

## PRACTICAL - 4

AIM:- Create the following plots using matplotlib pyplot.

- (a) Line graph for semester-wise grades.  
 import matplotlib.pyplot as plt.  
 import numpy as np.

Semesters = [1, 2, 3, 4, 5]

grades = [90, 85, 88, 92, 89]

programming\_languages = ['Python', 'Java', 'C++', 'JavaScript']

competencies = [98, 75, 80, 95]

percentages = np.array([75, 60, 70, 80, 65])

colors = ['blue', 'green', 'red', 'purple', 'orange']

plt.plot(semesters, grades, marker='o', color='b')

plt.title('Semester-Wise Grades')

plt.xlabel('Semester')

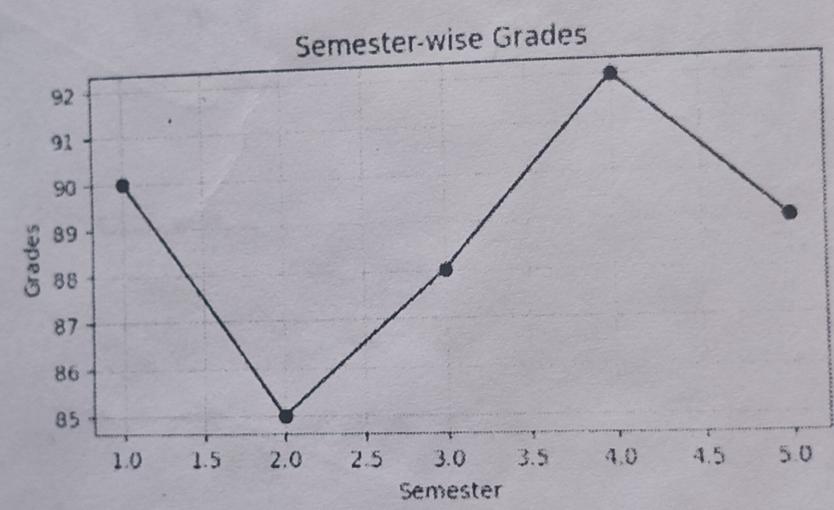
plt.ylabel('Grades')

plt.grid(True)

plt.show().

- (b) Pie-chart depicting your core competency in each of the following language learnt. Explode the wedges from the centre for each of the language.  
 Add labels to the wedges & a legend.

Output (a)



`explode = (0.5, 0, 0.3, 0)`

`plt.pie(competencies, labels=programming_languages,  
explode=explode, autopct='%1.1f%%')`

`plt.title('Core Competency in Programming Languages')`

`plt.legend(programming_languages, loc=3)  
plt.show().`

(c) Scatter plot depicting your %'s till Sem 5, use diff. colors for different semesters.

`plt.figure(figsize=(6, 3))`

for semester, color in zip(semesters, colors):

`plt.scatter(semester, percentages[semester - 1],  
label=f'Semester {semester}!', color=color)`

`plt.title('Scatter plot of Percentage by Semester')`

`plt.xlabel('Semester')`

`plt.ylabel('Percentage')`

`plt.legend()`

`plt.show().`

(d) Box chart depicting the % upto 75. Use height & width attribute, different colors for vertical & horizontal Barcharts.

`plt.figure(figsize=`

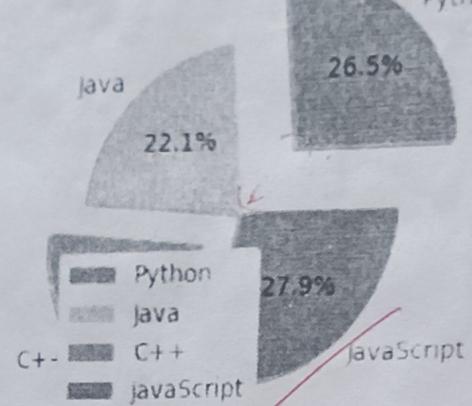
`plt.figure(figsize=(6, 3))`

# Horizontal Bar chart.

`plt.subplot(121)`

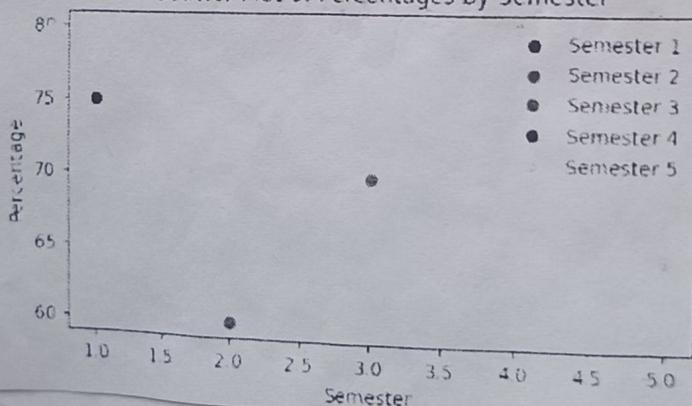
`plt.barh(semesters, percentages, color='skyblue',  
height=0.6)`

Core Competency in Programming Languages



~~Output (b)~~

Scatter Plot of Percentages by Semester



Output (c)

```

plt.title('Horizontal Bar Chart of Percentages')
plt.xlabel('Percentage')
plt.ylabel('Semester')
# Vertical bar chart.
plt.subplot(12)
plt.bar(semesters, percentages, color='lightcoral',
        width=0.6)
plt.title('Vertical Bar Chart of Percentages')
plt.xlabel('Semester')
plt.ylabel('Percentage')
plt.tight_layout()
plt.show().

```

(e) Create Subplots of  $3 \times 2$  dimensions.

~~fig, axes = plt.subplots(3, 2, figsize=(12, 10))~~  
~~fig.suptitle('Subplots (3x2)')~~

import matplotlib.pyplot as plt.

import numpy as np

$$x = np.linspace(0, 2 * np.pi, 100)$$

$$y1 = np.sin(x)$$

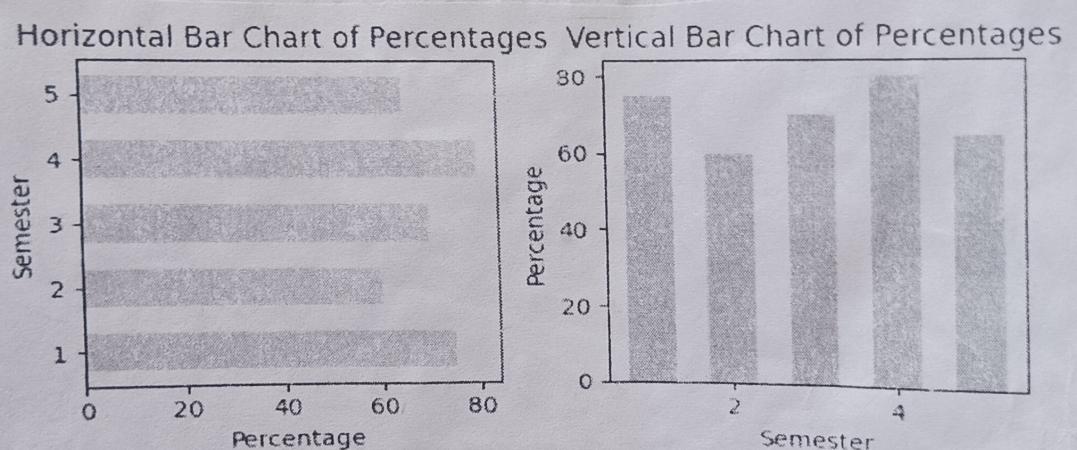
$$y2 = np.cos(x)$$

$$y3 = np.tan(x)$$

$$y4 = np.exp(x)$$

$$y5 = np.log(x+1)$$

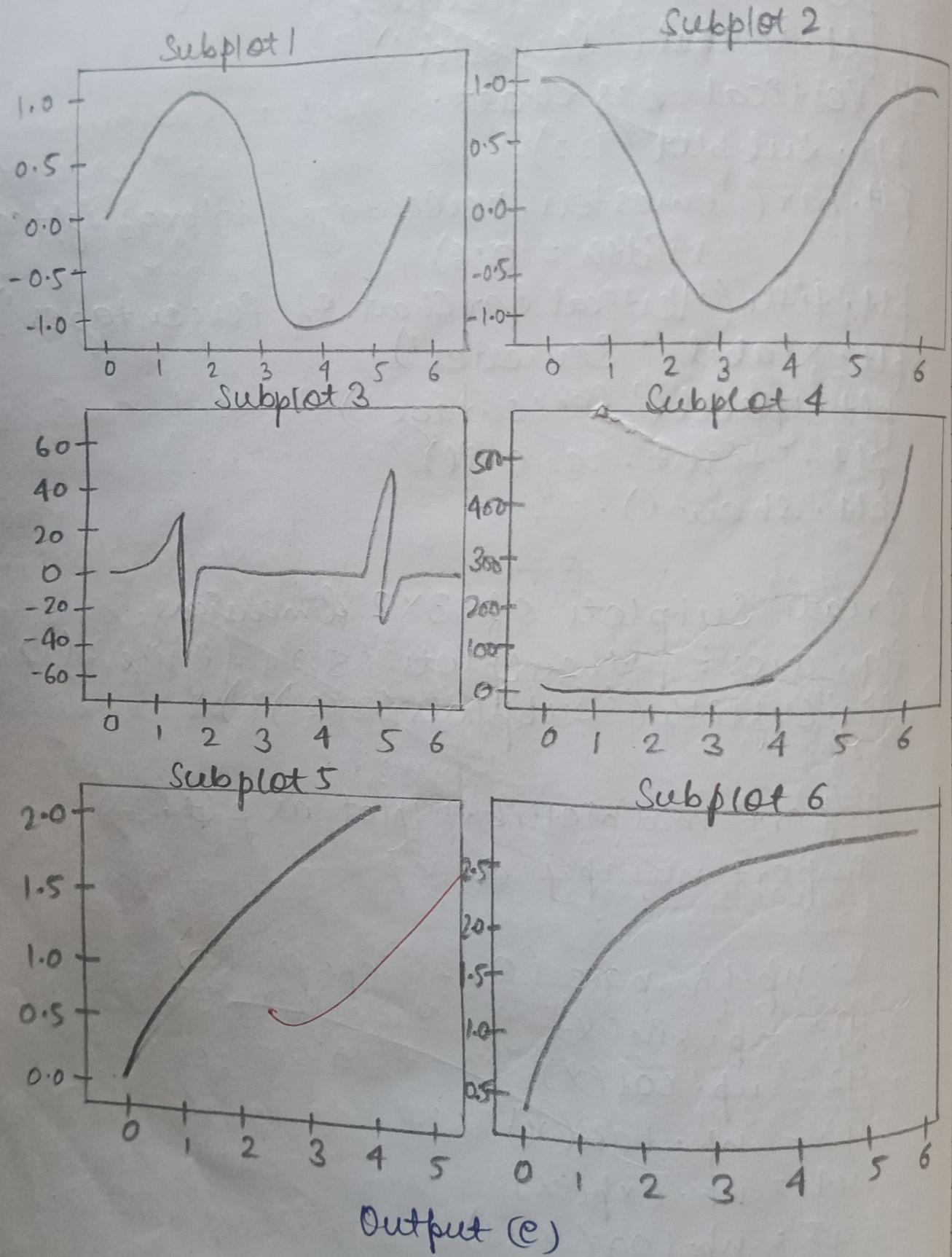
$$y6 = np.sqrt(x)$$



Output(d)

## OBSERVATION:-

Here, we have learnt different visualizations through Matplotlib like line plot, pie chart, subplot, etc. Matplotlib is one of the most popular & oldest plotting libraries in python which is used in Machine learning. In Machine learning, it helps to understand the huge amount of data through different visualizations e.g. Barplot, Histogram, Scatterplot, etc.



plt.figure(figsize=(10, 8))

plt.subplot(3, 2, 1)

plt.plot(x, y1)

plt.title('Subplot 1')

plt.subplot(3, 2, 2)

plt.plot(x, y2)

plt.title('Subplot 2')

plt.subplot(3, 2, 3)

plt.plot(x, y3)

plt.title('Subplot 3')

plt.subplot(3, 2, 4)

plt.plot(x, y4)

plt.title('Subplot 4')

plt.subplot(3, 2, 5)

plt.plot(x, y5)

plt.title('Subplot 5')

plt.subplot(3, 2, 6)

plt.plot(x, y6)

plt.title('Subplot 6')

plt.tight\_layout() # Add some space b/w subplots.  
plt.show().

## PRACTICAL - 5

AIM :- Working with a Seaborn Library to -

- (a) Create a distribution plot with histogram.
- ```
import matplotlib.pyplot as plt
import seaborn as sns
from numpy import random
import pandas as pd
```

```
sns.distplot([0, 1, 2, 3, 4, 5])
plt.show()
```

- (b) Create a distribution plot without histogram.
- ```
sns.distplot([0, 1, 2, 3, 4, 5], hist=False)
plt.show()
```

- (c) Create a Normal distribution plot.

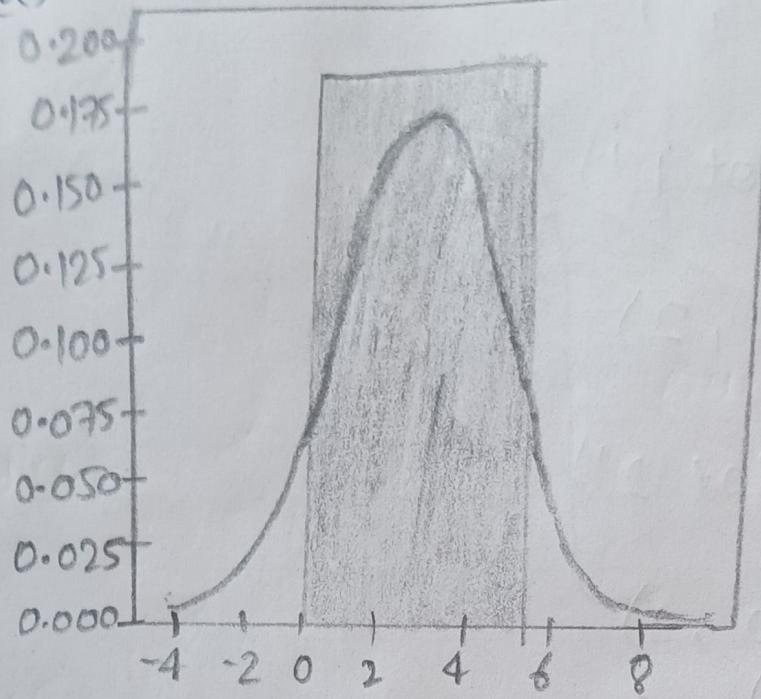
```
sns.distplot(random.normal(loc=50, scale=5,
size=1000), label='normal', hist=False)
plt.show()
```

- (d) Create a Binomial distribution plot & add the corresponding label.

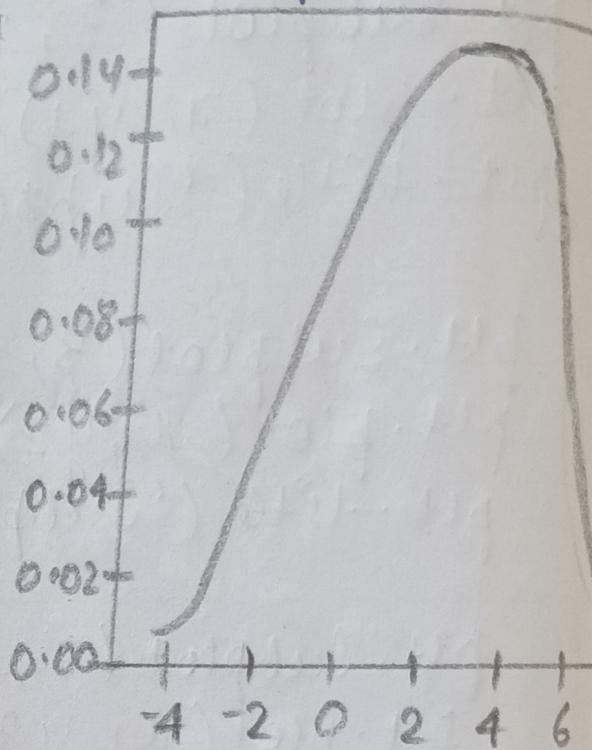
```
sns.distplot(random.binomial(n=10, p=0.5, size=1000),
label='Binomial', hist=False)
plt.show()
```

- (e) Use a CSV file to obtain lineplot b/w 2 parameters

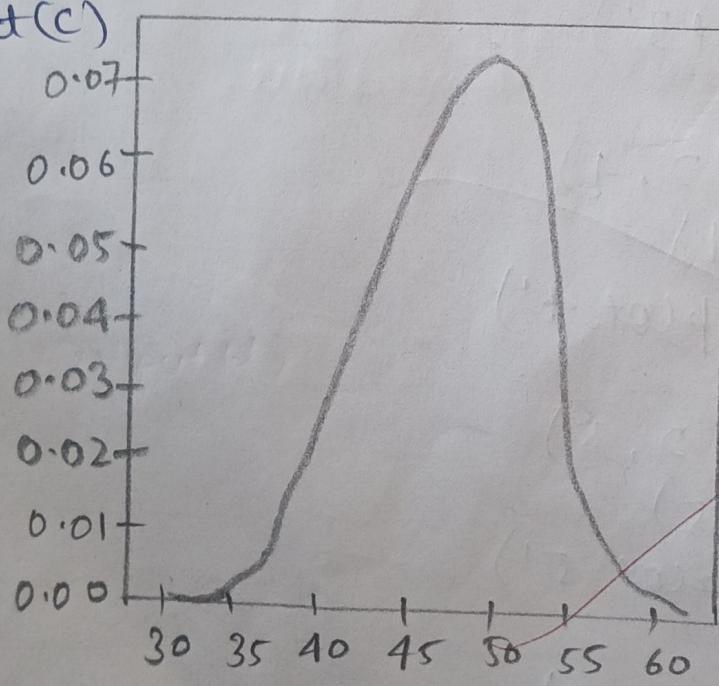
Output (a)



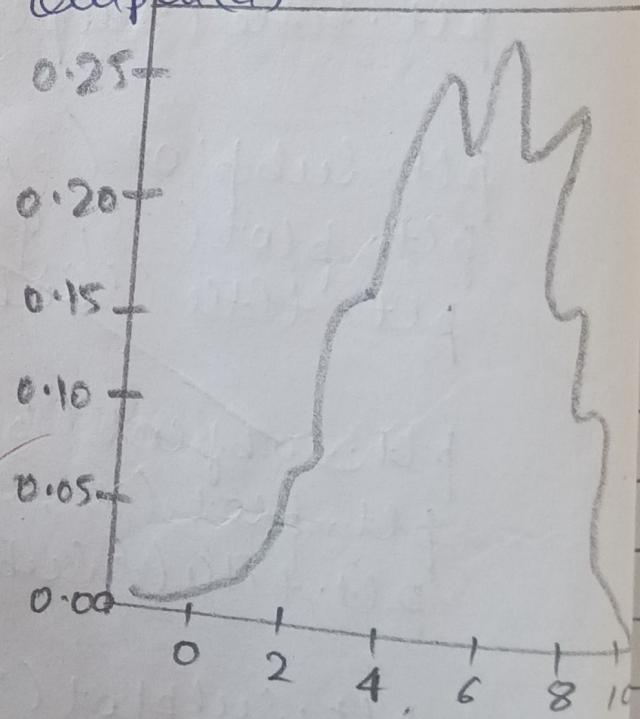
Output (b)



Output (c)



Output (d)



```
df = pd.read_csv("data.csv")
sns.lineplot(x=df['Col-A'], y=df['Col-B'])
plt.show()
```

(f) Draw a Scatterplot between 2 parameters.

```
sns.scatterplot(x=df['Col-A'], y=df['Col-B'])
plt.show()
```

(g) Draw a boxplot between 2 parameters.

```
sns.boxplot(x=df['Duration'], y=df['Maxpulse'])
plt.show()
```

(h) Draw a violinplot between the parameters

```
sns.violinplot(x=df['Duration'], y=df['Maxpulse'])
plt.show()
```

(i) Draw a swamplot between 2 parameters.

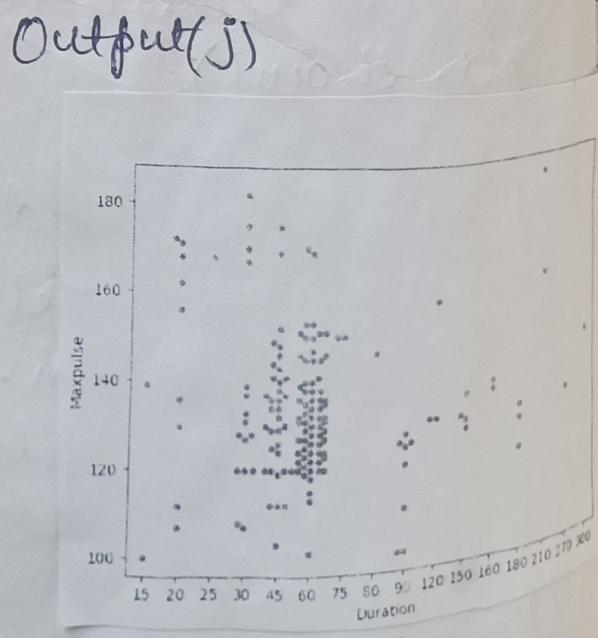
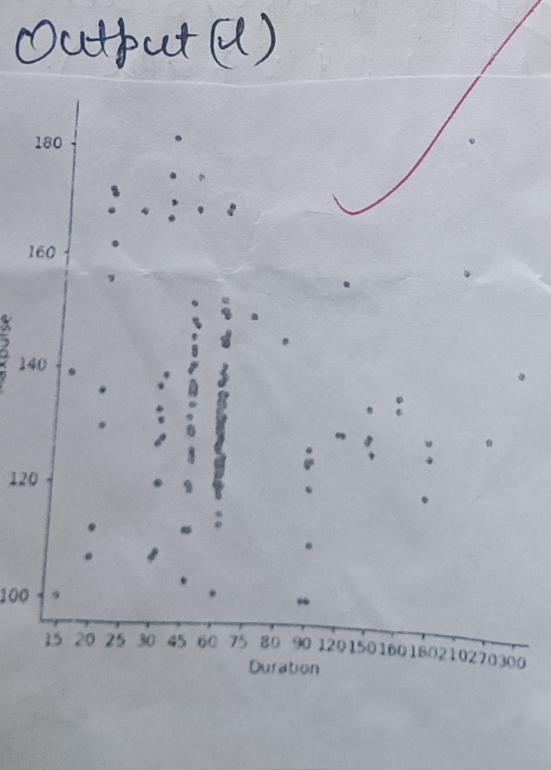
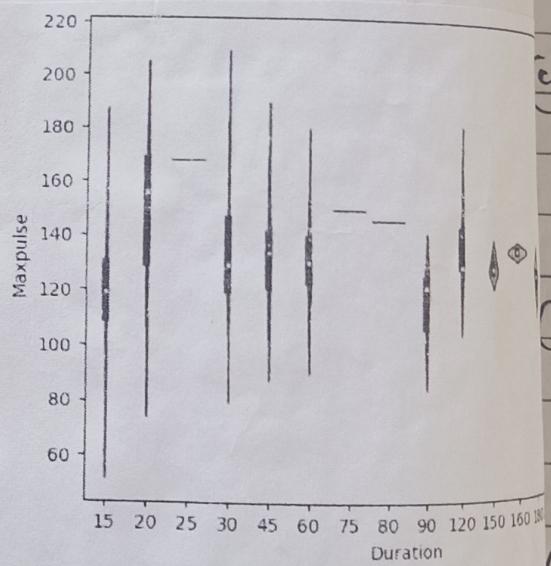
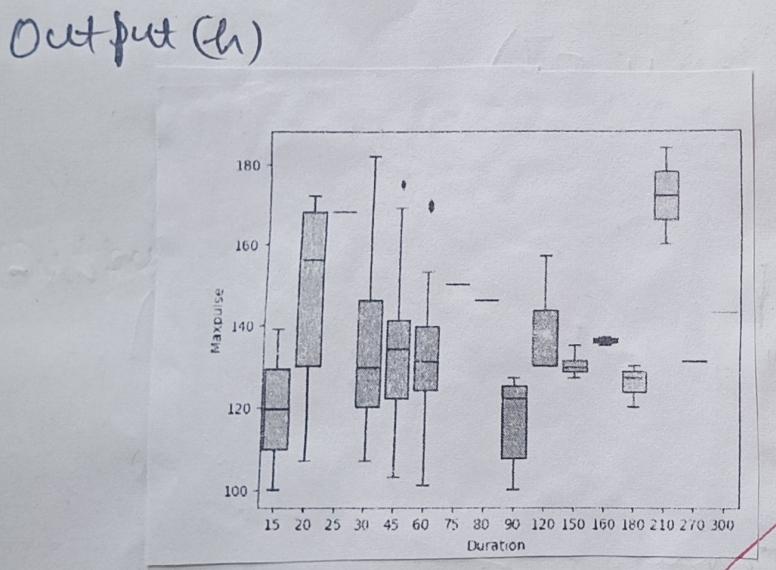
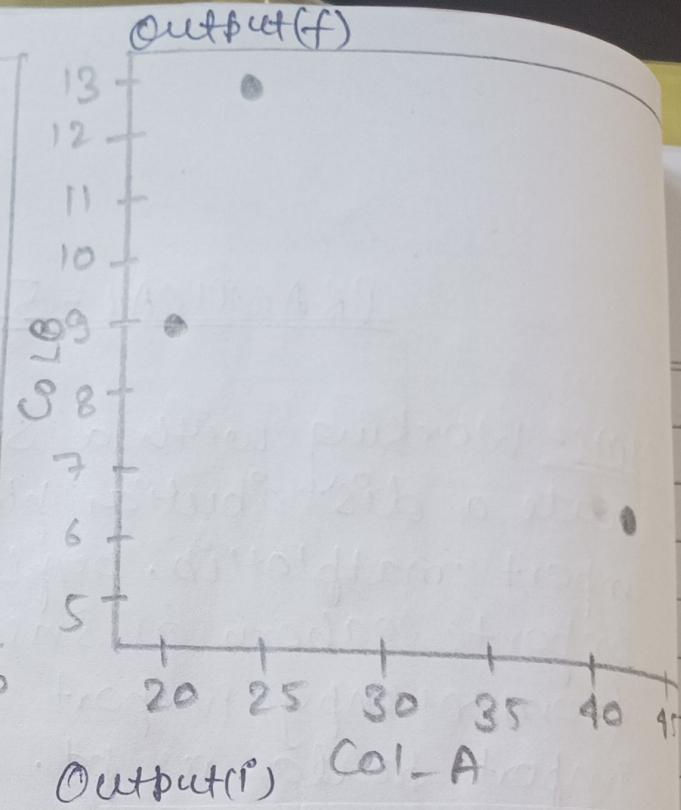
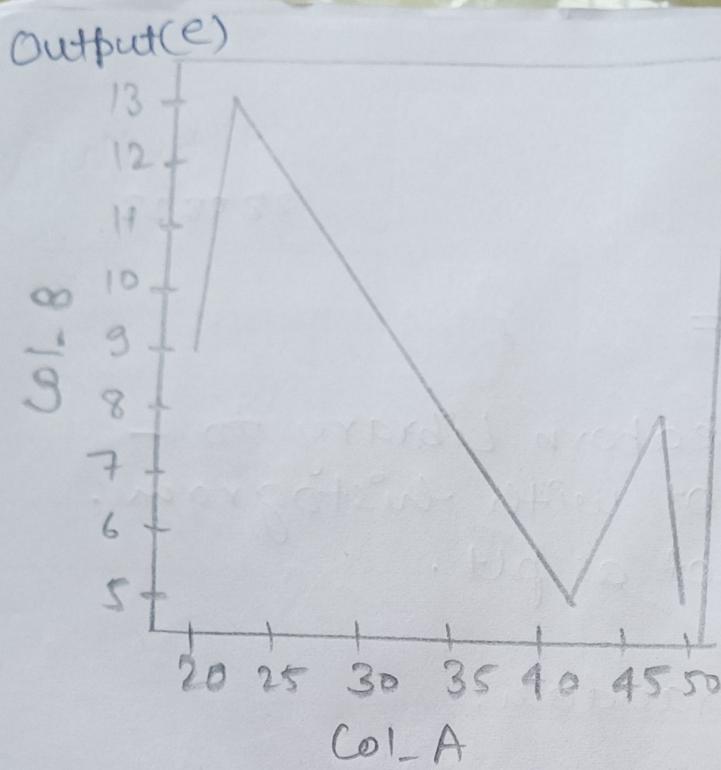
```
sns.swarmplot(x=df['Duration'], y=df['Maxpulse'])
plt.show()
```

(j) Draw a barplot between 2 parameters.

```
sns.barplot(x=df['Col-A'], y=df['Col-B'])
plt.show()
```

(k) Draw a catplot between 2 parameters.

```
sns.catplot(x=df['Duration'], y=df['Maxpulse'])
plt.show()
```



(m) Draw a pie chart between two parameters.

```
import matplotlib.pyplot as plt.
```

```
plt.pie(df = pd.read_csv('data.csv'))
```

```
duration_sum = df.groupby(['Duration'])['Calories'].sum()
```

```
plt.figure(figsize=(4,4))
```

```
plt.pie(duration_sum; labels=duration_sum_index,  
        autopct='%.1f %%', startangle=90).
```

```
plt.title('Pie chart between Duration & Calories')
```

```
plt.axis('equal')
```

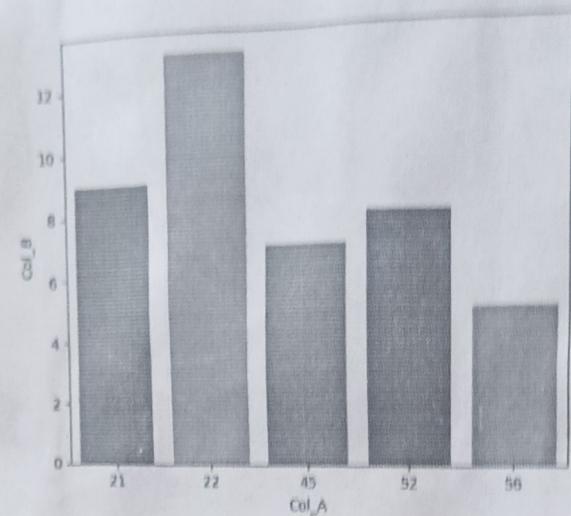
```
plt.show().
```

## OBSERVATION

~~boxplot ?~~ ~~- min max  
25th, 50th, 75th~~

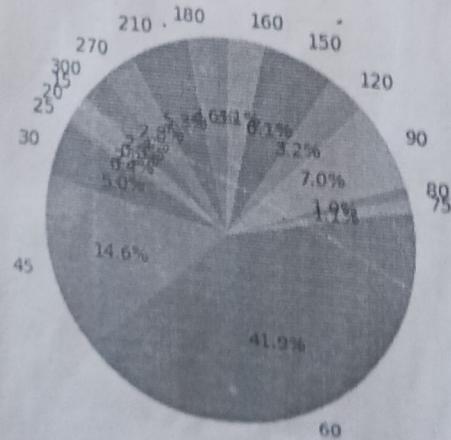
Seaborn with Matplotlib provide us with numerous allowing graphs through which one can easily analyse weak points, explore data with a deeper understanding & end up with a great insight into data & gaining the highest accuracy after training it. through different Algorithms.

Output (K)



Output(m)

Pie-Chart of Duration & Calories



PRACTICAL - 6

AIM: (a) Write to import load\_iris from sklearn datasets and describe all the parameters in load\_iris.

```
from sklearn.datasets import load_iris
import pandas as pd
```

```
iris = load_iris()
```

```
df = pd.DataFrame(iris.data, columns=iris.
```

feature\_names)

```
df['target'] = iris.target
```

df.

```
df.describe() # describe the dataframe.
```

# Attributes of load\_iris.

```
from sklearn.datasets import load_iris.
```

```
iris = load_iris()
```

x = iris.data

y = iris.target

feature\_names = iris.feature\_names

target\_names = iris.target\_names

print("Feature names:", feature\_names)

print("Target names:", target\_names)

	sepal length (cm)	sepal width (cm)	petal length (cm)	petal width (cm)	target
count	150.000000	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.057333	3.758000	1.199333	1.000000
std	0.828066	0.435866	1.765298	0.762238	0.819232
min	4.300000	2.000000	1.000000	0.100000	0.000000
25%	5.100000	2.800000	1.600000	0.300000	0.000000
50%	5.800000	3.000000	4.350000	1.300000	1.000000
75%	6.400000	3.300000	5.100000	1.800000	2.000000
max	7.900000	4.400000	6.900000	2.500000	2.000000

- (b) Draw a Scatter plot between the sepal-length and sepal-width.

```
data = sns.load_dataset('iris')
sns.scatterplot(data, x = 'sepal-length',
                 y = 'sepal-width', hue = 'sepal-length')
plt.show().
```

### PRACTICAL - 7

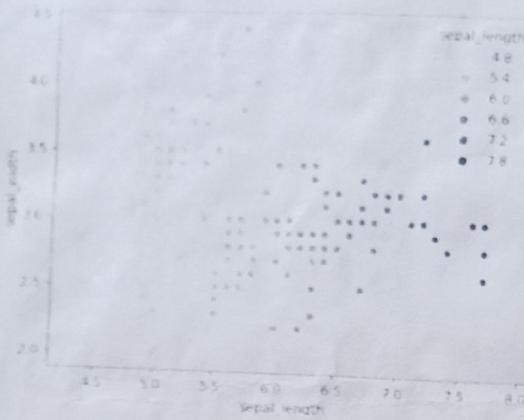
- (a) To perform Dataset loading using numpy, pandas & csv module.

- From Numpy

```
from numpy import loadtxt
data = r"C:\Users\hp\Downloads\pima-indians-diabetes.csv"
datapath = open(data, 'r')
data = loadtxt(datapath, delimiter = ',')
print(data.shape)
print(data[:3])
```

- Read a rsu file using header from pandas (without header)
- ```
from pandas import read_csv
path = r"C:\Users\hp\Downloads\pima-indians-diabetes.csv"
df = read_csv(path, delimiter = ',')
headernames = ['Sepal-length', 'Sepal-width', 'petal-length',
               'petal-width', 'target']
df = read_csv(path, names = headernames)
df.
```

## Output(b)



## Output (b) (768, 9)

[[6. 148. 72. 35. 0. 33.6 0.627 50.  
 [1. 85. 66. 29. 0. 26.6 0.351 31.  
 [8. 183. 64. 0. 0. 23.3 0.672 32.

|                      | Sepal-length | Sepal-width | petal-length | petal-width | Target |       |    |
|----------------------|--------------|-------------|--------------|-------------|--------|-------|----|
| 6                    | 148          | 72          | 35           | 0           | 33.6   | 0.627 | 50 |
| 1                    | 85           | 66          | 29           | 0           | 26.6   | 0.351 | 31 |
| 8                    | 183          | 64          | 0            | 0           | 23.3   | 0.672 | 32 |
| 1                    | 89           | 66          | 23           | 94          | 28.1   | 0.167 | 21 |
| 0                    | 137          | 40          | 35           | 168         | 43.1   | 2.288 | 33 |
| -                    | -            | -           | -            | -           | -      | -     | 1  |
| 10                   | 101          | 76          | 48           | 180         | 32.9   | 0.171 | 63 |
| 2                    | 122          | 70          | 27           | 0           | 36.8   | 0.340 | 27 |
| 6                    | 121          | 72          | 23           | 112         | 26.2   | 0.245 | 30 |
| 1                    | 120          | 60          | 0            | 0           | 30.1   | 0.349 | 47 |
| 93                   | 70           | 31          | 0            | 0           | 30.4   | 0.315 | 23 |
| 768 rows × 5 columns |              |             |              |             | 0      |       |    |

- Read a CSV file without header using pandas.

```
from pandas import read_csv
path = "C:/Users/hp/Downloads/diabetes.csv"
df = read_csv(path, delimiter = ", ")
df.
```

- Using CSV module.

```
import numpy as np
import csv
with open(path, 'r') as f:
    readdata = csv.reader(f)
    header = next(readdata)
    listdata = list(readdata)
    data = np.array(listdata)
print(header)
print(data[:3])
```

### OBSERVATION

~~Dataset structure of the groundwork of effective Machine-learning projects. Understanding the various kinds of dataset, the significance of data pre-processing, & the job of training & testing datasets are key stages towards building powerful models. By utilizing well-known sources such as kaggle, Aws, Data Researcher, Microsoft Datasets, government Datasets etc.~~

| Pregnancies | Glucose | BloodPressure | SkinThickness | Insulin | BMI | DiabetesPedigreeFunction | Age   | Outcome |   |
|-------------|---------|---------------|---------------|---------|-----|--------------------------|-------|---------|---|
| 0           | 8       | 148           | 72            | 35      | 0   | 33.6                     | 0.627 | 50      | 1 |
| 1           | 1       | 85            | 86            | 29      | 0   | 26.8                     | 0.351 | 31      | 0 |
| 2           | 8       | 183           | 64            | 0       | 0   | 23.3                     | 0.672 | 32      | 1 |
| 3           | 1       | 89            | 66            | 23      | 94  | 28.1                     | 0.167 | 21      | 0 |
| 4           | 0       | 137           | 40            | 36      | 168 | 43.1                     | 2.288 | 33      | 1 |
| ...         | ...     | ...           | ...           | ...     | ... | ...                      | ...   | ...     |   |
| 763         | 10      | 101           | 76            | 48      | 180 | 32.9                     | 0.171 | 63      | 0 |
| 764         | 2       | 122           | 70            | 27      | 0   | 36.8                     | 0.340 | 27      | 0 |
| 765         | 5       | 121           | 72            | 23      | 112 | 26.2                     | 0.245 | 30      | 0 |
| 766         | 1       | 126           | 60            | 0       | 0   | 30.1                     | 0.349 | 47      | 1 |
| 767         | 1       | 93            | 70            | 31      | 0   | 30.4                     | 0.315 | 23      | 0 |

768 rows x 9 columns

### F. ADDITIONAL

Output :- CSV module

```

['Pregnancies', 'Glucose', 'BloodPressure', 'SkinThickness',
 'Insulin', 'BMI', 'DiabetesPedigreeFunction', 'Age',
 [[{"Pregnancies": 1, "Glucose": 85, "BloodPressure": 66, "SkinThickness": 29, "Insulin": 0, "BMI": 26.8, "DiabetesPedigreeFunction": 0.351, "Age": 31, "Outcome": 0}, {"Pregnancies": 0, "Glucose": 148, "BloodPressure": 72, "SkinThickness": 35, "Insulin": 0, "BMI": 33.6, "DiabetesPedigreeFunction": 0.627, "Age": 50, "Outcome": 1}, {"Pregnancies": 1, "Glucose": 79, "BloodPressure": 64, "SkinThickness": 23, "Insulin": 94, "BMI": 28.1, "DiabetesPedigreeFunction": 0.167, "Age": 21, "Outcome": 0}, {"Pregnancies": 2, "Glucose": 183, "BloodPressure": 64, "SkinThickness": 0, "Insulin": 0, "BMI": 23.3, "DiabetesPedigreeFunction": 0.672, "Age": 32, "Outcome": 1}, {"Pregnancies": 3, "Glucose": 137, "BloodPressure": 40, "SkinThickness": 36, "Insulin": 168, "BMI": 43.1, "DiabetesPedigreeFunction": 2.288, "Age": 33, "Outcome": 1}, {"Pregnancies": 4, "Glucose": 101, "BloodPressure": 76, "SkinThickness": 48, "Insulin": 180, "BMI": 32.9, "DiabetesPedigreeFunction": 0.171, "Age": 63, "Outcome": 0}, {"Pregnancies": 5, "Glucose": 122, "BloodPressure": 70, "SkinThickness": 27, "Insulin": 0, "BMI": 36.8, "DiabetesPedigreeFunction": 0.340, "Age": 27, "Outcome": 0}, {"Pregnancies": 6, "Glucose": 121, "BloodPressure": 72, "SkinThickness": 23, "Insulin": 112, "BMI": 26.2, "DiabetesPedigreeFunction": 0.245, "Age": 30, "Outcome": 0}, {"Pregnancies": 7, "Glucose": 126, "BloodPressure": 60, "SkinThickness": 0, "Insulin": 0, "BMI": 30.1, "DiabetesPedigreeFunction": 0.349, "Age": 47, "Outcome": 1}, {"Pregnancies": 8, "Glucose": 93, "BloodPressure": 70, "SkinThickness": 31, "Insulin": 0, "BMI": 30.4, "DiabetesPedigreeFunction": 0.315, "Age": 23, "Outcome": 0}]]
```

## PRACTICAL-8

AIM:- To perform the Outlier Detection using statistical and Visualization Technique-i.e., (box plot and Scatter plot).

- Statistical technique i.e Z-score and IQR Range.

from sklearn.datasets import load\_iris

iris = load\_iris()

x = iris.data

y = iris.target

feature\_names = iris.feature\_names

target\_names = iris.target\_names

# Z-Score Method

import pandas as pd

threshold = 2

column\_name = ['sepal\_length', 'sepal\_width',  
'petal\_length', 'petal\_width']

df = pd.DataFrame(x, columns=column\_name)

sepal\_length = df["sepal\_length"]

z\_score = (sepal\_length - sepal\_length.mean()) /  
sepal\_length.std()

print("Z-score: ", z\_score)

outliers = df[abs(z\_score) > threshold]

print("Outliers: \n", outliers)

null\_values = df.isnull() # check Null Values

deleted\_df = df.drop(outliers.index, axis=0)

z-score: 0 -0.897674

1 -1.139200  
2 -1.380727  
3 -1.501490  
4 -1.018437

...  
145 1.034539  
146 0.551486  
147 0.793012  
148 0.430722  
149 0.068433

Name: sepal\_length, Length: 150, dtype: float64  
outliers:

|     | sepal_length | sepal_width | peta_length | petal_width |
|-----|--------------|-------------|-------------|-------------|
| 105 | 7.6          | 3.0         | 6.6         | 2.1         |
| 117 | 7.7          | 3.8         | 6.7         | 2.2         |
| 118 | 7.7          | 2.6         | 6.9         | 2.3         |
| 122 | 7.7          | 2.8         | 6.7         | 2.0         |
| 131 | 7.9          | 3.8         | 6.4         | 2.0         |
| 135 | 7.7          | 3.0         | 6.1         | 2.3         |

## # IQR (Inter-Quartile Range)

$Q_1 = df['sepal\_length'].quantile(0.25)$

$Q_3 = df['sepal\_length'].quantile(0.75)$

$IQR = Q_3 - Q_1$

`print("IQR: ", IQR)`

$threshold = 1.5 * IQR$

~~outliers = df[(df['sepal\_length']) < (q1 - threshold) | (df['sepal\_length'] > (q3 + threshold))]~~

`print("outliers are : ", outliers)`

## • Visualization Technique i.e., boxplot and Scatterplot.

`import seaborn as sns`

`import matplotlib.pyplot as plt`

`from sklearn.datasets import load_iris`

`iris = load_iris()`

`data = iris.data`

`target = iris.target`

`import pandas as pd`

~~`df = pd.DataFrame(data, columns=iris.feature_names)`~~

~~`df['Species'] = iris.target_names[target]`~~

`# Boxplot`

`plt.figure(figsize=(6, 3))`

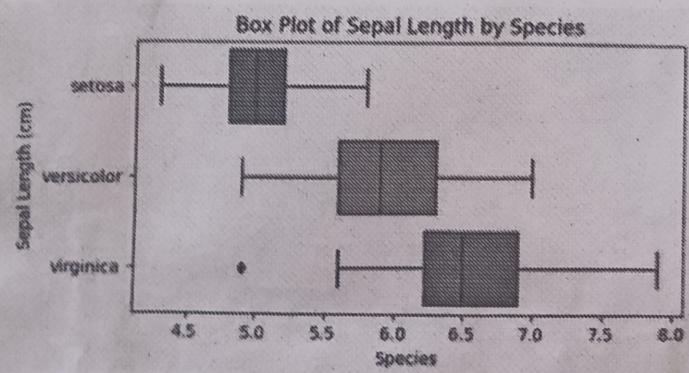
`sns.boxplot(x='sepal_length(cm)', y='Species', data=df)`

`plt.title("Boxplot of Sepal length by Species")`

`plt.xlabel("Species")`

`plt.ylabel("Sepal length(cm)")`

`plt.show()`



# Scatter plot

import seaborn as sns

import pandas as pd

import matplotlib.pyplot as plt

from sklearn.datasets import load\_iris

iris = load\_iris()

data = iris.data

target = iris.target

df = pd.DataFrame(data, columns=iris.feature\_names)

df['Species'] = iris.target\_names[target]

plt.figure(figsize=(6, 3))

sns.scatterplot(x='sepal.length(cm)', y='Species',  
data=df)

plt.title("Scatter plot of Sepal length by Species")

plt.xlabel("Sepal length (cm)")

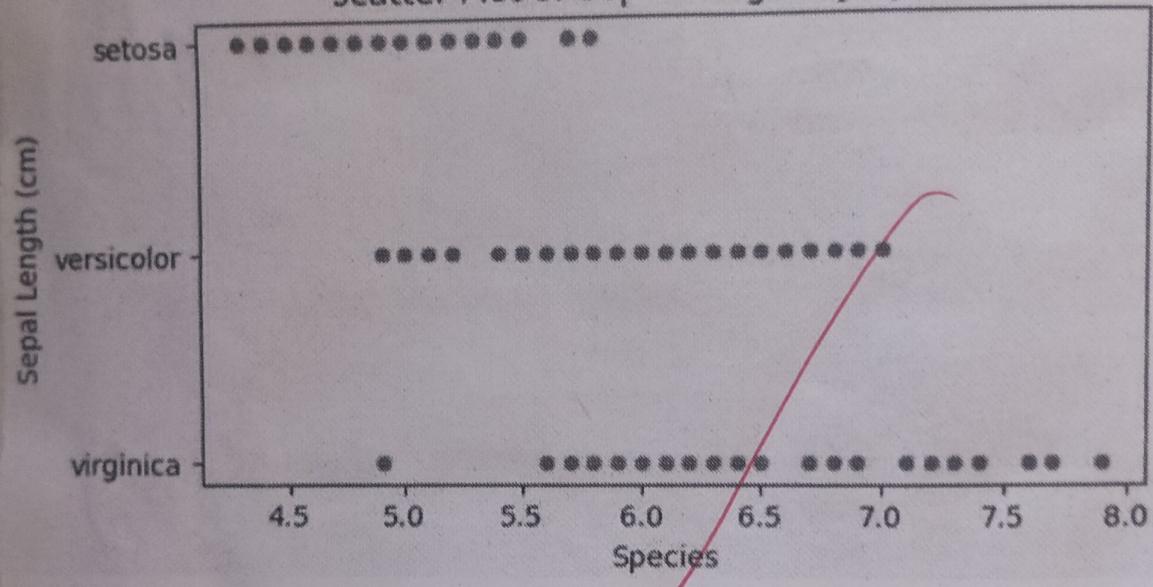
plt.ylabel("Species")

plt.show()

## OBSERVATION

In this, using the IQR method can find more "outlier" than using the z-score (standard deviation) method. It depends on the distribution of the data. Data that's peaked with long tails will have a comparatively low IQR, so the IQR method will find lot of outliers.

scatter Plot of Sepal Length by Species



## PRACTICAL-9

AIM:- To handle the missing numerical data and perform One hot Encoding on Categorical Data.

- Handling the missing numerical data :-  
 import pandas as pd  
 import numpy as np.

```
data = { 'Age': [20, 30, np.nan, 26, 35, np.nan],  

         'Salary': [40000, np.nan, 60000, 70000,  

                     np.nan, 80000],  

         'Rank-score': [82, np.nan, 78, np.nan, 81, 92]  

     }
```

```
df = pd.DataFrame(data)
```

```
print("dataframe : \n", df)
```

```
mean_data = df.mean()
```

```
df.fillna(mean_data, inplace=True)
```

```
print("after missing values handled : \n", df)
```

- One-Hot Encoding by get-dummies().

```
import numpy as np.
```

```
import pandas as pd.
```

```
data = pd.read_csv('Employee.csv')
```

```
one_hot_encoded_data = pd.get_dummies(data,  

   columns=['Gender', 'Remzeki'])
```

```
one_hot_encoded_data.
```

dataframe:  
 Age salary Rank\_score  
 0 20.0 40000.0 82.0  
 1 30.0 NaN NaN  
 2 NaN 60000.0 78.0  
 3 26.0 70000.0 NaN  
 4 35.0 NaN 81.0  
 5 NaN 80000.0 92.0

after missing values handled:  
 Age salary Rank\_score  
 0 20.00 40000.0 82.00  
 1 30.00 62500.0 83.25  
 2 27.75 60000.0 78.00  
 3 26.00 70000.0 83.25  
 4 35.00 62500.0 81.00  
 5 27.75 80000.0 92.00

|   | Employee_ID | Gender_Female | Gender_Male | Remarks_Good | Remarks_Great | Remarks_Nice |
|---|-------------|---------------|-------------|--------------|---------------|--------------|
| 0 | 45          | 0             | 1           | 0            | 0             | 1            |
| 1 | 12          | 1             | 0           | 1            | 0             | 0            |
| 2 | 78          | 1             | 0           | 0            | 1             | 0            |
| 3 | 35          | 1             | 0           | 0            | 0             | 1            |
| 4 | 7           | 0             | 1           | 0            | 1             | 0            |

• By sklearn.preprocessing module.

import pandas as pd.

import numpy as np.

from sklearn.preprocessing import OneHotEncoder

data = pd.read\_csv('Employee.csv')

data['Gender'] = data['Gender'].astype('category')

data['Remarks'] = data['Remarks'].astype('category')

data['Gen-new'] = data['Gender'].cat.codes

data['Rem-new'] = data['Remarks'].cat.codes

# Create an instance of one-hot encoder.

enc = OneHotEncoder()

# passing encoded column

enc\_data = pd.DataFrame(enc.fit\_transform(data[['Gen-new', 'Rem-new']]))

 toarray()

# Merge with main.

New\_df = data.join(enc\_data)

print("one-hot encoding: \n", New\_df)

### OBSERVATION:-

Here, we have converted the enc.fit\_transform() method to an array because the fit\_transform method of OneHotEncoder returns SciPy Sparse matrix so converting to an array first enables us to save space when we have a huge number of categorical variables.

Day 8/10/19

one-hot encoding:

|   | Employee_ID | Gender | Remarks | Gen_new | Rem_new | 0   | 1   | 2   | 3   | 4   |
|---|-------------|--------|---------|---------|---------|-----|-----|-----|-----|-----|
| 0 | 45          | Male   | Nice    | 1       | 2       | 0.0 | 1.0 | 0.0 | 0.0 | 1.0 |
| 1 | 12          | Female | Good    | 0       | 0       | 1.0 | 0.0 | 1.0 | 0.0 | 0.0 |
| 2 | 78          | Female | Great   | 0       | 1       | 1.0 | 0.0 | 0.0 | 1.0 | 0.0 |
| 3 | 35          | Female | Nice    | 0       | 2       | 1.0 | 0.0 | 0.0 | 0.0 | 1.0 |
| 4 | 7           | Male   | Great   | 1       | 1       | 0.0 | 1.0 | 0.0 | 1.0 | 0.0 |