

Practical No. 10

Name: Prathamesh Pawar | Roll No : B-23

Data Visualization III

Download the Iris flower dataset or any other dataset into a DataFrame. (e.g., <https://archive.ics.uci.edu/ml/datasets/Iris> (<https://archive.ics.uci.edu/ml/datasets/Iris>)). Scan the dataset and give the inference as:

1. List down the features and their types (e.g., numeric, nominal) available in the dataset.
2. Create a histogram for each feature in the dataset to illustrate the feature distributions.
3. Create a boxplot for each feature in the dataset.
4. Compare distributions and identify outliers.

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import seaborn as sns
import warnings
warnings.filterwarnings('ignore')
df = pd.read_csv('iris.csv')
df.head()
```

```
Out[1]:
```

	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	NaN	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 [14]: df.isnull().sum()
```

```
Out[14]: Id                0
SepalLengthCm            0
SepalWidthCm             0
PetalLengthCm            1
PetalWidthCm             0
Species                  0
dtype: int64
```

```
In [15]: df['PetalLengthCm']=df['PetalLengthCm'].fillna(np.mean(df['PetalLengthCm']))
```

```
In [16]: df.isnull().sum()
```

```
Out[16]: Id                0
SepalLengthCm            0
SepalWidthCm             0
PetalLengthCm            0
PetalWidthCm             0
Species                  0
dtype: int64
```

1. List down the features and their types (e.g., numeric, nominal) available in the dataset.

```
In [17]: 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
```

Hence the dataset contains 4 numerical columns and 1 object column

```
In [18]: np.unique(df["Species"])
```

```
Out[18]: array(['Iris-setosa', 'Iris-versicolor', 'Iris-virginica'], dtype=object)
```

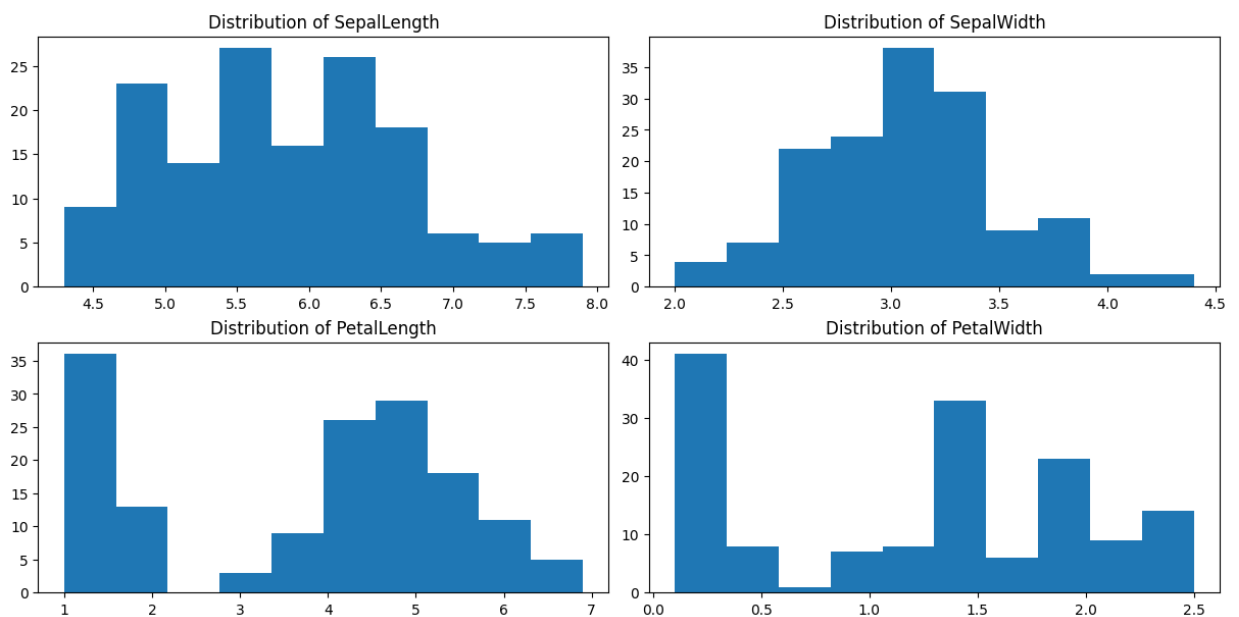
```
In [19]: df.describe()
```

```
Out[19]:
```

	Id	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	150.000000	150.000000	150.000000	150.000000	150.000000
mean	75.500000	5.843333	3.054000	3.775168	1.198667
std	43.445368	0.828066	0.433594	1.752808	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

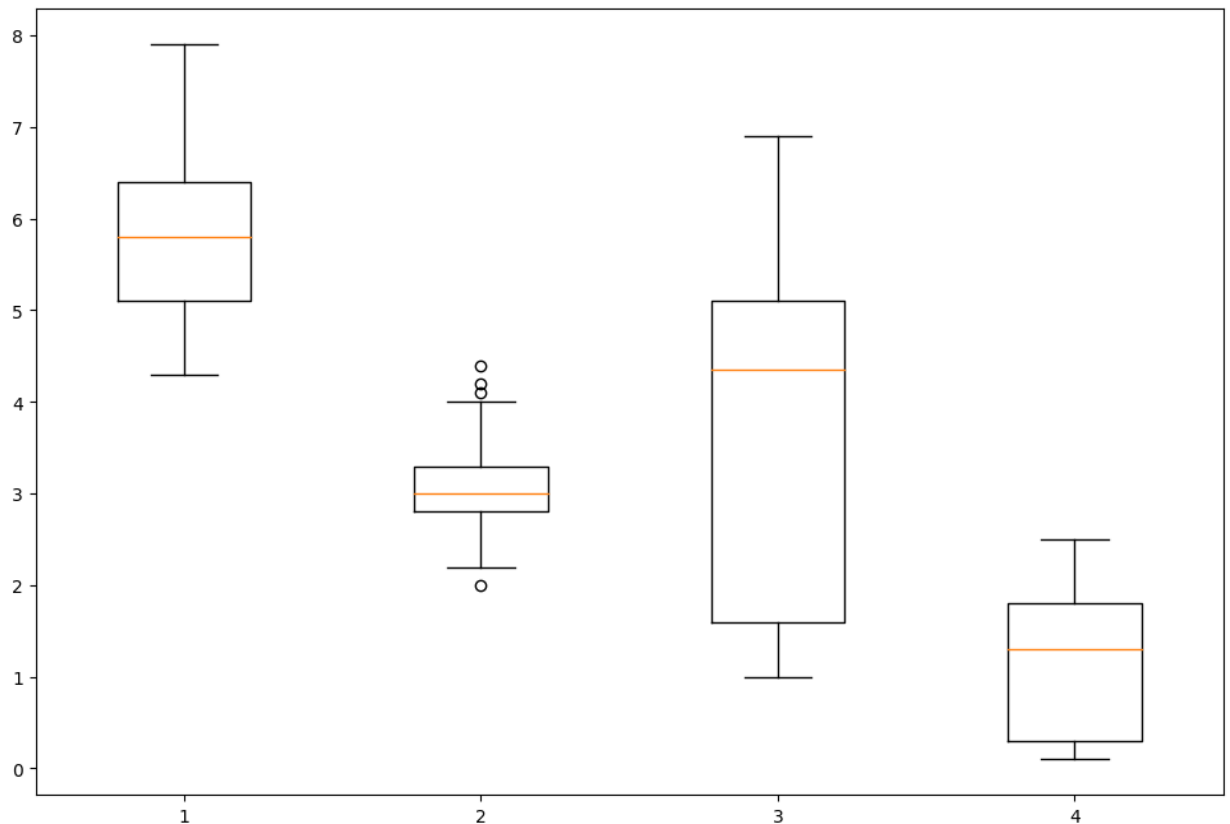
2. Create a histogram for each feature in the dataset to illustrate the feature distributions.

```
In [20]: fig, axes = plt.subplots(2, 2, figsize=(12, 6), constrained_layout = True)
for i in range(4):
    x, y = i // 2, i % 2
    axes[x, y].hist(df[df.columns[i + 1]])
    axes[x, y].set_title(f"Distribution of {df.columns[i + 1][:-2]}")
```



3. Create a boxplot for each feature in the dataset.

```
In [21]: data_to_plot = [df[x] for x in df.columns[1:-1]]
fig, axes = plt.subplots(1, figsize=(12,8))
bp = axes.boxplot(data_to_plot)
```



4. Compare distributions and identify outliers.

If we observe closely for the box 2, interquartile distance is roughly around 0.75 hence the values lying beyond this range of (third quartile + interquartile distance) i.e. roughly around 4.05 will be considered as outliers. Similarly outliers with other boxplots can be found.