# DS - MAJOR- SEPTEMBER-DS-09-MLB 5

## PYTHON MAJOR PROJECT

*Q.TAKE ANY DATASET OF YOUR CHOICE, PERFORM EDA(EXPLORATORY DATA ANALYSIS) AND APPLY A SUITABLE CLASSIFIER, REGRESSOR, OR CLUSTER AND CALCULATE THE ACCURACY OF THE MODEL.*

## NAME-AKSHITA SHARMA (SEPTEMBER-OCTOBER BATCH)

## IRIS DATASET

It is a flower. It contains five columns namely – Petal Length, Petal Width, Sepal Length, Sepal Width, and Species Type. Iris is a flowering plant, and many researchers have measured various features of the different iris flowers and recorded them digitally.

**FOLLOWING IS THE DATASET OF IRIS-**

| sepallength | sepalwidth | petallength | petalwidth | class |
|---|---|---|---|---|
| 5.1 | 3.5 | 1.4 | 0.2 | Iris-setosa |
| 4.9 | 3 | 1.4 | 0.2 | Iris-setosa |
| 4.7 | 3.2 | 1.3 | 0.2 | Iris-setosa |
| 4.6 | 3.1 | 1.5 | 0.2 | Iris-setosa |
| 5 | 3.6 | 1.4 | 0.2 | Iris-setosa |
| 5.4 | 3.9 | 1.7 | 0.4 | Iris-setosa |
| 4.6 | 3.4 | 1.4 | 0.3 | Iris-setosa |
| 5 | 3.4 | 1.5 | 0.2 | Iris-setosa |
| 4.4 | 2.9 | 1.4 | 0.2 | Iris-setosa |
| 4.9 | 3.1 | 1.5 | 0.1 | Iris-setosa |
| 5.4 | 3.7 | 1.5 | 0.2 | Iris-setosa |
| 4.8 | 3.4 | 1.6 | 0.2 | Iris-setosa |
| 4.8 | 3 | 1.4 | 0.1 | Iris-setosa |
| 4.3 | 3 | 1.1 | 0.1 | Iris-setosa |
| 5.8 | 4 | 1.2 | 0.2 | Iris-setosa |
| 5.7 | 4.4 | 1.5 | 0.4 | Iris-setosa |
| 5.1 | 3.5 | 1.4 | 0.3 | Iris-setosa |

NOW PERFORMING EDA (Exploratory Data Analysis)ON THE GIVEN DATASET-
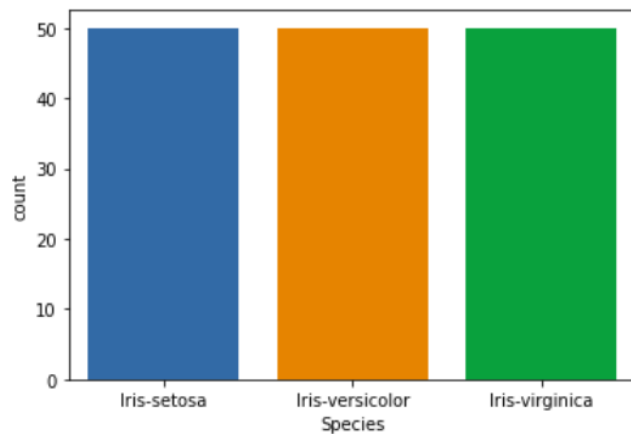
**Visualizing the target column:**

Iris-setosa ,iris versicolor species, iris virginica

**CODE:**

```python
# importing packages
import seaborn as sns
import matplotlib.pyplot as plt


sns.countplot(x='Species', data=df, )
plt.show()
```

## OUTPUT:



**CLUSTERING ON THE DATASET:**

**Relation between variables:**

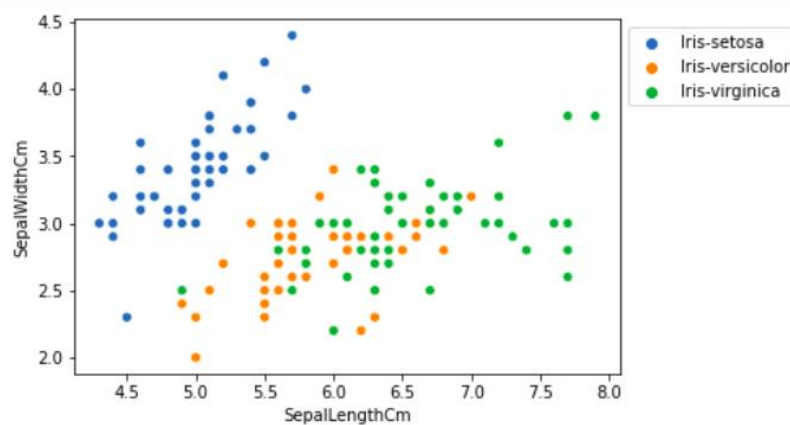Comparing Sepal Length and Sepal Width-

## CODE-

```python
# importing packages
import seaborn as sns
import matplotlib.pyplot as plt


sns.scatterplot(x='SepalLengthCm', y='SepalWidthCm',
                hue='Species', data=df, )

# Placing Legend outside the Figure
plt.legend(bbox_to_anchor=(1, 1), loc=2)

plt.show()
```

## OUTPUT-



From the above plot, we can infer that –

- Species Setosa has smaller sepal lengths but larger sepal widths.

- Versicolor Species lies in the middle of the other two species in terms of sepal length and width.

- Species Virginica has larger sepal lengths but smaller sepal widths.

Comparing Petal Length and Petal Width-
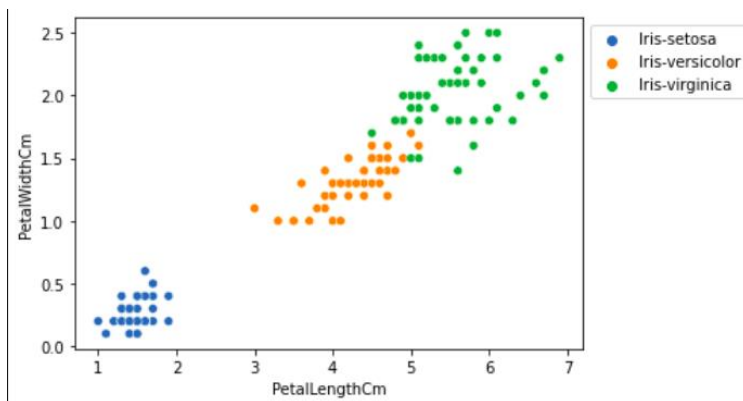
# CODE-

```
# importing packages
import seaborn as sns
import matplotlib.pyplot as plt


sns.scatterplot(x='PetalLengthCm', y='PetalWidthCm',
                hue='Species', data=df, )

# Placing Legend outside the Figure
plt.legend(bbox_to_anchor=(1, 1), loc=2)

plt.show()
```

# OUTPUT-



From the above plot, we can infer that –

- Species Setosa has smaller petal lengths and widths.

- Versicolor Species lies in the middle of the other two species in terms of petal length and width

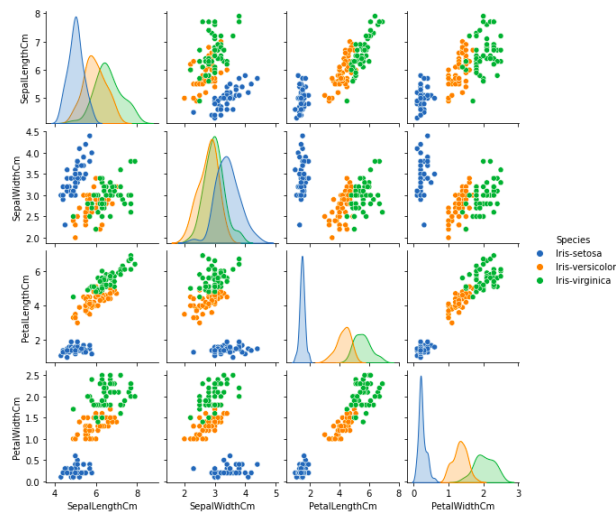- Species Virginica has the largest of petal lengths and widths.

Multivariate Analysis-

# CODE-

```
# importing packages
import seaborn as sns
import matplotlib.pyplot as plt


sns.pairplot(df.drop(['Id'], axis = 1),
             hue='Species', height=2)
```

## OUTPUT-



 As the species Seotsa has the smallest of petals widths and lengths. It also has the smallest sepal length but larger sepal widths. Such information can be gathered about any other species.

## CLASSIFIER :

## Uni as well as Bi-variate analysis-

## CODE:

```python
# importing packages
import seaborn as sns
import matplotlib.pyplot as plt


fig, axes = plt.subplots(2, 2, figsize=(10,10))

axes[0,0].set_title("Sepal Length")
axes[0,0].hist(df['SepalLengthCm'], bins=7)

axes[0,1].set_title("Sepal Width")
axes[0,1].hist(df['SepalWidthCm'], bins=5);

axes[1,0].set_title("Petal Length")
axes[1,0].hist(df['PetalLengthCm'], bins=6);

axes[1,1].set_title("Petal Width")
axes[1,1].hist(df['PetalWidthCm'], bins=6);
```
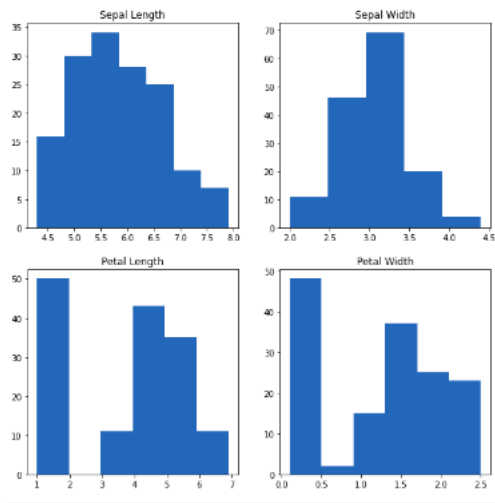
## OUTPUT:

From the above plot, we can see that –

- The highest frequency of the sepal length is between 30 and 35 which is between 5.5 and 6

- The highest frequency of the sepal Width is around 70 which is between 3.0 and 3.5

- The highest frequency of the petal length is around 50 which is between 1 and 2

- The highest frequency of the petal width is between 40 and 50 which is between 0.0 and 0.5

**The univariant set of observations and visualizing it through a histogram-**

## CODE-

```
# importing packages
import seaborn as sns
import matplotlib.pyplot as plt

plot = sns.FacetGrid(df, hue="Species")
plot.map(sns.distplot, "SepalLengthCm").add_legend()

plot = sns.FacetGrid(df, hue="Species")
plot.map(sns.distplot, "SepalWidthCm").add_legend()

plot = sns.FacetGrid(df, hue="Species")
plot.map(sns.distplot, "PetalLengthCm").add_legend()

plot = sns.FacetGrid(df, hue="Species")
plot.map(sns.distplot, "PetalWidthCm").add_legend()

plt.show()
```
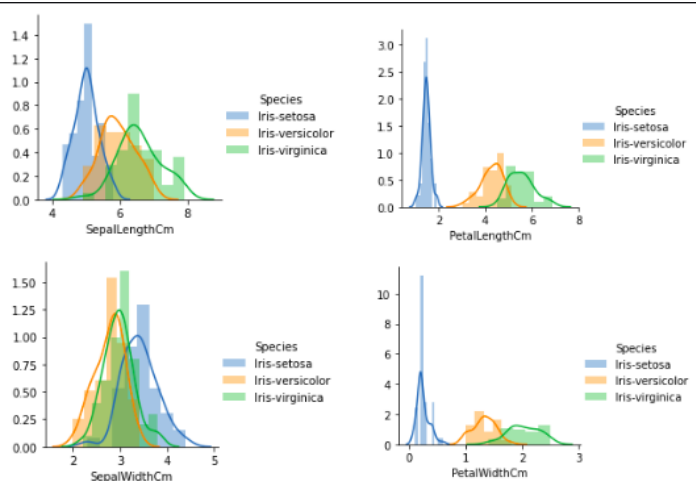
## OUTPUT-

From the above plots, we can see that –

- In the case of Sepal Length, there is a huge amount of overlapping.

- In the case of Sepal Width also, there is a huge amount of overlapping.

- In the case of Petal Length, there is a very little amount of overlapping.

- In the case of Petal Width also, there is a very little amount of overlapping.

**So we can use Petal Length and Petal Width as the classification feature.**

## REGRESSOR:

## CODE-

```python
# importing packages
import seaborn as sns
import matplotlib.pyplot as plt

def graph(y):
    sns.boxplot(x="Species", y=y, data=df)

plt.figure(figsize=(10,10))

# Adding the subplot at the specified
# grid position
plt.subplot(221)
graph('SepalLengthCm')

plt.subplot(222)
graph('SepalWidthCm')

plt.subplot(223)
graph('PetalLengthCm')

plt.subplot(224)
graph('PetalWidthCm')

plt.show()
```
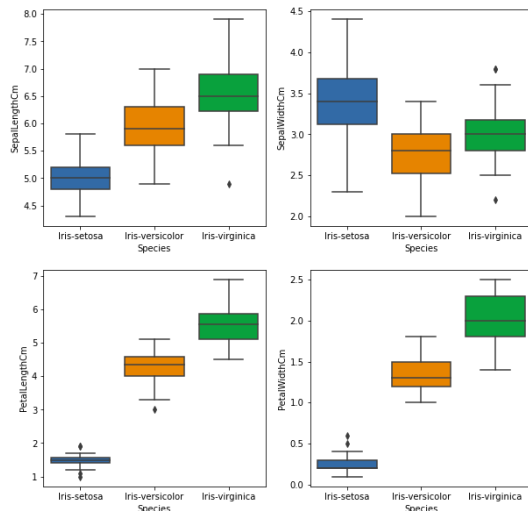
## OUTPUT-

From the above graph, we can see that –

- Species Setosa has the smallest features and less distributed with some outliers.

- Species Versicolor has the average features.

- Species Virginica has the highest feature

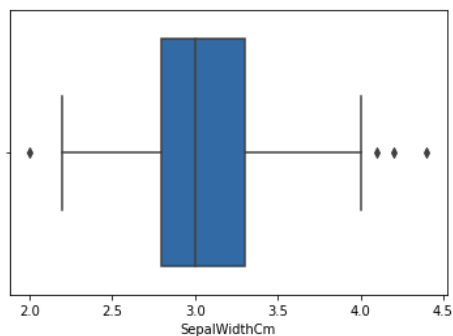**The iris dataset for the 'SepalWidthCm' column-**

## CODE-

```
# importing packages
import seaborn as sns
import matplotlib.pyplot as plt

# Load the dataset
df = pd.read_csv('Iris.csv')

sns.boxplot(x='SepalWidthCm', data=df)
```

## OUTPUT-



In the above graph, the values above 4 and below 2 are acting as outliers.

**Removing Outliers**

## CODE-

```
# Importing
import sklearn
from sklearn.datasets import load_boston
import pandas as pd
import seaborn as sns

# Load the dataset
df = pd.read_csv('Iris.csv')

# IQR
Q1 = np.percentile(df['SepalWidthCm'], 25,
                interpolation = 'midpoint')

Q3 = np.percentile(df['SepalWidthCm'], 75,
                interpolation = 'midpoint')
IQR = Q3 - Q1

print("Old Shape: ", df.shape)

# Upper bound
upper = np.where(df['SepalWidthCm'] >= (Q3+1.5*IQR))

# Lower bound
lower = np.where(df['SepalWidthCm'] <= (Q1-1.5*IQR))

# Removing the Outliers
df.drop(upper[0], inplace = True)
df.drop(lower[0], inplace = True)

print("New Shape: ", df.shape)
```
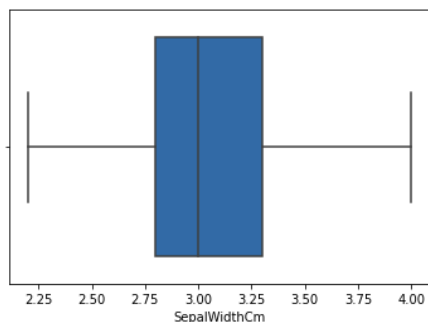
## OUTPUT-

```
Old Shape:  (150, 6)
New Shape:  (146, 6)

<AxesSubplot:xlabel='SepalWidthCm'>
```



**ACCURACY OF THE MODEL:**

**A correlation between all numerical variables in the dataset-**

## CODE-

```
# importing packages
import seaborn as sns
import matplotlib.pyplot as plt


sns.heatmap(df.corr(method='pearson').drop(
    ['Id'], axis=1).drop(['Id'], axis=0),
            annot = True);

plt.show()
```
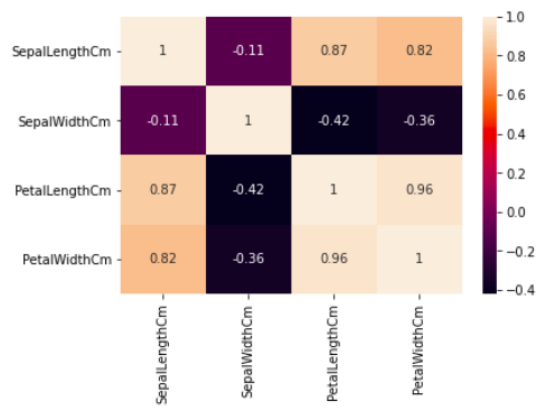
## OUTPUT-

From the above graph, we can see that –

- Petal width and petal length have high correlations.

- Petal length and sepal width have good correlations.

- Petal Width and Sepal length have good correlations.