

Tanvi Bhosle

Data Science And Business Analytics Intern @ TSF GRIP JULY2021

Task 6: Prdiction using Decision Tree Algorithm

Dataset: <https://bit.ly/3kXTdox>

```
In [ ]: 
```

Import Libraries¶

```
In [103]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn import metrics
import graphviz
from sklearn import tree
```

Import Dataset

```
In [8]: data=pd.read_csv("C:\\Users\\Admin\\Desktop\\Iris.csv",encoding=("ISO-8859-1"),low_memory=False)
data.head(5)
```

```
Out[8]:
```

	Id	SepalLengthCm	SepalWidthCm	Petal.LengthCm	Petal.WidthCm	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 [65]: data.shape

<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   Petal.LengthCm     150 non-null   float64
 4   Petal.WidthCm      150 non-null   float64
 5   Species            150 non-null   object
dtypes: float64(4), int64(1), object(1)
memory usage: 7.2+ KB
```

```
In [11]: data["Species"].value_counts()
```

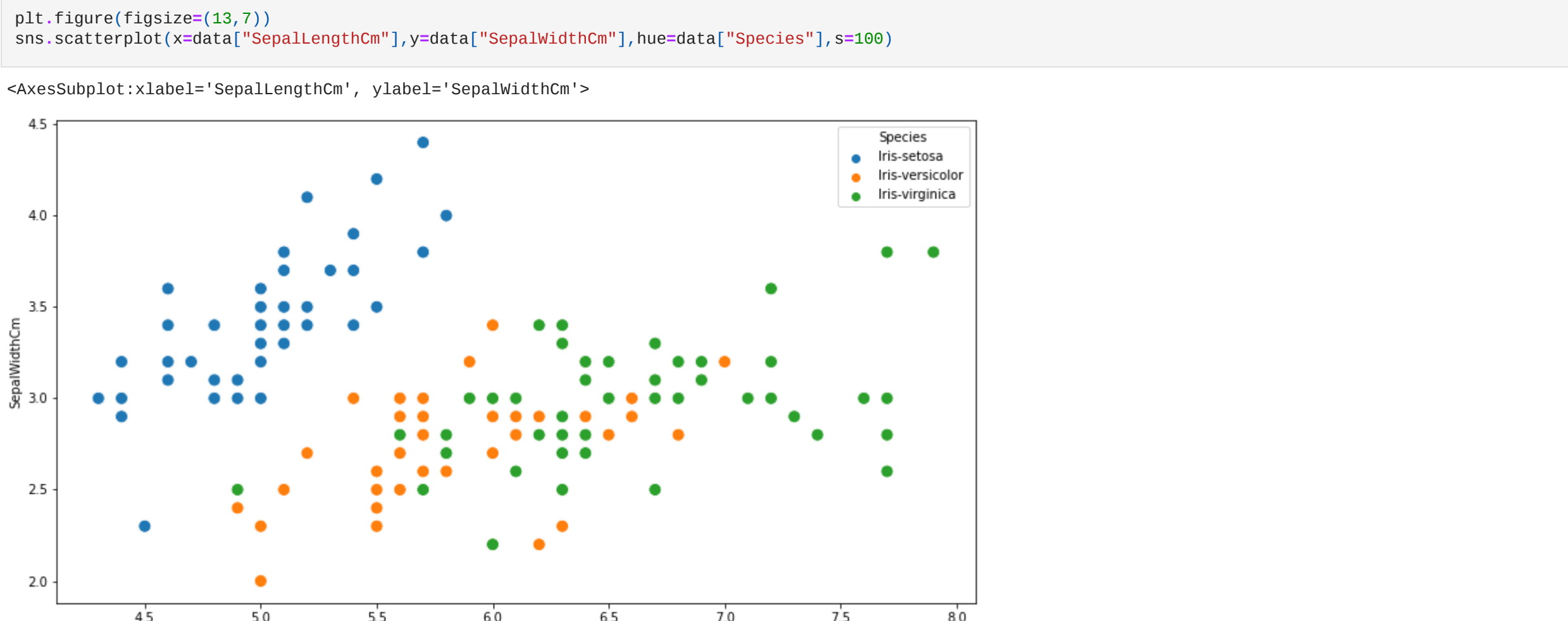
```
Out[11]: Iris-setosa      50
Iris-versicolor      50
Iris-virginica       50
Name: Species, dtype: int64
```

```
In [12]: data.isnull().sum()
```

```
Out[12]: Id                0
SepalLengthCm            0
SepalWidthCm             0
Petal.LengthCm           0
Petal.WidthCm            0
Species                 0
dtype: int64
```

Sepal Dimensions

```
In [15]: plt.figure(figsize=(13,7))
sns.scatterplot(x=data["SepalLengthCm"],y=data["SepalWidthCm"],hue=data["Species"],s=100)
```



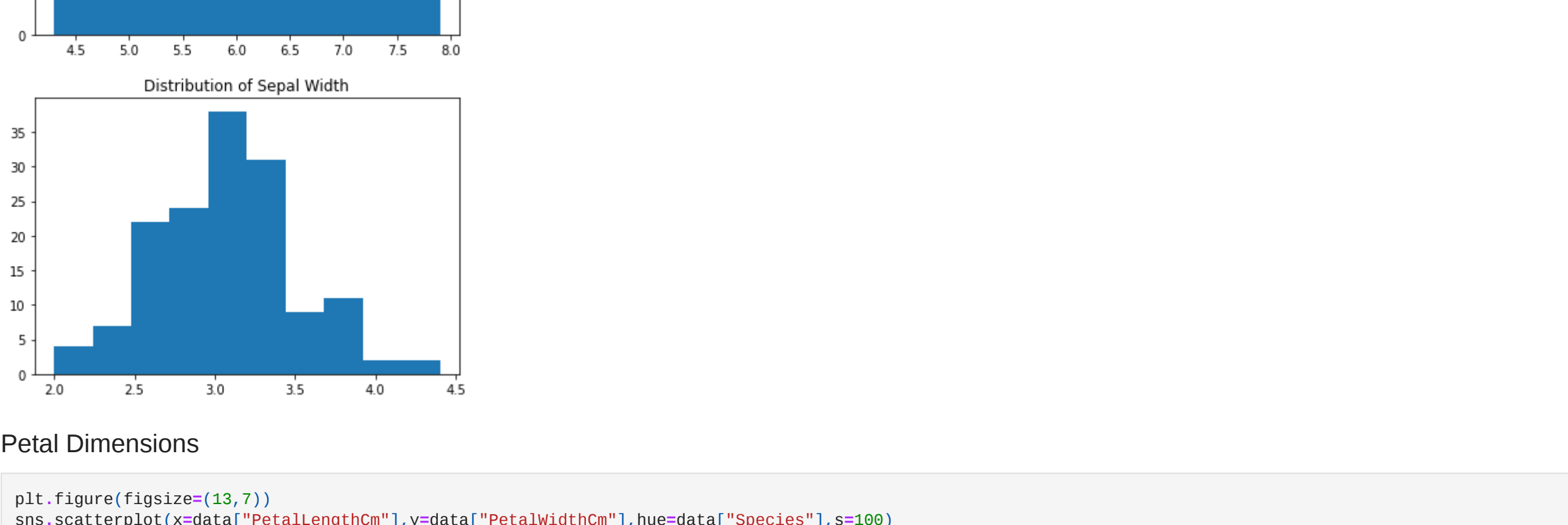
```
In [28]: fig,axes=plt.subplots(figsize=(6,4))
plt.title("Distribution of Sepal Length")
plt.hist(data["SepalLengthCm"])
fig,axes=plt.subplots(figsize=(6,4))
plt.title("Distribution of Sepal Width")
plt.hist(data["SepalWidthCm"])
```

```
Out[28]: (array([ 4.,  7., 22., 24., 38., 31.,  9., 11.,  2.,  2.]),
array([ 2.,  2.24, 2.48, 2.72, 2.96,  3.2,  3.44, 3.68, 3.92, 4.16, 4.4 ]),
<BarContainer object of 10 artists>)
```



Petal Dimensions

```
In [29]: plt.figure(figsize=(13,7))
sns.scatterplot(x=data["Petal.LengthCm"],y=data["Petal.WidthCm"],hue=data["Species"],s=100)
```



```
In [30]: fig,axes=plt.subplots(figsize=(6,4))
plt.title("Distribution of Petal Length")
plt.hist(data["Petal.LengthCm"])
fig,axes=plt.subplots(figsize=(6,4))
plt.title("Distribution of Petal Width")
plt.hist(data["Petal.WidthCm"])
```

```
Out[30]: (array([ 4.1,  8.,  1.,  7.,  8., 33.,  6., 23.,  9., 14.]),
array([ 0.1,  0.34, 0.58, 0.82, 1.06, 1.3,  1.54, 1.78, 2.02, 2.26, 2.5 ]),
<BarContainer object of 10 artists>)
```



Correlation between Dependent and Independent variables

```
In [70]: data=data.drop("Id",axis=1)
data.head()
data["Species"].replace("Iris-setosa","0",inplace=True)
data["Species"].replace("Iris-versicolor","1",inplace=True)
data["Species"].replace("Iris-virginica","2",inplace=True)
data["Species"]=data.Species.astype(int)
data.head()
```

```
Out[70]:
```

	SepalLengthCm	SepalWidthCm	Petal.LengthCm	Petal.WidthCm	Species
0	5.1	3.5	1.4	0.2	0
1	4.9	3.0	1.4	0.2	0
2	4.7	3.2	1.3	0.2	0
3	4.6	3.1	1.5	0.2	0
4	5.0	3.6	1.4	0.2	0

```
In [73]: plt.figure(figsize=(8,4))
sns.heatmap(data1.corr(),annot=True)
```

Out[73]: <AxesSubplot: >



Model Building

```
In [75]: train,test=train_test_split(data1,test_size=0.3)
```

```
In [76]: train.shape,test.shape
```

```
Out[76]: ((105, 5), (45, 5))
```

```
In [79]: train_x=train[["SepalLengthCm","SepalWidthCm","Petal.LengthCm","Petal.WidthCm"]]
train_y=train.Species
test_x=test[["SepalLengthCm","SepalWidthCm","Petal.LengthCm","Petal.WidthCm"]]
test_y=test.Species
```

```
In [83]: dtree=DecisionTreeClassifier()
dtree.fit(train_x,train_y)
predictions=dtree.predict(test_x)
accuracy=metrics.accuracy_score(predictions,test_y)
accuracy
```

```
Out[83]: 0.9555555555555556
```

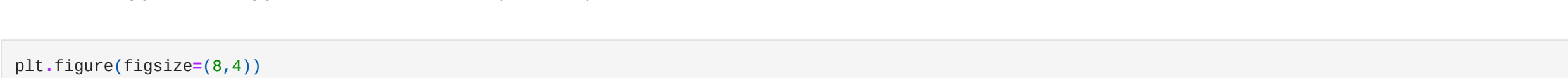
```
In [85]: x=data1[["SepalLengthCm","SepalWidthCm","Petal.LengthCm","Petal.WidthCm"]]
y=data1.Species
dtree1=DecisionTreeClassifier()
dtree1.fit(x,y)
```

```
Out[85]: DecisionTreeClassifier()
```

Visualization

```
In [109]: fig=plt.figure(figsize=(25,20))
tree=plt_tree(dtree1,feature_names=["SepalLengthCm","SepalWidthCm","Petal.LengthCm","Petal.WidthCm"],class_names=["0","1","2"],
filled=True)
```

```
Out[109]: [Text(697.5, 996.6, 'Petal.LengthCm <= 2.45\n gini = 0.667\n samples = 150\n value = [50, 50, 50]\n class = 0'),
Text(590.1923076923077, 815.4080000000001, 'gini = 0.0\n samples = 50\n value = [50, 0, 0]\n class = 0'),
Text(804.8076923076923, 815.4080000000001, 'Petal.WidthCm <= 1.75\n gini = 0.5\n samples = 100\n value = [0, 50, 50]\n class = 1'),
Text(429.2307692307692, 634.2, 'Petal.LengthCm <= 4.95\n gini = 0.168\n samples = 54\n value = [0, 49, 5]\n class = 1'),
Text(214.6153846153846, 453.0, 'Petal.WidthCm <= 1.65\n gini = 0.0\n samples = 48\n value = [0, 47, 1]\n class = 1'),
Text(107.3076923076923, 271.79999999999995, 'gini = 0.0\n samples = 47\n value = [0, 47, 0]\n class = 1'),
Text(321.9230769230769, 271.79999999999995, 'gini = 0.0\n samples = 1\n value = [0, 0, 1]\n class = 2'),
Text(643.8461538461538, 453.0, 'Petal.WidthCm <= 1.55\n gini = 0.444\n samples = 6\n value = [0, 2, 0]\n class = 2'),
Text(536.5384615384615, 271.79999999999995, 'gini = 0.0\n samples = 3\n value = [0, 0, 3]\n class = 2'),
Text(751.5384615384615, 271.79999999999995, 'Sepal.LengthCm <= 6.95\n gini = 0.444\n samples = 3\n value = [0, 1, 2]\n class = 1'),
Text(643.8461538461538, 90.59999999999999, 'gini = 0.0\n samples = 2\n value = [0, 2, 0]\n class = 1'),
Text(858.4615384615385, 90.59999999999999, 'gini = 0.0\n samples = 1\n value = [0, 0, 1]\n class = 2'),
Text(1180.3846153846155, 634.2, 'Petal.LengthCm <= 4.85\n gini = 0.043\n samples = 46\n value = [0, 1, 45]\n class = 2'),
Text(1073.676923076923, 453.0, 'Sepal.WidthCm <= 3.1\n gini = 0.444\n samples = 3\n value = [0, 1, 2]\n class = 2'),
Text(965.7692307692307, 271.79999999999995, 'gini = 0.0\n samples = 2\n value = [0, 0, 2]\n class = 2'),
Text(1180.3846153846155, 271.79999999999995, 'gini = 0.0\n samples = 1\n value = [0, 1, 0]\n class = 1'),
Text(1287.6923076923076, 453.0, 'gini = 0.0\n samples = 43\n value = [0, 0, 43]\n class = 2')]
```



```
In [116]: dtree.predict([[6.8,2.5,1.5]])
```

```
Out[116]: array([2])
```

```
In [118]: dtree.predict([[2,2,2,2]])
```

```
Out[118]: array([0])
```

```
In [124]: dtree.predict([[0.1,2.5,0.9,2]])
```

```
Out[124]: array([2])
```

```
In [125]: dtree.predict([[5,3.4,1.5,0.2]])
```

```
Out[125]: array([0])
```

```
In [126]: dtree.predict([[5.7,2.8,4.5,1.3]])
```

```
Out[126]: array([1])
```

```
In [ ]: 
```