize = 17);
          plt.title('Accuracy Score: {}'.format(score), size = 17);
          plt.tick_params(labelsize= 15)
```



```
In [40]:  # List of values to try for max_depth:
          max_depth_range = list(range(1, 6))

          # List to store the average RMSE for each value of max_depth:
          accuracy = []

          for depth in max_depth_range:

              clf = DecisionTreeClassifier(max_depth = depth,
                                           random_state = 0)
              clf.fit(x_train, y_train)

              score = clf.score(x_test, y_test)
              accuracy.append(score)
```
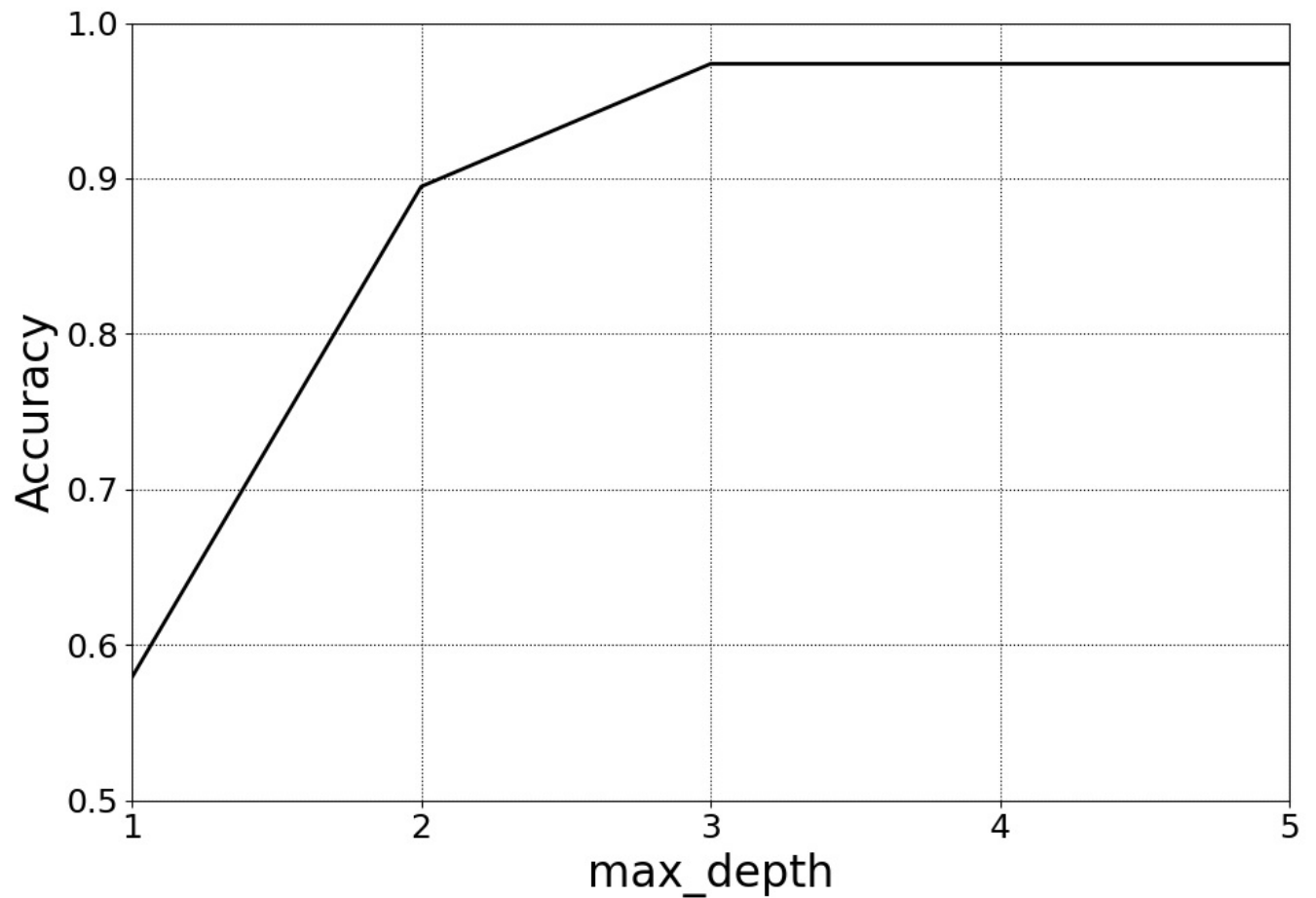
```
In [41]:  #ploting accuracy score depth wise
          fig, ax = plt.subplots(nrows = 1, ncols = 1, figsize = (10,7));

          ax.plot(max_depth_range,
                  accuracy,
                  lw=2,
                  color='k')

          ax.set_xlim([1, 5])
          ax.set_ylim([.50, 1.00])
```
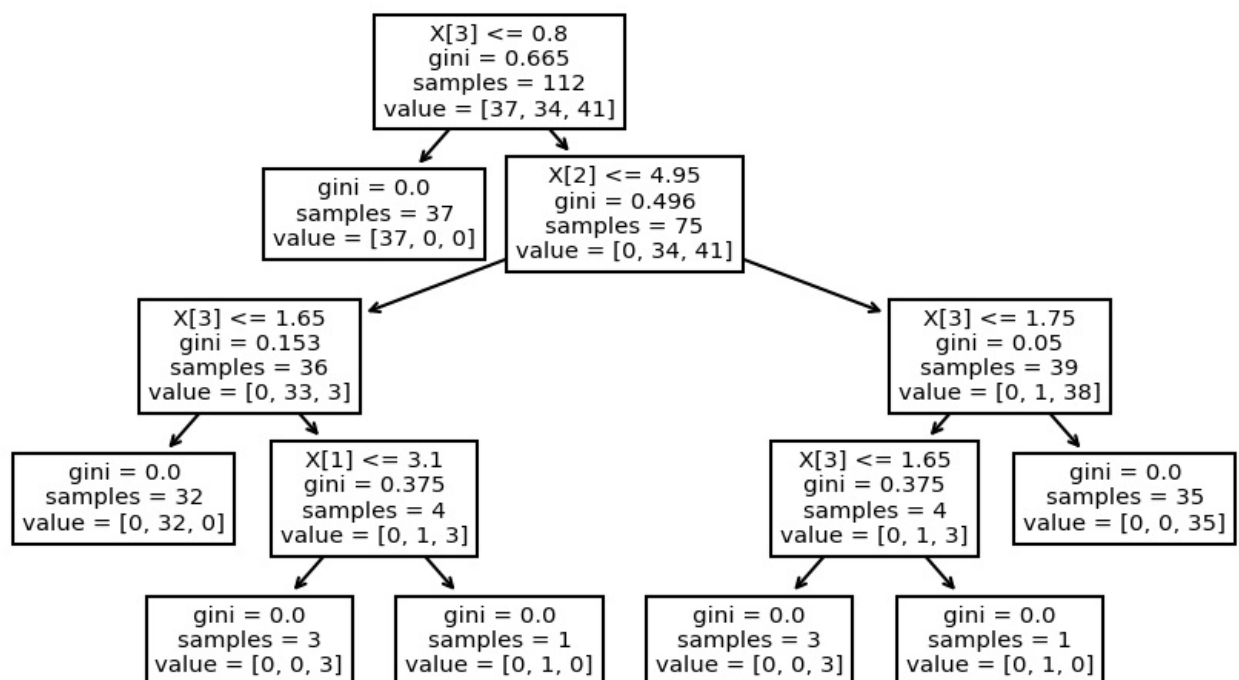
```
ax.grid(True,
        axis = 'both',
        zorder = 0,
        linestyle = ':',
        color = 'k')


ax.tick_params(labelsize = 18)
ax.set_xticks([1,2,3,4,5])
ax.set_xlabel('max_depth', fontsize = 24)
ax.set_ylabel('Accuracy', fontsize = 24)
fig.tight_layout()
#fig.savefig('images/max_depth_vs_accuracy.png', dpi = 300)
```
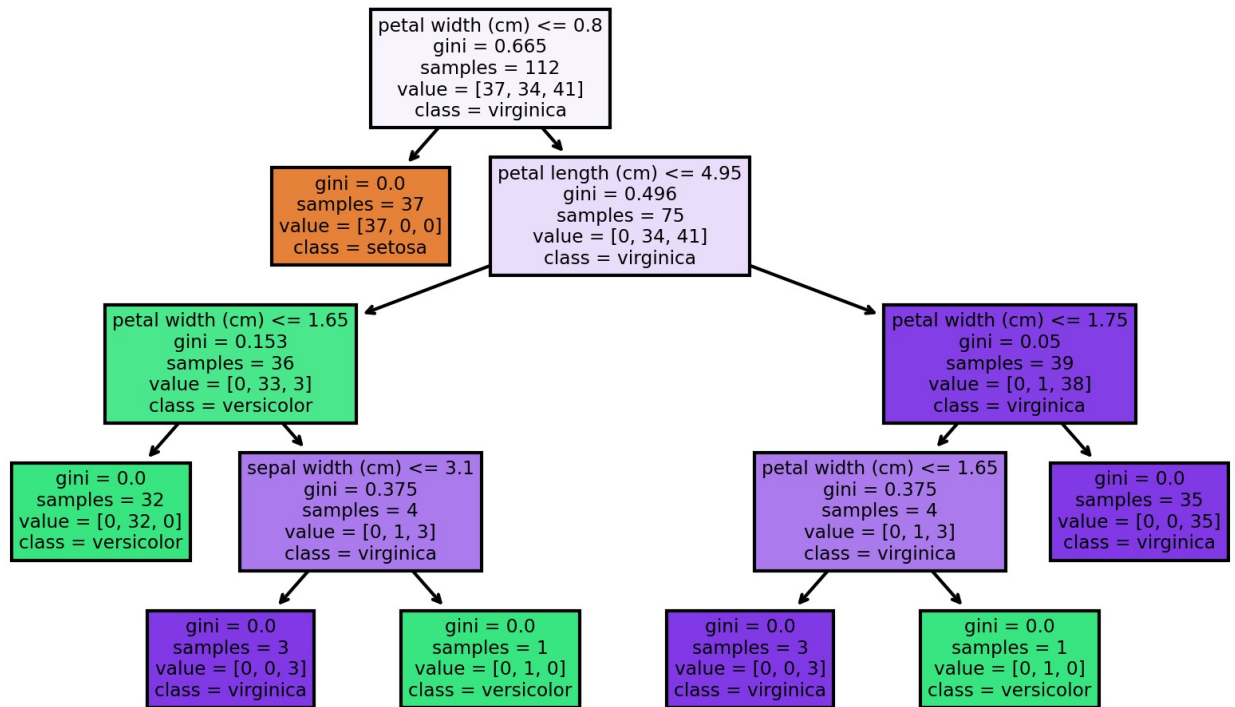


In [42]:
```
fig, axes = plt.subplots(nrows = 1, ncols = 1, figsize = (7,4), dpi = 150)

tree.plot_tree(clf);
```

```python
# Putting the feature names and class names into variables
fn = ['sepal length (cm)','sepal width (cm)','petal length (cm)','petal width (cm)']
cn = ['setosa', 'versicolor', 'virginica']
```

In [44]:
```python
fig, axes = plt.subplots(nrows = 1, ncols = 1, figsize = (7,4), dpi = 300)

tree.plot_tree(clf,
               feature_names = fn,
               class_names=cn,
               filled = True);
#fig.savefig('images/plottreefncn.png')
```



## Conclusion :

- After Importing, Fit our dataset in our model, accuracy is 89.47%.

- We can clearly see model performance by confusion matrix and classification report.

- By ploting accuracy score depth wise graph, optimal depth for model is 3.

In [ ]:

Loading [MathJax]/jax/output/CommonHTML/fonts/TeX/fontdata.js