

In [ ]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

In [ ]:

```
df = pd.read_csv('Iris.csv')
df.head()
```

Out[ ]:

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

In [ ]:

```
df.drop(columns = ['Id'], axis = 1, inplace = True)
```

In [ ]:

```
df['Species'].value_counts()
```

Out[ ]:

```
Iris-setosa      50
Iris-versicolor 50
Iris-virginica  50
Name: Species, dtype: int64
```

In [ ]:

```
#Categorizing data species wise
setosa_df = df[df['Species'] == 'Iris-setosa']
versicolor_df = df[df['Species'] == 'Iris-versicolor']
virginica_df = df[df['Species'] == 'Iris-virginica']
```

In [ ]:

```
#Iris - setosa
setosa_df.head()
```

Out[ ]:

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

In [ ]:

```
#Iris - setosa statistical description
setosa_df.describe()
```

Out[ ]:

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
count	50.00000	50.000000	50.000000	50.00000
mean	5.00600	3.418000	1.464000	0.24400
std	0.35249	0.381024	0.173511	0.10721
min	4.30000	2.300000	1.000000	0.10000
25%	4.80000	3.125000	1.400000	0.20000
50%	5.00000	3.400000	1.500000	0.20000
75%	5.20000	3.675000	1.575000	0.30000
max	5.80000	4.400000	1.900000	0.60000

In [ ]:

```
#Iris - versicolor
versicolor_df.head()
```

Out[ ]:

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
50	7.0	3.2	4.7	1.4	Iris-versicolor
51	6.4	3.2	4.5	1.5	Iris-versicolor
52	6.9	3.1	4.9	1.5	Iris-versicolor
53	5.5	2.3	4.0	1.3	Iris-versicolor
54	6.5	2.8	4.6	1.5	Iris-versicolor

In [ ]:

```
#Iris - versicolor statistical description
versicolor_df.describe()
```

Out[ ]:

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
<b>count</b>	50.000000	50.000000	50.000000	50.000000
<b>mean</b>	5.936000	2.770000	4.260000	1.326000
<b>std</b>	0.516171	0.313798	0.469911	0.197753
<b>min</b>	4.900000	2.000000	3.000000	1.000000
<b>25%</b>	5.600000	2.525000	4.000000	1.200000
<b>50%</b>	5.900000	2.800000	4.350000	1.300000
<b>75%</b>	6.300000	3.000000	4.600000	1.500000
<b>max</b>	7.000000	3.400000	5.100000	1.800000

In [ ]:

```
#Iris - virginica
virginica_df.head()
```

Out[ ]:

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm	Species
<b>100</b>	6.3	3.3	6.0	2.5	Iris-virginica
<b>101</b>	5.8	2.7	5.1	1.9	Iris-virginica
<b>102</b>	7.1	3.0	5.9	2.1	Iris-virginica
<b>103</b>	6.3	2.9	5.6	1.8	Iris-virginica
<b>104</b>	6.5	3.0	5.8	2.2	Iris-virginica

In [ ]:

```
#Iris - virginica statistical description  
virginica_df.describe()
```

Out[ ]:

	SepalLengthCm	SepalWidthCm	PetalLengthCm	PetalWidthCm
<b>count</b>	50.00000	50.000000	50.000000	50.00000
<b>mean</b>	6.58800	2.974000	5.552000	2.02600
<b>std</b>	0.63588	0.322497	0.551895	0.27465
<b>min</b>	4.90000	2.200000	4.500000	1.40000
<b>25%</b>	6.22500	2.800000	5.100000	1.80000
<b>50%</b>	6.50000	3.000000	5.550000	2.00000
<b>75%</b>	6.90000	3.175000	5.875000	2.30000
<b>max</b>	7.90000	3.800000	6.900000	2.50000

In [ ]:

# Data Science and Big Data Analytics

Group : A

## Lab Assignment : 04

The screenshot shows two Jupyter Notebook sessions side-by-side.

**Session 1 (Left):**

- In [1]:

```
import numpy as np
import pandas as pd
from sklearn import datasets
from sklearn.linear_model import LinearRegression
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
```

- In [2]:

```
#loading the dataset directly from sklearn
boston = datasets.load_boston()

C:\Users\rushil\AppData\Local\Programs\Python\Python310\lib\site-packages\sklearn\utils\deprecation.py:87: FutureWarning: Function load_boston is deprecated; load_boston is deprecated in 1.0 and will be removed in 1.2.

The Boston housing prices dataset has an ethical problem. You can refer to the documentation of this function for further details.

The scikit-learn maintainers therefore strongly discourage the use of this dataset unless the purpose of the code is to study and educate about ethical issues in data science and machine learning.

In this special case, you can fetch the dataset from the original source:

import pandas as pd
import numpy as np

data_url = "http://lib.stat.cmu.edu/datasets/boston"
raw_df = pd.read_csv(data_url, sep="\s+", skiprows=22, header=None)
data = np.hstack([raw_df.values[:,2:], raw_df.values[:,1:2, :2]])
target = raw_df.values[:,1:2, 2]
```

**Session 2 (Right):**

- In [1]:

```
data_url = "http://lib.stat.cmu.edu/datasets/boston"
raw_df = pd.read_csv(data_url, sep="\s+", skiprows=22, header=None)
data = np.hstack([raw_df.values[:,2:], raw_df.values[:,1:2, :2]])
target = raw_df.values[:,1:2, 2]

Alternative datasets include the California housing dataset (i.e., :func:`~sklearn.datasets.fetch_california_housing`) and the Ames housing dataset. You can load the datasets as follows:

from sklearn.datasets import fetch_california_housing
housing = fetch_california_housing()

for the California housing dataset and:

from sklearn.datasets import fetch_openml
housing = fetch_openml(name="house_prices", as_frame=True)

for the Ames housing dataset.

warnings.warn(msg, category=FutureWarning)
```

- In [2]:

```
print(boston)
```

- In [3]:

```
[{'data': array([[6.3200e-03, 1.8000e+01, 2.3100e+00, ..., 1.5300e+01, 3.9690e+02,
   4.9800e+00],
  [2.7310e-02, 0.0000e+00, 7.0700e+00, ..., 1.7800e+01, 3.9690e+02,
  9.1400e+00],
  [2.7290e-02, 0.0000e+00, 7.0700e+00, ..., 1.7800e+01, 3.9283e+02,
  4.0300e+00],
  ...,
  [6.0760e-02, 0.0000e+00, 1.1930e+01, ..., 2.1000e+01, 3.9690e+02,
  5.6400e+00],
```



A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A4 - Jupyter Notebook | +

localhost:8888/notebooks/DSBDA\_Lab\_A4.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel) O

**jupyter DSBDA\_Lab\_A4 Last Checkpoint: 4 hours ago (autosaved)**

Code

```

LSTAT, INDUS, NOX, RM, AGE, DIS, RAD, TAX, PTRATIO, B, LSTAT]], dtype='<U7' ), 'DESCR': '.. boston dataset:\nBoston house prices dataset\n-----\nData Set Characteristics:-- \n :Number of Instances: 506 \n :Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target.\n :Attribute Information (in order):\n - CRIM per capita crime rate by town\n - ZN proportion of residential land zoned for lots over 25,000 sq.ft.\n - INDUS proportion of non-retail business acres per town\n - CHAS Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)\n - NOX nitric oxides concentration (parts per 10 million)\n - RM average number of rooms per dwelling\n - AGE proportion of owner-occupied units built prior to 1940\n - DIS weighted distances to five employment centres\n - RAD index of accessibility to radial highways\n - TAX full-value property-tax rate per $10,000\n - PTRATIO pupil-teacher ratio by town\n - B 1000(Bk - 0.63)^2 w here Bk is the proportion of black people by town\n - LSTAT % lower status of the population\n :Missing Attribute Values: None\n :Creator: Harrison, D. and Rubinfeld, D.L.\n :Source: This is a copy of UCI ML housing dataset. (https://archive.ics.uci.edu/ml/machine-learning-databases/housing/)\nThis dataset was taken from the Statistician which is maintained at Carnegie Mellon University.\n\nThe Boston house-prices data of Harrison, D., and Rubinfeld, D.L., 'Hedonic Prices and the Demand for Clean Air', J. Environ. Economics & Management, vol.5, 81-102, 1978. Used in Belsley, Kuh & Welsch, 'Regression diagnostics', Wiley, 1980. N.B. Various transformations are used in the table on pages 244-261 of the latter.\n\nThe Boston house-price data has been used in many machine learning papers to address regression problems.\n... topic: References\n... Belsley, Kuh & Welsch, 'Regression diagnostics: Identifying Influential Data and Sources of Collinearity', Wiley, 1980. 244-261.\n... Quinlan, R. (1993). Combining Instance-Based and Model-Based Learning. In Proceedings on the Tenth International Conference of Machine Learning, 236-243, University of Massachusetts, Amherst. Morgan Kaufmann.\n', 'filename': 'boston_house_prices.csv', 'data_module': 'sklearn.datasets.data'}

```

**sklearn returns Dictionary-like object, the interesting attributes are: 'data', the data to learn, 'target', the regression targets, 'DESCR', the full description of the dataset, and 'filename', the physical location of boston csv dataset.**

In [4]:

```

print(type(boston))
print('\n')
print(boston.keys())
print('\n')
print(boston.data.shape)

```

Out[4]:

```

<class 'sklearn.utils.Bunch'>

dict_keys(['data', 'target', 'feature_names', 'DESCR', 'filename', 'data_module'])

(506, 13)
['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX', 'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT']

```

In [5]:

```

#the details about the features and more information about the dataset can be seen by using boston.DESCR
print(boston.DESR)
.. _boston_dataset:
Boston house prices dataset
-----
**Data Set Characteristics:**

:Number of Instances: 506

:Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target.

```

In [5]: #The details about the features and more information about the dataset can be seen by using boston.DESCRIPTOR  
**print(boston.DESCRIPTOR)**

```
.. _boston_dataset:
-----
Boston house prices dataset
-----
**Data Set Characteristics:**
:Number of Instances: 506
:Number of Attributes: 13 numeric/categorical predictive. Median Value (attribute 14) is usually the target.
:Attribute Information (in order):
- CRIM per capita crime rate by town
- ZN proportion of residential land zoned for lots over 25,000 sq.ft.
- INDUS proportion of non-retail business acres per town
- CHAS Charles River dummy variable (= 1 if tract bounds river; 0 otherwise)
- NOX nitric oxide concentration (parts per 10 million)
- RM average number of rooms per dwelling
- AGE proportion of owner-occupied units built prior to 1940
- DIS weighted distances to five Boston employment centres
- RAD index of accessibility to radial highways
- TAX full-value property-tax rate per $10,000
- PTRATIO pupil-teacher ratio by town
- B 1000(Bk - 0.63)^2 where Bk is the proportion of black people by town
- LSTAT % lower status of the population
- MEDV Median value of owner-occupied homes in $1000's

:Missing Attribute Values: None
:Creator: Harrison, D. and Rubinfeld, D.L.
```

In [6]: #Before applying any model we have to convert this to a pandas dataframe,
#which we can do by calling the dataframe on boston.data. We also adds the target variable to the dataframe from boston.target

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A4 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A4.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [6]: #Before applying any model we have to convert this to a pandas dataframe, which we can do by calling the dataframe on boston.data. We also adds the target variable to the dataframe from boston.target  
bos = pd.DataFrame(boston.data, columns = boston.feature\_names)  
bos['PRICE'] = boston.target  
  
print(bos.head())

```

CRIM      ZN     INDUS    CHAS    NOX      RM     AGE      DIS      RAD      TAX
0  0.00632  18.0   2.31  0.0  0.538  6.575  65.2  4.0900  1.0  296.0
1  0.02731  0.0   7.07  0.0  0.469  6.421  78.9  4.9671  2.0  242.0
2  0.02729  0.0   7.07  0.0  0.469  7.185  61.1  4.9671  2.0  242.0
3  0.03237  0.0   2.18  0.0  0.458  6.998  45.8  6.0622  3.0  222.0
4  0.06905  0.0   2.18  0.0  0.458  7.147  54.2  6.0622  3.0  222.0

PTRATIO      B      LSTAT      PRICE
0  15.3  396.90  9.14  24.0
1  17.8  396.90  9.14  21.6
2  17.8  392.83  4.03  34.7
3  18.7  394.63  2.94  33.4
4  18.7  396.90  5.33  36.2

```

In [7]: #Get some statistics from dataset  
print(bos.describe())

```

CRIM      ZN     INDUS    CHAS    NOX      RM
count  506.000000  506.000000  506.000000  506.000000  506.000000
mean   3.613524  11.363636  11.136779  0.069170  0.554695  6.284634
std    8.601545  23.322453  6.869353  0.253994  0.115878  0.702617
min    0.006320  0.000000  0.469000  0.000000  0.385000  3.561000
25%   0.082045  0.000000  5.190000  0.000000  0.449000  5.885500
50%   0.256510  0.000000  9.690000  0.000000  0.538000  6.208500
75%   3.677083  12.500000  18.100000  0.000000  0.624000  6.623500
max   88.976200  108.000000  27.740000  1.000000  0.871000  8.780000

AGE      DIS      RAD      TAX      PTRATIO      B
count  506.000000  506.000000  506.000000  506.000000  506.000000
mean   68.574901  3.795043  9.549407  408.237154  18.455534  356.674032
std    28.148861  2.105710  8.707259  168.537116  2.164946  91.294864
min    2.900000  1.129600  1.000000  187.000000  12.600000  0.320000
25%   45.025000  2.100175  4.000000  279.000000  17.400000  375.377500
50%   77.500000  3.207450  5.000000  330.000000  19.050000  391.440000
75%   94.075000  5.188425  24.000000  666.000000  20.200000  596.225000
max   100.000000  12.126500  24.000000  711.000000  22.000000  596.900000

LSTAT      PRICE
count  506.000000  506.000000
mean   12.653063  22.532806
std    7.141062  9.197104
min    1.730000  5.000000
25%   6.950000  17.025000
50%   11.360000  21.200000
75%   16.955000  25.000000
max   37.970000  50.000000

```

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A4 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A4.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [8]: `#initialize linear regression model  
reg=LinearRegression()`

In [9]: `#split into training-67% & testing data-33%  
X_train, X_test, Y_train, Y_test = train_test_split(bos['PRICE'], test_size = 0.33, random_state=42)  
  
print(X_train.shape)  
print(X_test.shape)  
print(Y_train.shape)  
print(Y_test.shape)`

(339, 14)  
(167, 14)  
(339,)  
(167,)

In [10]: `#train model with our training data  
reg.fit(X_train,Y_train)`

Out[10]: `LinearRegression()`

In [11]: `print(reg.coef_)`

[ 6.05979592e-16 4.09394740e-16 -8.79071121e-16 2.64338231e-15  
3.01350318e-14 1.05777474e-16 -2.95770353e-16 4.68212708e-16  
9.69276742e-16 -1.17961196e-16 5.02744547e-16 -1.26634814e-16  
9.27318118e-16 1.00000000e+00]

In [12]: `#print predictions on our test data  
y_pred=reg.predict(X_test)  
print(y_pred)`

[23.6 32.4 13.6 22.8 16.1 20. 17.8 14. 19.6 16.8 21.5 18.9 7. 21.2  
18.5 29.8 18.8 10.2 50. 14.1 25.2 29.1 12.7 22.4 14.2 13.8 20.3 14.9  
21.7 18.3 23.1 23.8 15. 20.8 19.1 19.4 34.7 19.5 24.4 23.4 19.7 28.2  
50. 17.4 22.6 15.1 13.1 24.2 19.9 24. 18.9 35.4 15.2 26.5 43.5 21.2  
18.4 28.5 23.9 18.5 25. 35.4 31.5 20.2 24.1 20. 13.1 24.8 30.8 12.7  
20. 23.7 10.8 20.6 20.8 5. 20.1 48.5 10.9 7. 20.9 17.2 20.9 9.7  
19.4 29. 16.4 25. 25. 17.1 23.2 10.4 19.6 17.2 27.5 23. 50. 17.9  
9.6 17.2 22.5 21.4 12. 19.9 19.4 13.4 18.2 24.6 21.1 24.7 8.7 27.5  
20.7 36.2 31.6 11.7 39.8 13.9 21.8 23.7 17.6 24.4 8.8 19.2 25.3 20.4  
23.1 37.9 15.6 45.3 15.7 22.6 14.5 18.7 17.8 16.1 20.6 31.6 29.1 15.6  
17.5 22.5 19.4 19.3 8.5 20.6 17. 17.1 14.5 50. 14.3 12.6 28.7 21.2  
19.3 23.1 19.1 25. 33.4 5. 29.6 18.7 21.7 23.1 22.8 21. 48.8]

In [13]: `#actual values  
print(Y_test)`

173 23.6  
274 32.4  
491 13.6  
72 22.8  
452 16.1  
...  
110 21.7  
321 23.1  
265 22.8  
29 21.0  
262 48.8  
Name: PRICE, Length: 167, dtype: float64

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A4 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A4.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [14]:

```
from sklearn.metrics import mean_squared_error
y_pred = reg.predict(X_test)
rmse = np.sqrt(mean_squared_error(Y_test, y_pred))
r2 = round(reg.score(X_test, Y_test), 2)

print("The model performance for training set")
print("-----")
print("Root Mean Squared Error: {}".format(rmse))
print("R^2: {}".format(r2))
print("\n")
```

The model performance for training set

-----

Root Mean Squared Error: 2.294818025447708e-14

In [15]:

```
plt.scatter(Y_test, y_pred)
plt.show()
```

In [15]:

```
plt.scatter(Y_test, y_pred)
plt.show()
```

In [ ]:

# Data Science and Big Data Analytics

Group : A

## Lab Assignment : 05

The screenshot shows a Jupyter Notebook interface with two code cells. The first cell contains code for loading a dataset and splitting it into training and test sets. The second cell contains code for fitting a Logistic Regression model to the training set and predicting results for the test set. The output of the second cell shows the predicted results.

```
In [4]: import numpy as np
import matplotlib.pyplot as plt
import pandas as pd

In [5]: dataset = pd.read_csv('Social_Network_Ads.csv')
X = dataset.iloc[:, [2, 3]].values
y = dataset.iloc[:, 4].values

In [6]: # Splitting the dataset into the Training set and Test set
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25, random_state = 0)

In [7]: # Feature Scaling
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)

In [8]: # Fitting Logistic Regression to the Training set
from sklearn.linear_model import LogisticRegression
log_reg = LogisticRegression(random_state = 0)
log_reg.fit(X_train, y_train)

Out[8]: LogisticRegression(random_state=0)

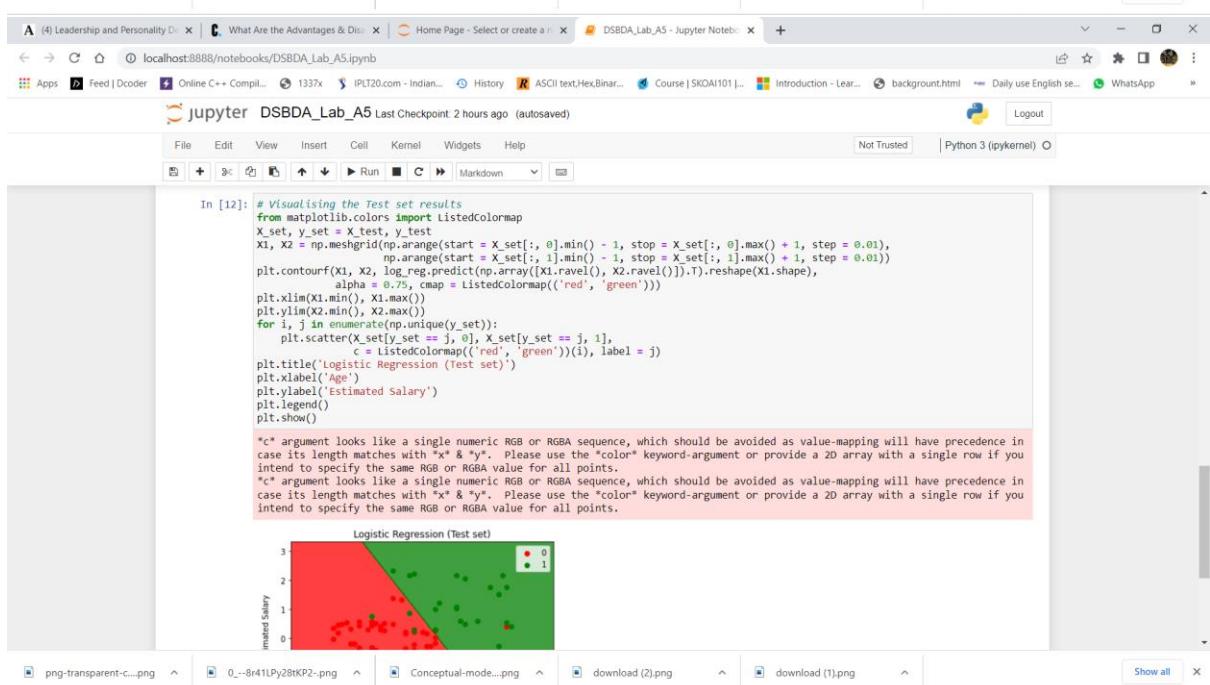
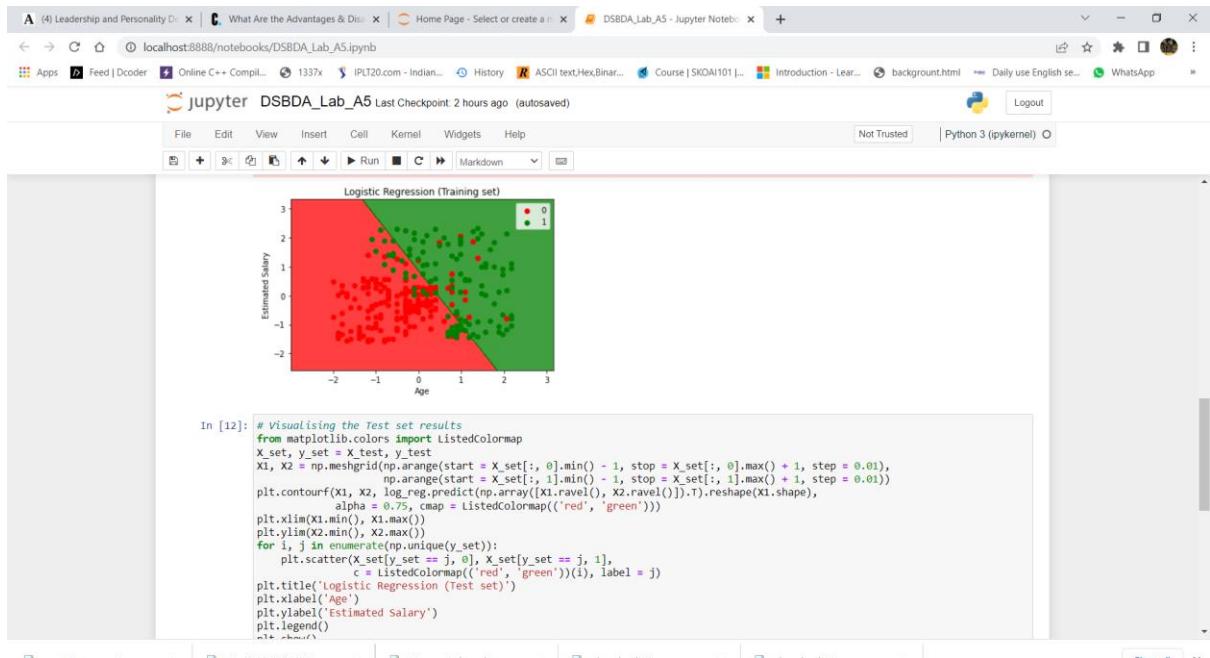
In [9]: # Predicting the Test set results
```

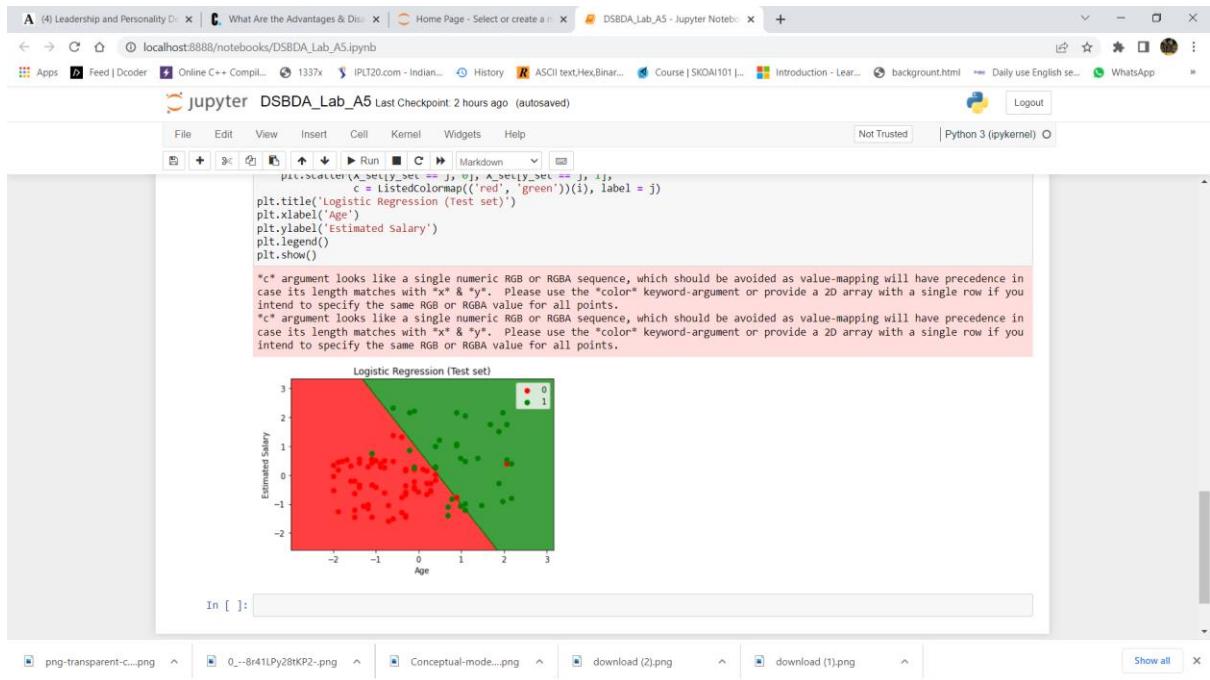
The second part of the screenshot shows another code cell for visualizing the training set results using a scatter plot. The plot shows the relationship between Age and Estimated Salary, with points colored by the predicted class (red for one class, green for the other).

```
In [9]: # Predicting the Test set results
y_pred = log_reg.predict(X_test)

In [10]: # Making the Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)

In [11]: # Visualising the Training set results
from matplotlib.colors import ListedColormap
X_set, y_set = X_train, y_train
X1, X2 = np.meshgrid(np.arange(start = X_set[:, 0].min() - 1, stop = X_set[:, 0].max() + 1, step = 0.01),
                     np.arange(start = X_set[:, 1].min() - 1, stop = X_set[:, 1].max() + 1, step = 0.01))
plt.contourf(X1, X2, log_reg.predict(np.array([X1.ravel(), X2.ravel()]).T).reshape(X1.shape),
             alpha = 0.75, cmap = ListedColormap(('red', 'green')))
plt.xlim(X1.min(), X1.max())
plt.ylim(X2.min(), X2.max())
for i, j in enumerate(np.unique(y_set)):
    plt.scatter(X_set[y_set == j, 0], X_set[y_set == j, 1],
                c = ListedColormap(('red', 'green'))(i), label = j)
plt.title('Logistic Regression (Training set)')
plt.xlabel('Age')
plt.ylabel('Estimated Salary')
plt.legend()
plt.show()
```





# Data Science and Big Data Analytics

Group : A

## Lab Assignment : 06

The image shows two separate sessions of a Jupyter Notebook running on a Windows operating system. Both sessions are titled 'DSBDA\_Lab\_A6 - Jupyter Notebook' and are connected to a Python 3 (ipykernel) kernel.

**Session 1 (Top):**

- In [1]:

```
# Importing the libraries
import numpy as np
import matplotlib.pyplot as plt
import matplotlib.image as mpimg
import pandas as pd
```
- In [2]:

```
dataset = pd.read_csv('iris.csv')
```
- In [3]:

```
dataset.head()
```
- Out[3]:

	sepal.length	sepal.width	petal.length	petal.width	variety
0	5.1	3.5	1.4	0.2	Setosa
1	4.9	3.0	1.4	0.2	Setosa
2	4.7	3.2	1.3	0.2	Setosa
3	4.6	3.1	1.5	0.2	Setosa
4	5.0	3.6	1.4	0.2	Setosa
- In [5]:

```
#Splitting the dataset in independent and dependent variables
X = dataset.iloc[:,4:5].values
y = dataset['variety'].values
```
- In [6]:

```
# Splitting the dataset into the Training set and Test set
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.20, random_state = 82)
```

**Session 2 (Bottom):**

- In [6]:

```
# splitting the dataset into the Training set and Test set
from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.20, random_state = 82)
```
- In [7]:

```
# Feature Scaling to bring the variable in a single scale
from sklearn.preprocessing import StandardScaler
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.transform(X_test)
```
- In [8]:

```
# Fitting Naive Bayes Classification to the Training set with Linear kernel
from sklearn.naive_bayes import GaussianNB
nvclassifier = GaussianNB()
nvclassifier.fit(X_train, y_train)
```
- Out[8]:

```
GaussianNB()
```
- In [9]:

```
# Predicting the Test set results
y_pred = nvclassifier.predict(X_test)
print(y_pred)
```
- [Output]:

```
'Virginica' 'Virginica' 'Setosa' 'Setosa' 'Virginica'
'Virginica' 'Versicolor' 'Versicolor' 'Versicolor' 'Versicolor'
'Virginica' 'Setosa' 'Setosa' 'Setosa' 'Virginica' 'Versicolor'
'Setosa' 'Versicolor' 'Setosa' 'Virginica' 'Setosa' 'Virginica'
'Virginica' 'Versicolor' 'Virginica' 'Setosa' 'Virginica' 'Versicolor'
```
- In [10]:

```
#lets see the actual and predicted value side by side
y_compare = np.vstack((y_test,y_pred)).T
#actual value on the left side and predicted value on the right hand side
#printing the top 5 values
y_compare[:5,:]
```

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A6 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A6.ipynb

jupyter DSBDA\_Lab\_A6 Last Checkpoint: 3 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel) O

```
In [10]: #lets see the actual and predicted value side by side
y_compare = np.vstack((y_test,y_pred)).T
#actual value on the left side and predicted value on the right hand side
#printing the top 5 values
y_compare[:5,:]
```

```
Out[10]: array([['Virginica', 'Virginica'],
   ['Virginica', 'Virginica'],
   ['Setosa', 'Setosa'],
   ['Setosa', 'Setosa'],
   ['Setosa', 'Setosa']], dtype=object)
```

```
In [11]: # Making the Confusion Matrix
from sklearn.metrics import confusion_matrix
cm = confusion_matrix(y_test, y_pred)
print(cm)
```

```
[11  0  0]
[ 0  8  1]
[ 0  1  9]]
```

```
In [12]: #finding accuracy from the confusion matrix.
a = cm.shape
corrPred = 0
falsePred = 0

for row in range(a[0]):
    for c in range(a[1]):
        if row == c:
            corrPred += cm[row,c]
        else:
            falsePred += cm[row,c]
print('Correct predictions:', corrPred)
print('False predictions:', falsePred)
```

[[11 0 0]
 [ 0 8 1]
 [ 0 1 9]]

```
In [12]: #finding accuracy from the confusion matrix.
a = cm.shape
corrPred = 0
falsePred = 0

for row in range(a[0]):
    for c in range(a[1]):
        if row == c:
            corrPred += cm[row,c]
        else:
            falsePred += cm[row,c]
print('Correct predictions:', corrPred)
print('False predictions:', falsePred)
print ('\n\nAccuracy of the Naive Bayes Clasification is: ', corrPred/(cm.sum()))

Correct predictions: 28
False predictions 2

Accuracy of the Naive Bayes Clasification is:  0.9333333333333333
```

In [ ]:

png-transparent-c...png ^ 0\_-8r41Py28tKP2-.png ^ Conceptual-mode...png ^ download (2).png ^ download (1).png ^ Show all X

Data Science and Big Data Analytics

## Group : A

## Lab Assignment : 07

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of... | Home Page - Select or create a new notebook | DSBDA\_Lab\_A7 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A7.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel)

jupyter DSBDA\_Lab\_A7 Last Checkpoint: 3 hours ago (autosaved)

Logout

```

print("This is the unclean version:")
[nltk_data] Downloading package stopwords to C:/Users/Rushi/AppData/Roaming/nltk_data...
[nltk_data] Unzipping corpora/stopwords.zip.

This is the unclean version:
['CSI-DYPIEMR', 'is', 'the', 'Student', 'Chapter', 'of', 'Computer', 'Society', 'of', 'India', 'in', 'Dr.', 'D.', 'Y.', 'Patil', 'J.', 'Pratishthan', 'is', 'Dr.', 'D.', 'Patil', 'Institute', 'of', 'Engineering', 'Management', 'and', 'Research', 'Computer', 'Society', 'India', 'body', 'Computer', 'Society', 'of', 'India', 'is', 'a', 'body', 'of', 'computer', 'professionals', 'in', 'India', 'It', 'was', 'started', 'on', '6', 'March', '1965', 'by', 'a', 'few', 'computer', 'professionals', 'and', 'has', 'now', 'grown', 't', 'o', 'be', 'the', 'national', 'body', 'representing', 'computer', 'professionals', 'and', 'it', 'has', '72', 'chapters', 'across', 'India', '511', 'student', 'branches', 'and', '100,000', 'members', '']

This is the cleaned version:
['CSI-DYPIEMR', 'Student', 'Chapter', 'Computer', 'Society', 'India', 'Dr.', 'D.', 'Y.', 'Patil', 'Pratishthan', 'is', 'Dr.', 'D.', 'Y.', 'Patil', 'Institute', 'Engineering', 'Management', 'Research', 'Computer', 'Society', 'India', 'body', 'computer', 'professionals', 'India', 'It', 'started', '6', 'March', '1965', 'computer', 'professionals', 'grown', 'na', 'tional', 'body', 'representing', 'computer', 'professionals', 'It', '72', 'chapters', 'across', 'India', '511', 'stud', 'branches', '100,000', 'members', '']

In [5]: # A sentence may contain words that convey the same meaning but are written in different forms, taking for example
# verbs in different forms of tenses like 'running', 'ran', 'run', 'runs' ultimately convey the same meaning but
# are written in different forms to suit the different types of tenses, for computer analysis these types of words
# are kept under same section as their base form (which in this case will be "run") and this process is known as stemming

In [6]: from nltk.stem import PorterStemmer
stemmer = nltk.PorterStemmer()
words = ['rain', 'rained', 'raining', 'rains']
stemmed = [stemmer.stem(word) for word in words]
print(stemmed)

['rain', 'rain', 'rain', 'rain']

```

Show all

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of... | Home Page - Select or create a new notebook | DSBDA\_Lab\_A7 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A7.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel)

jupyter DSBDA\_Lab\_A7 Last Checkpoint: 3 hours ago (autosaved)

Logout

```

In [6]: from nltk.stem import PorterStemmer
stemmer = nltk.PorterStemmer()
words = ['rain', 'rained', 'raining', 'rains']
stemmed = [stemmer.stem(word) for word in words]
print(stemmed)

['rain', 'rain', 'rain', 'rain']

In [7]: from nltk import pos_tag
nltk.download('averaged_perceptron_tagger')
tagged = nltk.pos_tag(cleaned_token)
print(tagged)

[nltk_data] Downloading package averaged_perceptron_tagger to
[nltk_data] C:/Users/Rushi/AppData/Roaming/nltk_data...
[nltk_data] Unzipping taggers/averaged_perceptron_tagger.zip.

[('CSI-DYPIEMR', 'JJ'), ('Student', 'NNP'), ('Chapter', 'NNP'), ('Computer', 'NNP'), ('Society', 'NNP'), ('India', 'NNP'), ('Dr.', 'NNP'), ('D.', 'NNP'), ('Y.', 'NNP'), ('Patil', 'NNP'), ('Institute', 'NNP'), ('Engineering', 'NNP'), ('Management', 'NNP'), ('Research', 'NNP'), ('Computer', 'NNP'), ('Society', 'NNP'), ('India', 'NNP'), ('body', 'NN'), ('computer', 'NN'), ('professionals', 'NN'), ('India', 'NNP'), ('.', '.'), ('It', 'PRP'), ('started', 'VBD'), ('6', 'CD'), ('March', 'NNP'), ('1965', 'CD'), ('computer', 'NN'), ('professionals', 'NNS'), ('grew', 'VBP'), ('national', 'JJ'), ('body', 'NN'), ('representing', 'VBG'), ('computer', 'NN'), ('professionals', 'NNS'), ('.', '.'), ('It', 'PRP'), ('72', 'CD'), ('chapters', 'NNS'), ('across', 'IN'), ('India', 'NNP'), ('.', '.'), ('511', 'CD'), ('student', 'NN'), ('branches', 'NNS'), ('.', '.'), ('100,000', 'CD'), ('members', 'NNS'), ('.', '.')]

In [8]: # Lemmatization is the process of finding the form of the related word in the dictionary. Lemmatization does not
# simply chop off inflections. Instead, it uses lexical knowledge bases to get the correct base forms of words.

In [9]: from nltk.stem import WordNetLemmatizer
nltk.download('wordnet')
nltk.download('omw-1.4')

```

Show all

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Leadership? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A7 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A7.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel)

In [9]:

```
from nltk.stem import WordNetLemmatizer
nltk.download('wordnet')
nltk.download('omw-1.4')
lemmatizer = nltk.WordNetLemmatizer()
lemmatized = [lemmatizer.lemmatize(word) for word in cleaned_token]
print(lemmatized)
```

[nltk\_data] Downloading package wordnet to C:/Users/Rushil/AppData/Roaming/nltk\_data...
[nltk\_data] Unzipping corpora/wordnet.zip.
[nltk\_data] Downloading package omw-1.4 to C:/Users/Rushil/AppData/Roaming/nltk\_data...
[nltk\_data] Unzipping corpora/omw-1.4.zip.

['CSI-DYPIEKM', 'Student', 'Chapter', 'Computer', 'Society', 'India', 'Dr.', 'D.', 'Y.', 'Patil', 'Pratishthan', 's', 'Dr.', 'D.', 'Y.', 'Patil', 'Institute', 'Engineering', 'Management', 'Research', 'Computer', 'Society', 'India', 'body', 'computer', 'professional', 'India', 'It', 'started', '6', 'March', '1965', 'computer', 'professional', 'grown', 'national', 'body', 'representing', 'computer', 'professional', 'It', '72', 'chapter', 'across', 'India', 'sii', 'student', 'branch', '100,000', 'member', 't']

In [10]: #TF\_IDF PART 2

In [11]:

```
import pandas as pd
import sklearn as sk
import math
```

In [12]:

```
block_1 = "Our aim is to develop a good work culture among students, a culture where students from various technical backgrounds
block_2 = "Keeping in mind the interest of the IT professionals and computer enthusiasts, CSI works towards making the professionals
#split so each word have their own string
first_block = block_1.split(" ")
second_block = block_2.split(" ")
```

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Leadership? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A7 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A7.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel)

In [12]:

```
block_1 = "Our aim is to develop a good work culture among students, a culture where students from various technical backgrounds
block_2 = "Keeping in mind the interest of the IT professionals and computer enthusiasts, CSI works towards making the professionals
#split so each word have their own string
first_block = block_1.split(" ")
second_block = block_2.split(" ")
#join them to remove common duplicate words
total = set(first_block).union(set(second_block))
print(total)
{'computer', 'sections', 'ensures', 'a', 'good', 'backgrounds', 'various', 'The', 'making', 'projects', 'of', 'professionals.', 'work', 'technical', 'works', 'promotion', 'each', 'at', 'together', 'interest', 'where', 'Technology', 'organizes', 'award', 'To', 'conventions', 'is', 'organized', 'priority', 'the', 'And', 'is', 'top', 'today', 'towards', 'projects', 'with', 'also', 'same', 'on', 'amongst', 'all', 'training', 'develop', 'IT', 'in', 'conferences', 'from', 'objective', 'choice', 'keeping', 'that', 'among', 'come', 'teach', 'it', 'society', 'for', 'enthusiasts', 'together', 'updating', 'future', 'Information', 'regular', 'culture', 'fulfill', 'an', 'other', 'skill', 'mind', 'and', 'regularly', 'time', 'aim', 'our', 'guide', 'collaborate', 'grow', 'are', 'CSI', 'students', 'as', 'to', 'lectures', 'profession', 'this', 'area', 'students', 'professionals'}
```

In [13]:

```
wordDictA = dict.fromkeys(total, 0)
wordDictB = dict.fromkeys(total, 0)

for word in first_block:
    wordDictA[word]+=1

for word in second_block:
    wordDictB[word]+=1
```

In [14]:

```
pd.DataFrame([wordDictA, wordDictB])
```

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A7 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A7.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel)

```
In [14]: pd.DataFrame([wordDictA, wordDictB])
Out[14]:
   computer sections ensures a good backgrounds various The making projects ... CSI students. as to lectures, profession this area students
0          0        0      0  2     1       1    2  0     0      1 ... 0      1  0  2     0      0      0      0      1
1          1        1      1  1     0       0      0  1     1      0 ... 3      0  1  0     1      2  1  1     0
2 rows x 88 columns
```

```
In [15]: def computeTF(wordDict, doc):
    tfDict = {}
    corpusCount = len(doc)
    for word, count in wordDict.items():
        tfDict[word] = count/float(corpusCount)
    return(tfDict)

#running our sentences through the tf function:
tfFirst = computeTF(wordDictA, first_block)
tfSecond = computeTF(wordDictB, second_block)

#Converting to dataframe for visualization
tf = pd.DataFrame([tfFirst, tfSecond])
```

```
In [16]: # Now we'll remove stopwords from the list
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords

stop_words = set(stopwords.words('english'))
```

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A7 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A7.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel)

```
In [16]: # Now we'll remove stopwords from the list
import nltk
nltk.download('stopwords')
from nltk.corpus import stopwords

stop_words = set(stopwords.words('english'))
filtered_sentence = [w for w in wordDictA if not w in stop_words]

print(filtered_sentence)
[nltk_data] Downloading package stopwords to C:\Users\Rushi\AppData\Roaming\nltk_data...
[nltk_data]   C:\Users\Rushi\AppData\Roaming\nltk_data...
[nltk_data]   Package stopwords is already up-to-date!
```

```
In [17]: # Now moving onto the IDF Part
def computeIDF(docList):
    idfDict = {}
    N = len(docList)
    idfDict = dict.fromkeys(docList[0].keys(), 0)
    for word, val in idfDict.items():
        idfDict[word] = math.log10(N / (float(val) + 1))
    return(idfDict)

#inputting our sentences in the Log file
idfs = computeIDF([wordDictA, wordDictB])
```

In [17]: # Now moving onto the IDF Part

```
def computeIDF(docList):
    idfDict = {}
    N = len(docList)
    idfDict = dict.fromkeys(docList[0].keys(), 0)
    for word, val in idfDict.items():
        idfDict[word] = math.log10(N / (float(val) + 1))
    return(idfDict)

#inputting our sentences in the log file
idfs = computeIDF([wordDictA, wordDictB])
```

In [18]: # Now we'll implement the IDF formula

```
def computeTFIDF(tfBow, idfs):
    tfidf = {}
    for word, val in tfBow.items():
        tfidf[word] = val*idfs[word]
    return(tfidf)

#running our two sentences through the IDF:
idfFirst = computeTFIDF(tfFirst, idfs)
idfSecond = computeTFIDF(tfSecond, idfs)

#putting it in a dataframe
idf = pd.DataFrame([idfFirst, idfSecond])

print(idf)
```

	computer	sections	ensures	a	good	backgrounds	various
0	0.000000	0.000000	0.000000	0.017202	0.008601	0.008601	0.017202
1	0.003859	0.003859	0.003859	0.003859	0.000000	0.000000	0.000000

	The	making	projects	...	CSI	students,	as	to
0	0.000000	0.000000	0.008601	...	0.000000	0.008601	0.000000	0.017202
1	0.003859	0.003859	0.000000	...	0.011578	0.000000	0.003859	0.000000

[2 rows x 88 columns]

In [19]: # Above way was the generic/formalistic way of implementing TFIDF, This process can be made way more simpler by using sklearn library, example given below

```
from sklearn.feature_extraction.text import TfidfVectorizer
#make sure all words are in lowercase
version_1 = "Developing a competitive culture where the students polish technical and professional attributes, gain experience and version_2 = "Personalized career guidance, Regular Logic and aptitude building activities, Industrial level project collaboration"
#calling the Tfidfvectorizer
vectorize= TfidfVectorizer()
#fitting the model and passing our sentences right away:
response= vectorize.fit_transform([version_1.lower(), version_2.lower()])
```

In [20]: print(response)

(0, 61)	0.13915271943780658
(0, 31)	0.13915271943780658
(0, 40)	0.13915271943780658
(0, 4)	0.13915271943780658
(0, 48)	0.13915271943780658

A (4) Leadership and Personality D x C. What Are the Advantages & Disadvantages of Leadership? x Home Page - Select or create a new notebook x DSBDA\_Lab\_A7 - Jupyter Notebook x +

localhost:8888/notebooks/DSBDA\_Lab\_A7.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [20]: `print(response)`

```
(0, 61) 0.13915271943780658
(0, 31) 0.13915271943780658
(0, 48) 0.13915271943780658
(0, 4) 0.13915271943780658
(0, 18) 0.13915271943780658
(0, 30) 0.13915271943780658
(0, 58) 0.13915271943780658
(0, 29) 0.13915271943780658
(0, 7) 0.13915271943780658
(0, 46) 0.13915271943780658
(0, 54) 0.13915271943780658
(0, 9) 0.13915271943780658
(0, 60) 0.13915271943780658
(0, 67) 0.13915271943780658
(0, 28) 0.13915271943780658
(0, 65) 0.13915271943780658
(0, 23) 0.13915271943780658
(0, 59) 0.27830543887561315
(0, 24) 0.2970249178760062
(0, 53) 0.13915271943780658
(0, 44) 0.13915271943780658
(0, 3) 0.13915271943780658
(0, 62) 0.13915271943780658
(0, 64) 0.13915271943780658
: :
(1, 21) 0.16649349332910351
(1, 37) 0.16649349332910351
(1, 41) 0.16649349332910351
(1, 26) 0.16649349332910351
(1, 0) 0.16649349332910351
(1, 35) 0.16649349332910351
(1, 1) 0.16649349332910351
(1, 66) 0.16649349332910351
(1, 38) 0.16649349332910351
(1, 12) 0.16649349332910351
(1, 47) 0.16649349332910351
(1, 35) 0.16649349332910351
(1, 32) 0.16649349332910351
(1, 2) 0.16649349332910351
(1, 10) 0.3298698665820703
(1, 6) 0.16649349332910351
(1, 36) 0.16649349332910351
(1, 49) 0.16649349332910351
(1, 27) 0.16649349332910351
(1, 11) 0.16649349332910351
(1, 42) 0.16649349332910351
(1, 24) 0.11846149176425531
(1, 52) 0.11846149176425531
(1, 5) 0.35538447529276596
(1, 57) 0.11846149176425531
```

In [ ]:

png-transparent-c...png ^ 0\_-8r41LPy28tKP2-.png ^ Conceptual-mode...png ^ download (2).png ^ download (1).png ^ Show all x

# Data Science and Big Data Analytics

Group : A

## Lab Assignment : 08

A (4) Leadership and Personality D... | C, What Are the Advantages & Disadvantages of Big Data? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A8 - Jupyter Notebook | +

localhost:8888/notebooks/DSBDA\_Lab\_A8.ipynb

Apps Feed | Docker Online C++ Compiler 1337x IPLT20.com - Indian... History ASCII text,Hex,Binary Course | SKOAI101 | ... Introduction - Lear... background.html Daily use English se... WhatsApp Logout

jupyter DSBDA\_Lab\_A8 Last Checkpoint: 3 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help Not Trusted Python 3 (ipykernel) O

In [1]:

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

In [2]:

```
titanic = sns.load_dataset('titanic')
```

In [3]:

```
titanic
```

Out[3]:

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True	Nan	Southampton	no	False
1	1	1	female	38.0	1	0	71.2833	C	First	woman	False	C	Cherbourg	yes	False
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	False	Nan	Southampton	yes	True
3	1	1	female	35.0	1	0	53.1000	S	First	woman	False	C	Southampton	yes	False
4	0	3	male	35.0	0	0	8.0500	S	Third	man	True	Nan	Southampton	no	True
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
886	0	2	male	27.0	0	0	13.0000	S	Second	man	True	Nan	Southampton	no	True
887	1	1	female	19.0	0	0	30.0000	S	First	woman	False	B	Southampton	yes	True
888	0	3	female	NaN	1	2	23.4500	S	Third	woman	False	Nan	Southampton	no	False
889	1	1	male	26.0	0	0	30.0000	C	First	man	True	C	Cherbourg	yes	True
890	0	3	male	32.0	0	0	7.7500	Q	Third	man	True	Nan	Queenstown	no	True

891 rows × 15 columns

In [4]:

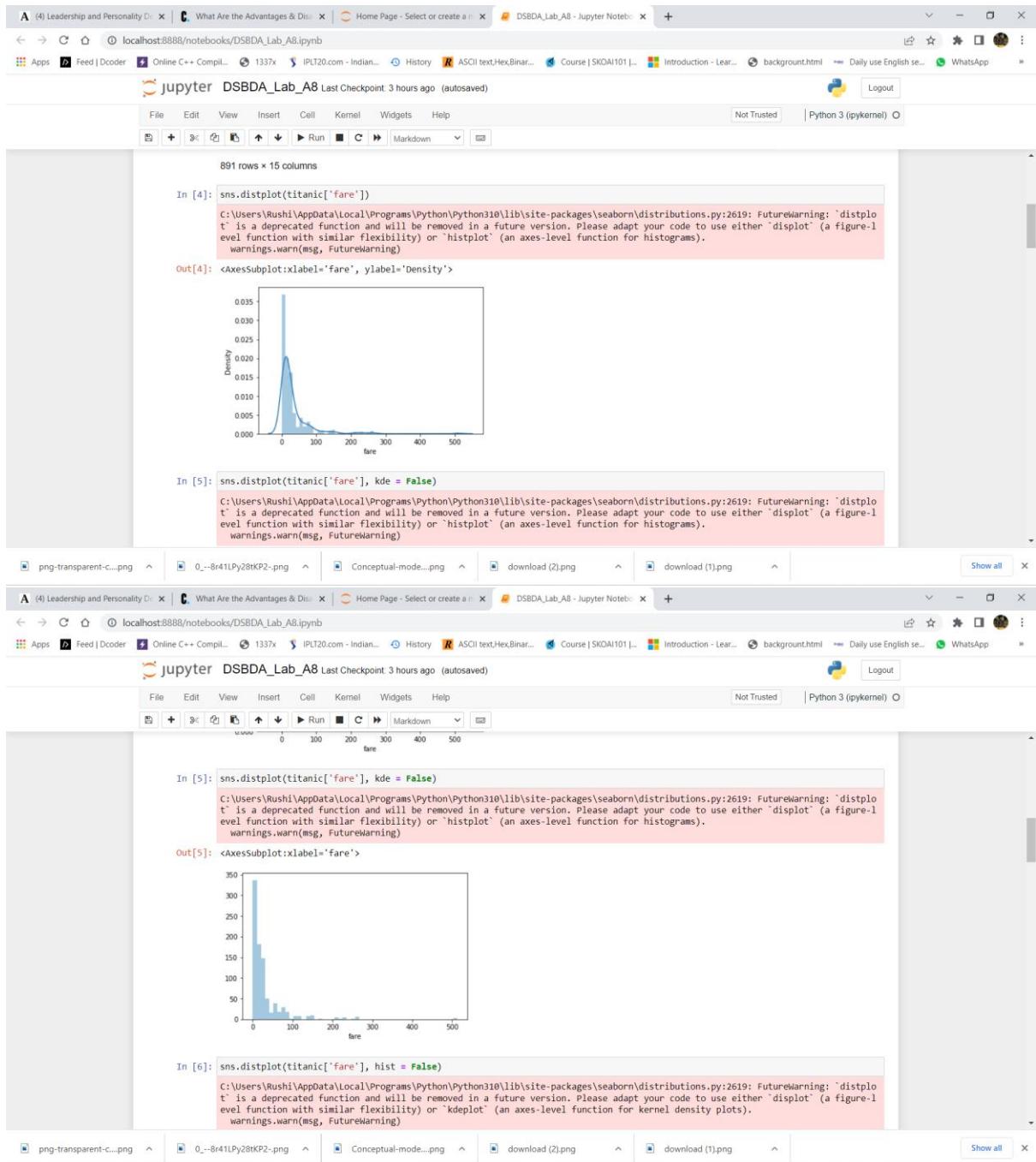
```
sns.distplot(titanic['fare'])
```

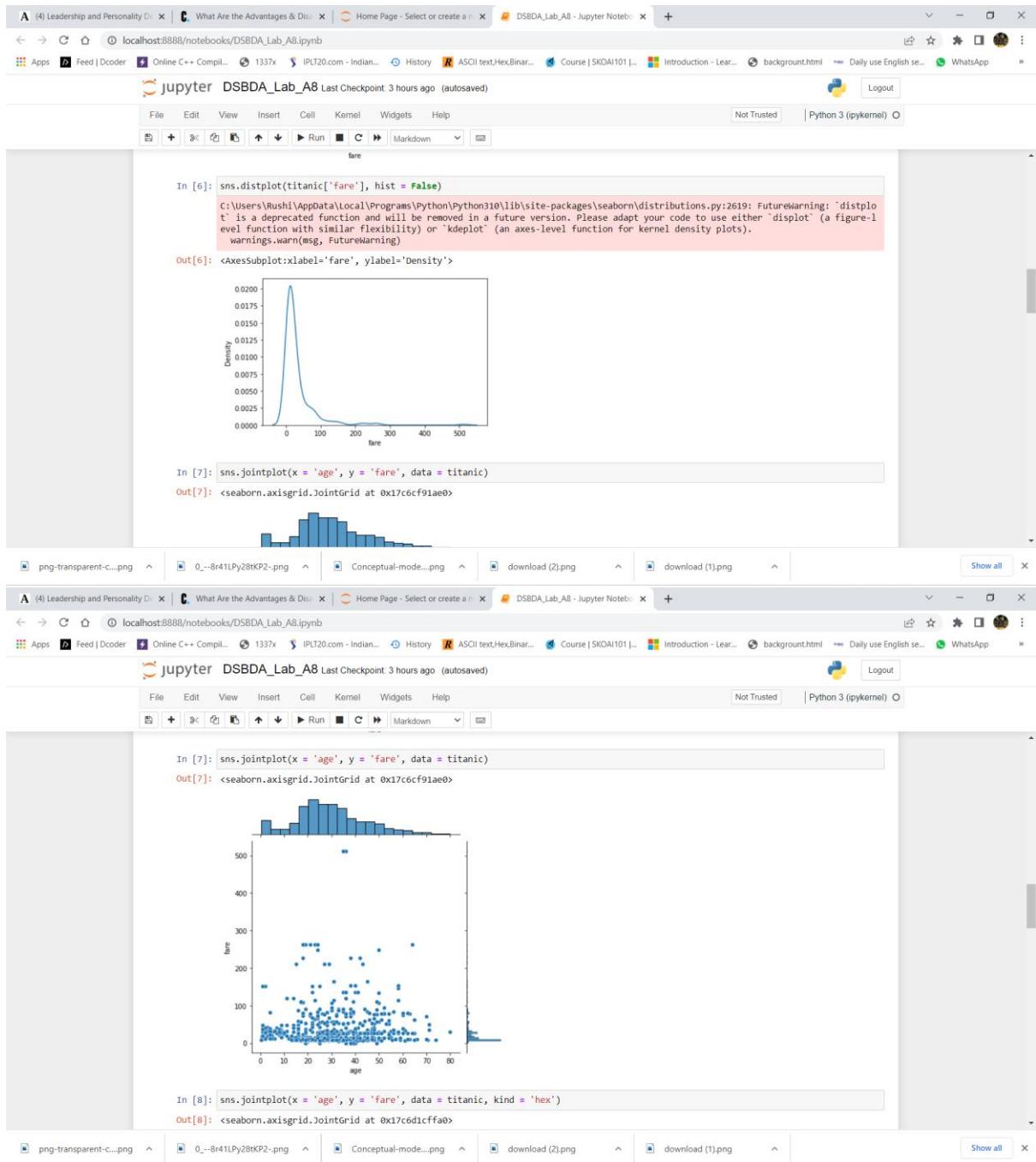
C:\Users\Rushi\AppData\Local\Programs\Python\Python310\lib\site-packages\seaborn\distributions.py:2619: FutureWarning: 'distplot' is a deprecated function and will be removed in a future version. Please adapt your code to use either 'displot' (a figure-level function with similar flexibility) or 'histplot' (an axes-level function for histograms).  
warnings.warn(msg, FutureWarning)

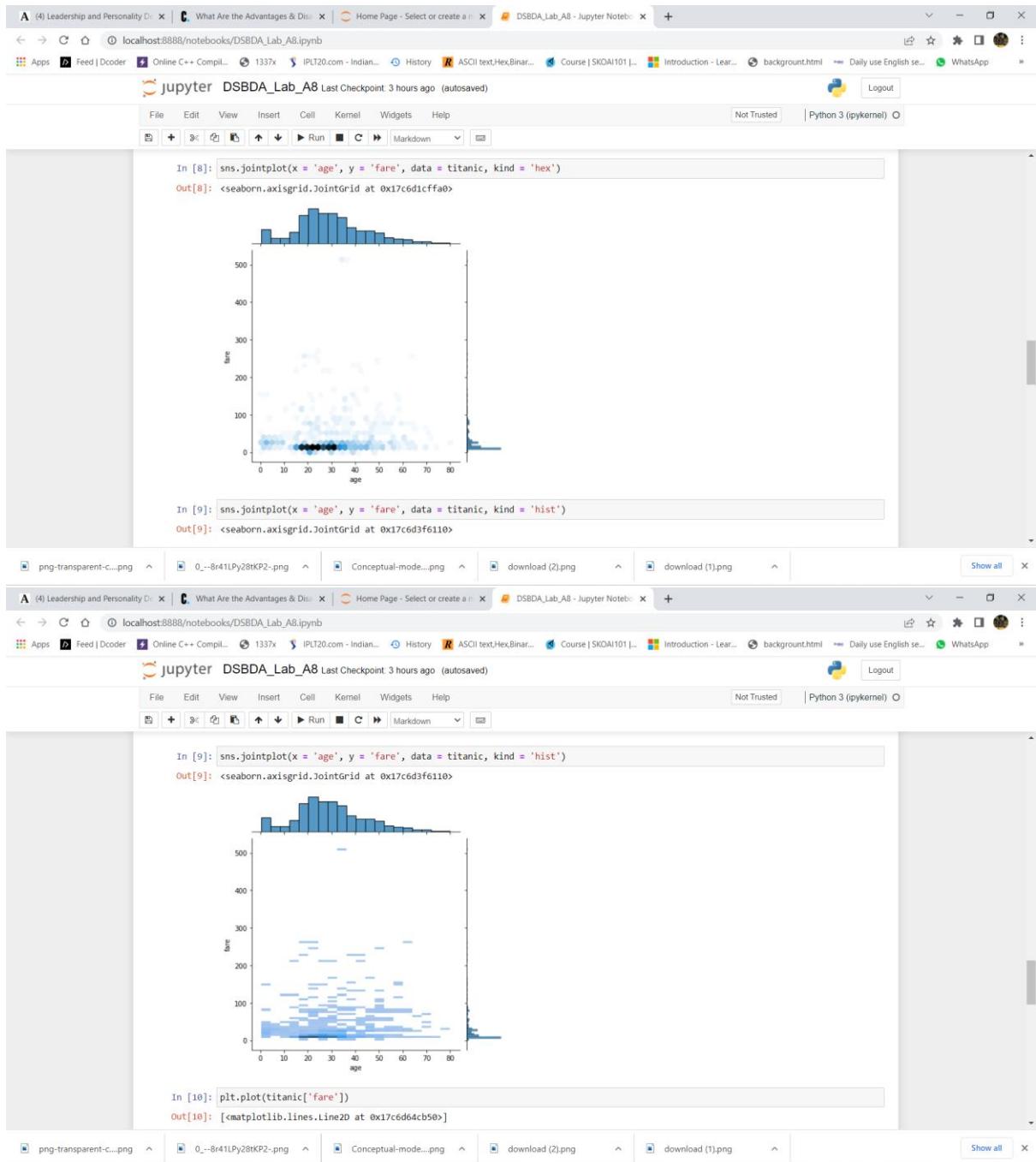
Out[4]:

```
<AxesSubplot:xlabel='fare', ylabel='Density'>
```

png-transparent-c...png ^ 0\_--8r41LPy28tKP2-.png ^ Conceptual-mode...png ^ download (2).png ^ download (1).png ^ Show all X







A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A8 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A8.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [10]: `plt.plot(titanic['fare'])`

Out[10]: [`[<matplotlib.lines.Line2D at 0x17c6d64cb50>]`]

In [11]: `plt.hist(titanic['fare'])`

Out[11]: (`array([732, 106, 31, 2, 11, 6, 0, 0, 0, 3.]), array([ 0, 51.23292, 102.46584, 153.69876, 204.93168, 256.1646, 307.39752, 358.63044, 409.86336, 461.09628, 512.3292]), <BarContainer object of 10 artists>`)

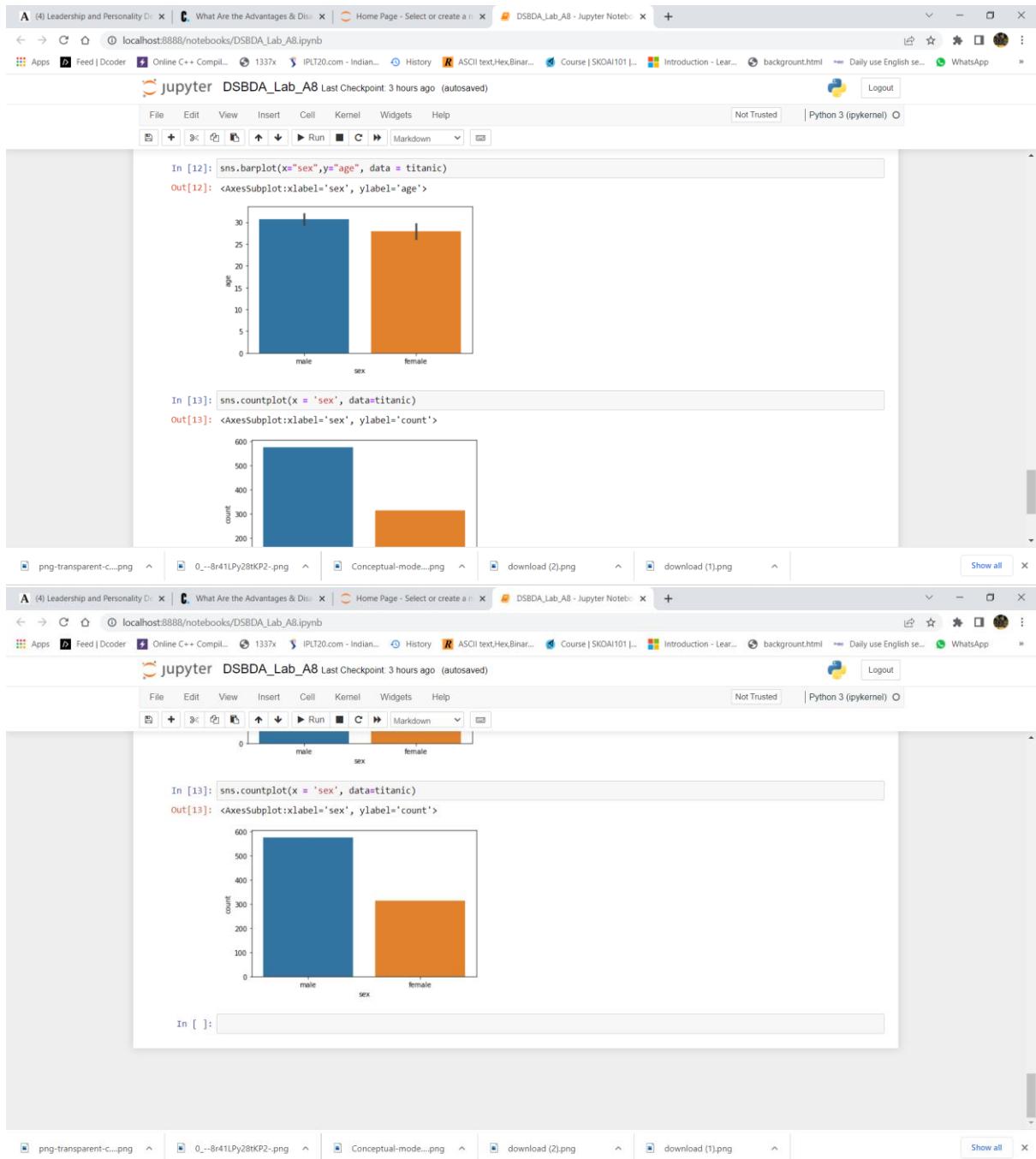
In [11]: `plt.hist(titanic['fare'])`

Out[11]: (`array([732, 106, 31, 2, 11, 6, 0, 0, 0, 3.]), array([ 0, 51.23292, 102.46584, 153.69876, 204.93168, 256.1646, 307.39752, 358.63044, 409.86336, 461.09628, 512.3292]), <BarContainer object of 10 artists>`)

In [12]: `sns.barplot(x="sex",y="age", data = titanic)`

Out[12]: [`<AxesSubplot:xlabel='sex', ylabel='age'>`]

Show all



# Data Science and Big Data Analytics

Group : A

Lab Assignment : 09

jupyter DSBDA\_Lab\_A9 Last Checkpoint: 3 hours ago (autosaved)

In [2]:

```
import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

In [4]:

```
titanic = sns.load_dataset('titanic')
```

In [5]:

```
titanic
```

Out[5]:

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True	NaN	Southampton	no	False
1	1	1	female	38.0	1	0	71.2833	C	First	woman	False	C	Cherbourg	yes	False
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	False	NaN	Southampton	yes	True
3	1	1	female	35.0	1	0	53.1000	S	First	woman	False	C	Southampton	yes	False
4	0	3	male	35.0	0	0	8.0500	S	Third	man	True	NaN	Southampton	no	True
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
886	0	2	male	27.0	0	0	13.0000	S	Second	man	True	NaN	Southampton	no	True
887	1	1	female	19.0	0	0	30.0000	S	First	woman	False	B	Southampton	yes	True
888	0	3	female	NaN	1	2	23.4500	S	Third	woman	False	NaN	Southampton	no	False
889	1	1	male	26.0	0	0	30.0000	C	First	man	True	C	Cherbourg	yes	True
890	0	3	male	32.0	0	0	7.7500	Q	Third	man	True	NaN	Queenstown	no	True

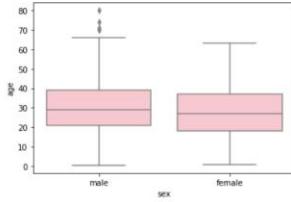
891 rows × 15 columns

In [6]:

```
sns.boxplot(x = 'sex', y = 'age', data = titanic, color = 'pink')
```

Out[6]:

```
<AxesSubplot:xlabel='sex', ylabel='age'>
```



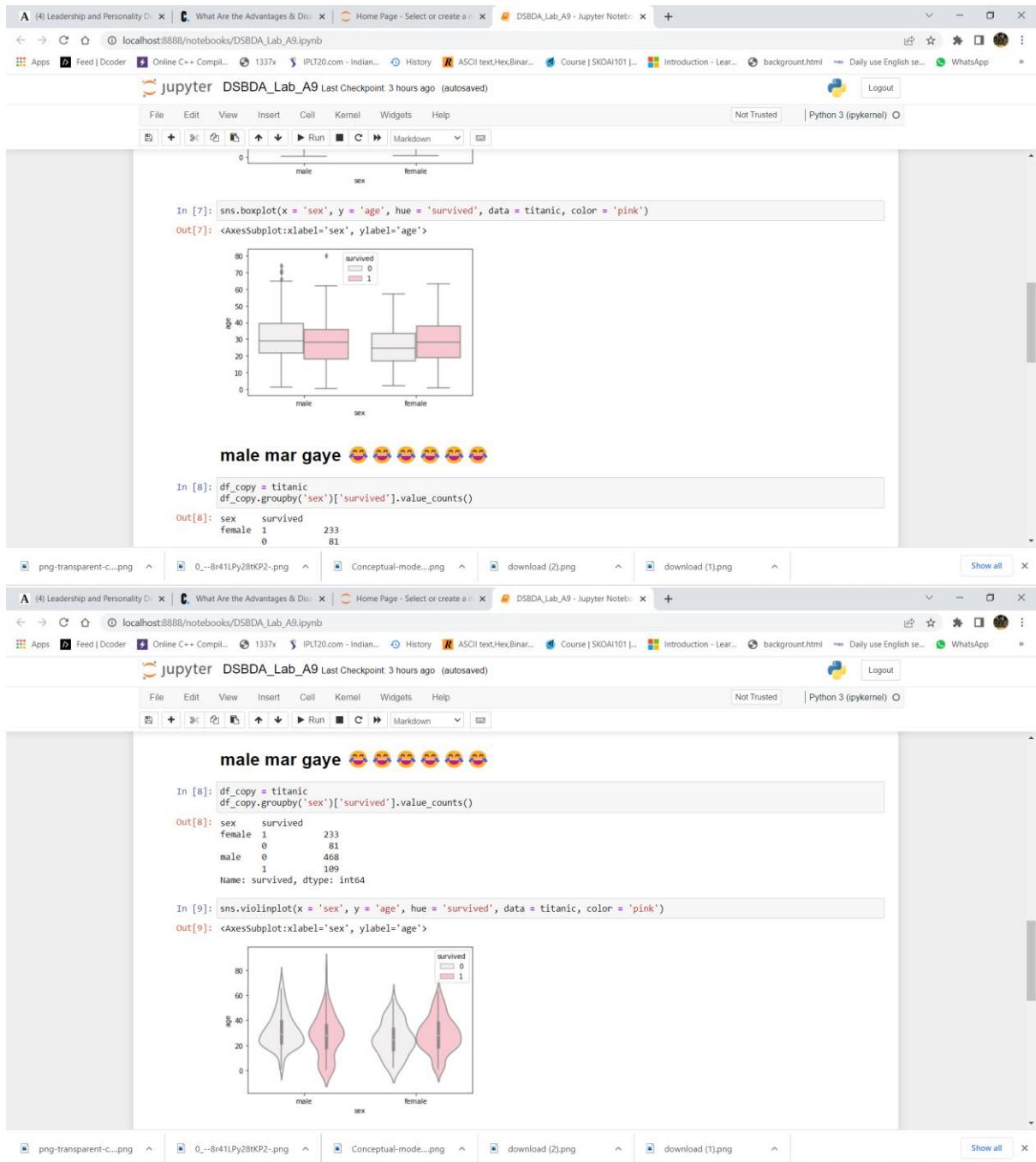
In [7]:

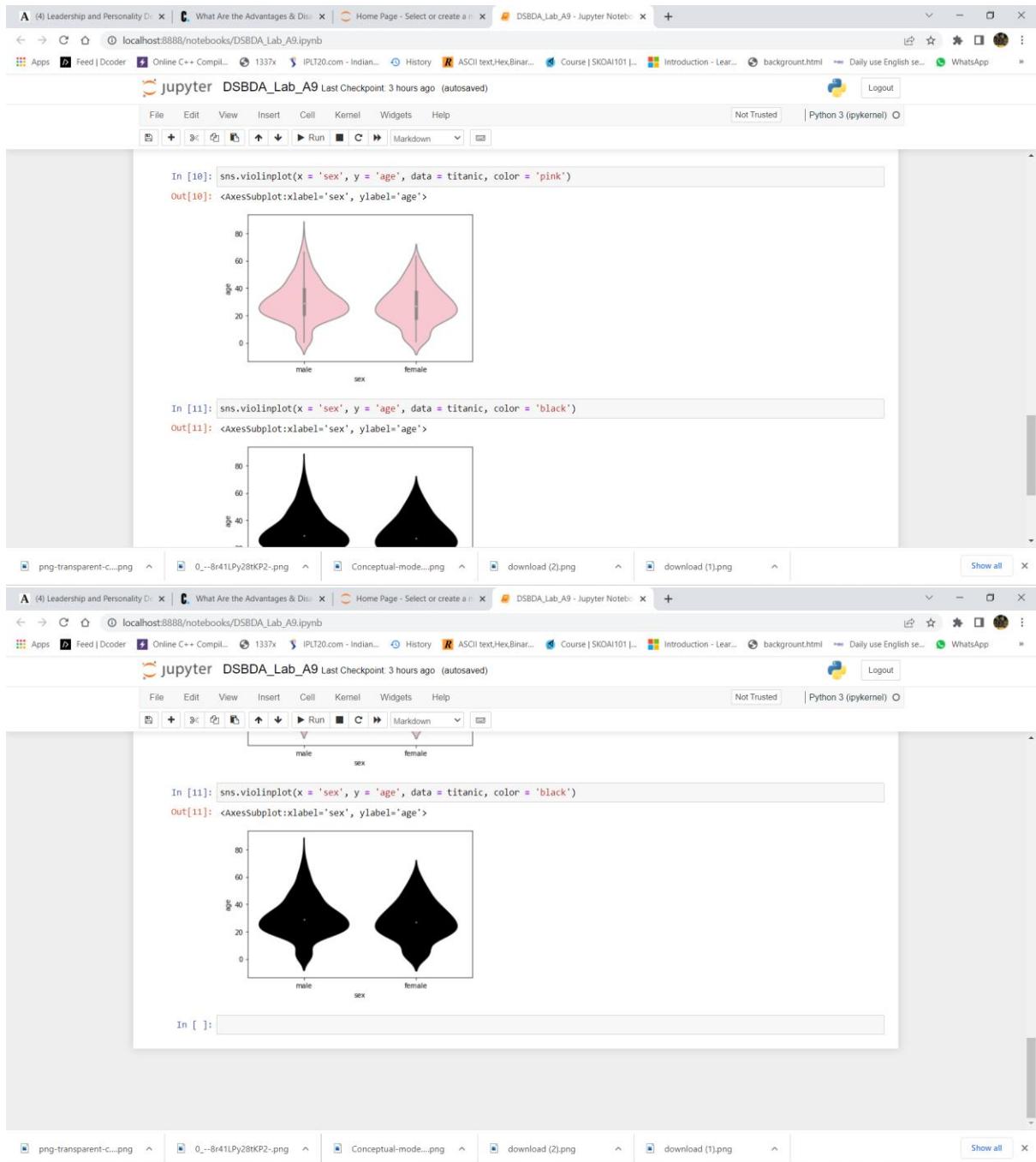
```
sns.boxplot(x = 'sex', y = 'age', hue = 'survived', data = titanic, color = 'pink')
```

Out[7]:

```
<AxesSubplot:xlabel='sex', ylabel='age'>
```







# Data Science and Big Data Analytics

Group : A

Lab Assignment : 10

A (4) Leadership and Personality D... x C, What Are the Advantages & Disadvantages of Big Data? x Home Page - Select or create a new notebook x DSBDA\_Lab\_A9 - Jupyter Notebook x +

localhost:8888/notebooks/DSBDA\_Lab\_A9.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel) O

jupyter DSBDA\_Lab\_A9 Last Checkpoint: 3 hours ago (autosaved)

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel) O

## Data Visualisation II

```
In [2]: import pandas as pd
import numpy as np
import seaborn as sns
import matplotlib.pyplot as plt
```

```
In [4]: titanic = sns.load_dataset('titanic')
```

```
In [5]: titanic
```

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True	Nan	Southampton	no	False
1	1	1	female	38.0	1	0	71.2833	C	First	woman	False	C	Cherbourg	yes	False
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	False	Nan	Southampton	yes	True
3	1	1	female	35.0	1	0	53.1000	S	First	woman	False	C	Southampton	yes	False
4	0	3	male	35.0	0	0	8.0500	S	Third	man	True	Nan	Southampton	no	True
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
886	0	2	male	27.0	0	0	13.0000	S	Second	man	True	Nan	Southampton	no	True
887	1	1	female	19.0	0	0	30.0000	S	First	woman	False	B	Southampton	yes	True
888	0	3	female	NaN	1	2	23.4500	S	Third	woman	False	Nan	Southampton	no	False
889	1	1	male	26.0	0	0	30.0000	C	First	man	True	C	Cherbourg	yes	True
890	0	3	male	32.0	0	0	7.7500	Q	Third	man	True	Nan	Queenstown	no	True

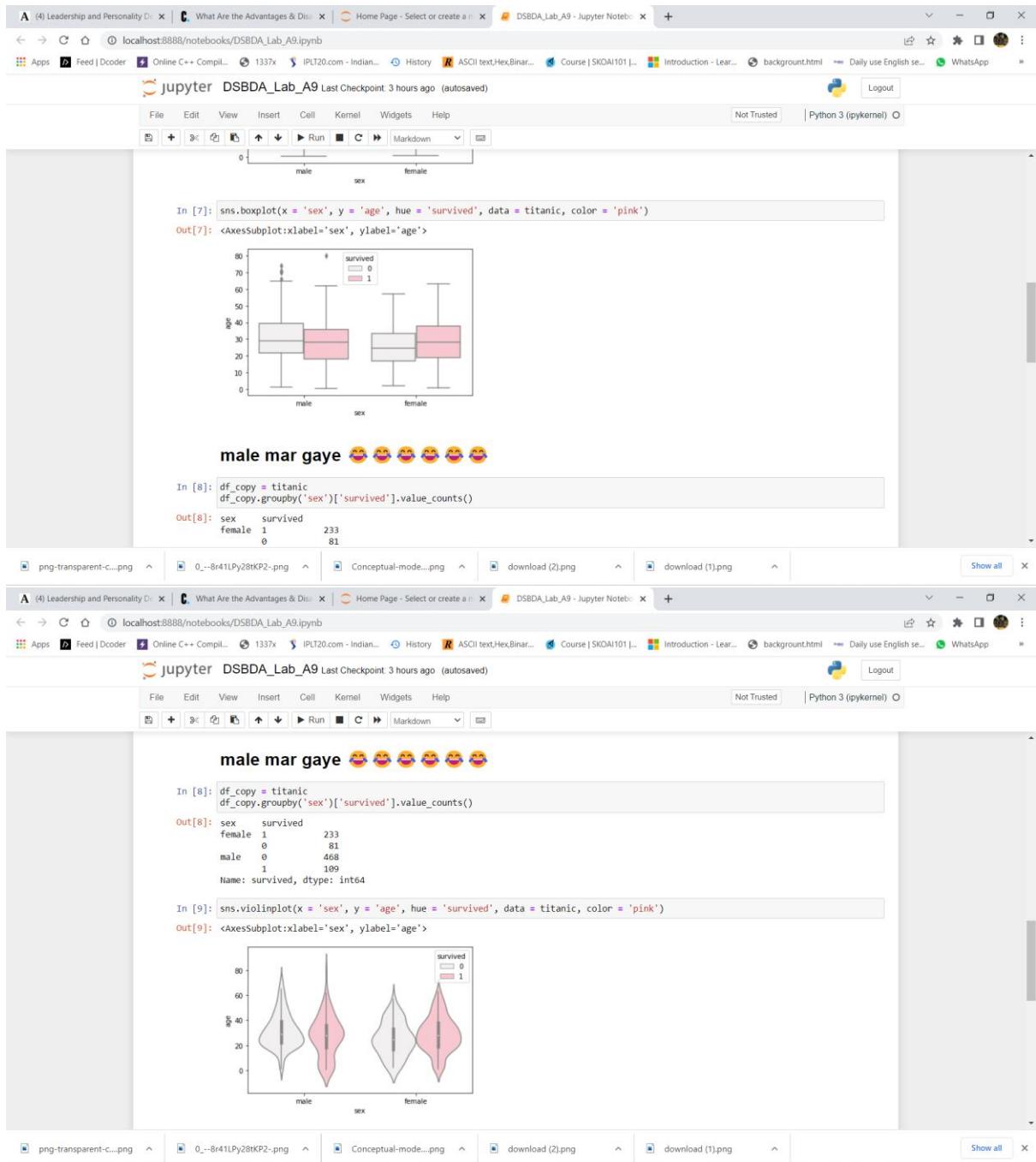
891 rows × 15 columns

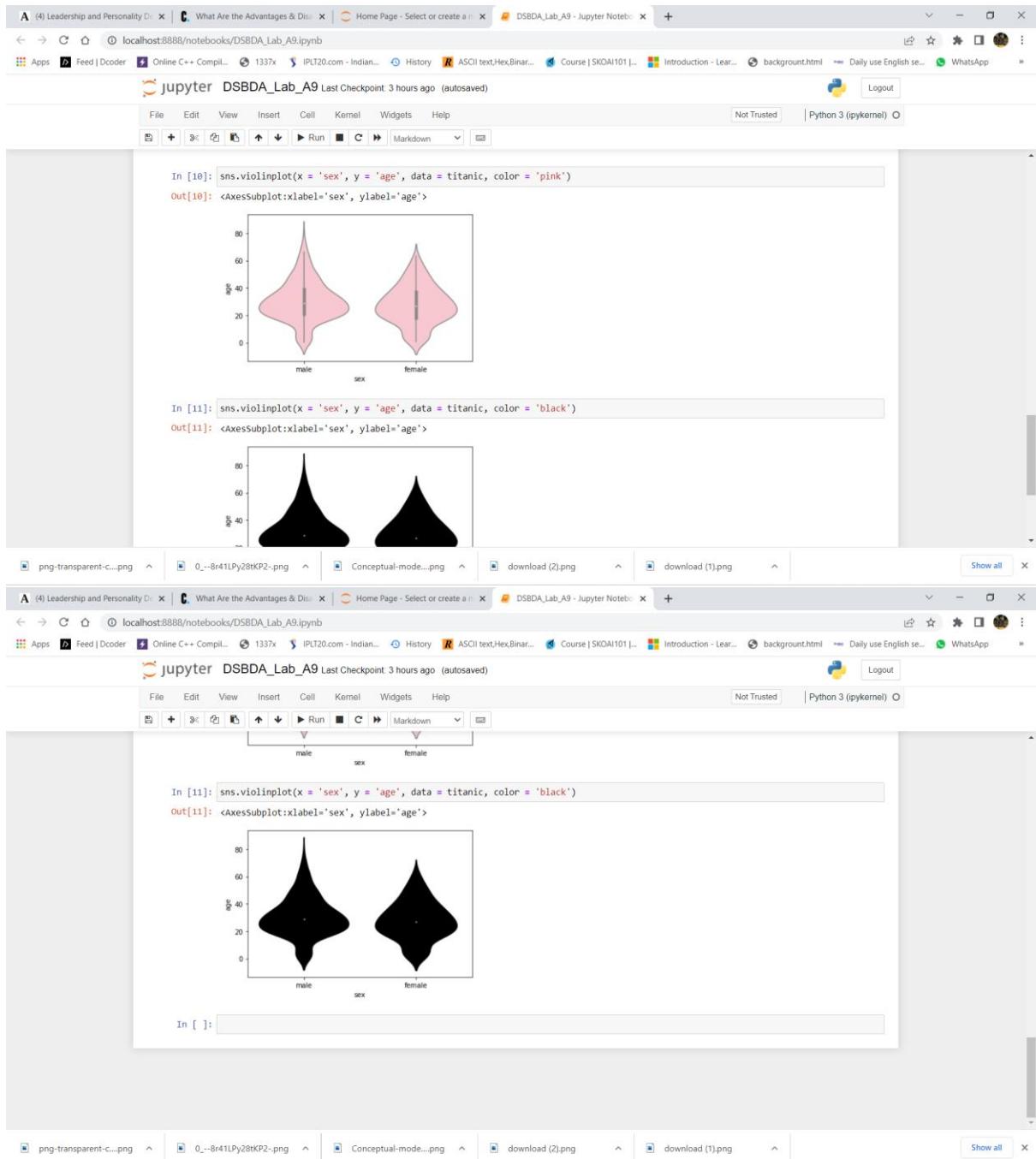
```
In [6]: sns.boxplot(x = 'sex', y = 'age', data = titanic, color = 'pink')
```

```
Out[6]: <AxesSubplot:xlabel='sex', ylabel='age'>
```

```
In [7]: sns.boxplot(x = 'sex', y = 'age', hue = 'survived', data = titanic, color = 'pink')
```

```
Out[7]: <AxesSubplot:xlabel='sex', ylabel='age'>
```





**Data Visualization III**

```
In [2]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [3]: data1 = pd.read_csv("Iris.csv")
data1.head()

Out[3]:   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 [4]: print(data1.columns)
Index(['Id', 'SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm',
       'Species'],
      dtype='object')

In [5]: #anotherway
column = list(data1)
print(column)

In [5]: #anotherway
column = list(data1)
print(column)

In [6]: data1.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 [7]: data1.dtypes

Out[7]: Id          int64
SepallengthCm  float64
SepalwidthCm   float64
PetallengthCm  float64
PetalwidthCm   float64
Species        object
dtype: object

In [8]: data1.hist()
```

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A10 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A10.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [8]: `data1.hist()`

Out[8]: array([[<AxesSubplot:title={'center':'Id'}>,<AxesSubplot:title={'center':'SepallengthCm'}>,<AxesSubplot:title={'center':'SepalwidthCm'}>,<AxesSubplot:title={'center':'PetallengthCm'}>,<AxesSubplot:title={'center':'PetalwidthCm'}>],<AxesSubplot:title={'center': 'PetalwidthCm'}>],<AxesSubplot:>]],  
dtype=object)

In [9]: `fig, axes = plt.subplots(2, 2, figsize=(16, 8))`

axes[0,0].set\_title("Distribution of First Column")  
axes[0,0].hist(data1["SepallengthCm"]);  
axes[0,1].set\_title("Distribution of Second Column")  
axes[0,1].hist(data1["SepalwidthCm"]);

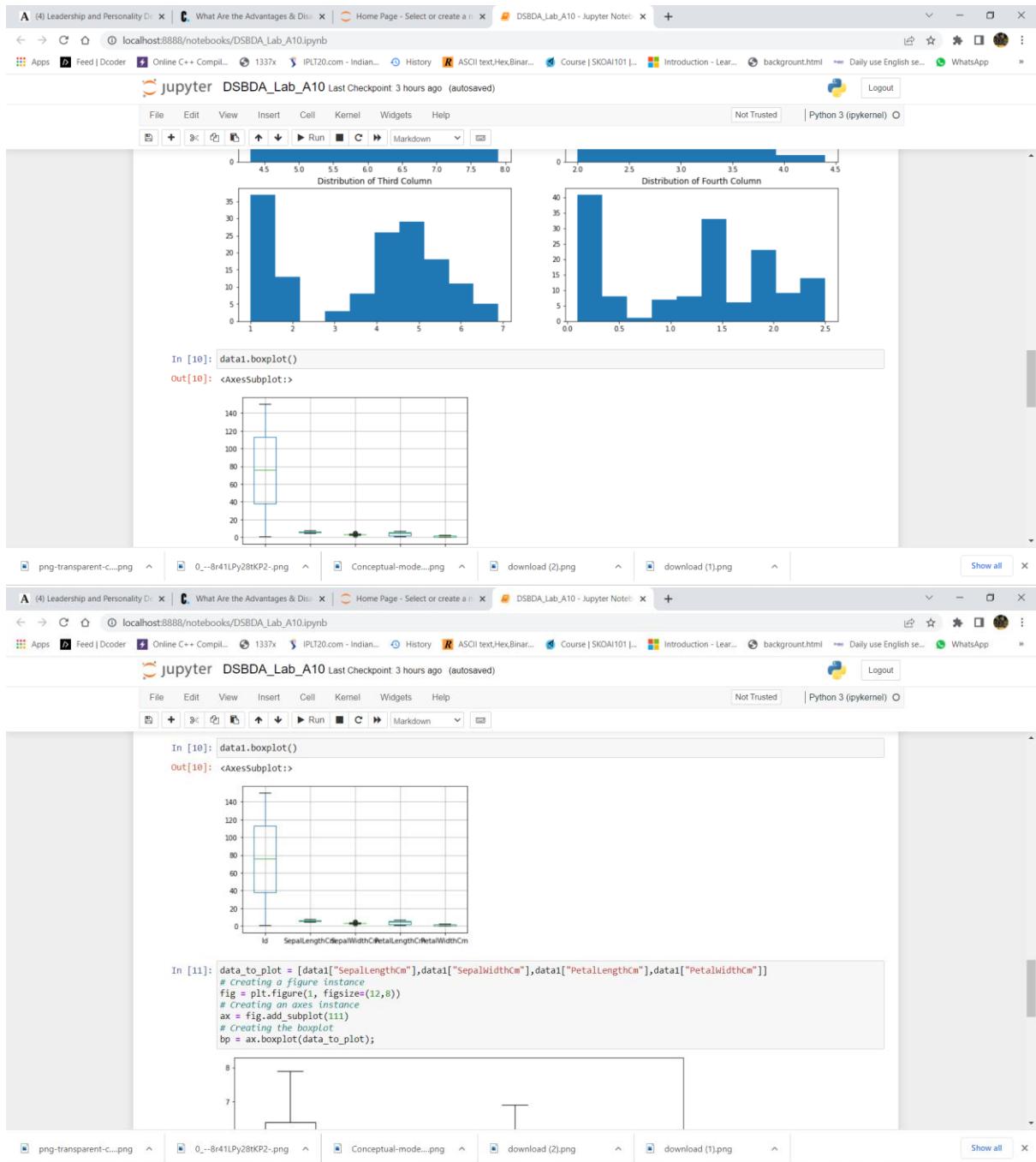
A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A10 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A10.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [9]: `fig, axes = plt.subplots(2, 2, figsize=(16, 8))`

axes[0,0].set\_title("Distribution of First Column")  
axes[0,0].hist(data1["SepallengthCm"]);  
axes[0,1].set\_title("Distribution of Second Column")  
axes[0,1].hist(data1["SepalwidthCm"]);  
axes[1,0].set\_title("Distribution of Third Column")  
axes[1,0].hist(data1["PetallengthCm"]);  
axes[1,1].set\_title("Distribution of Fourth Column")  
axes[1,1].hist(data1["PetalwidthCm"]);



In [11]:

```
data_to_plot = [data1["SepalLengthCm"],data1["SepalWidthCm"],data1["PetalLengthCm"],data1["PetalWidthCm"]]
# Creating a figure instance
fig = plt.figure(1, figsize=(12,8))
# Creating an axes instance
ax = fig.add_subplot(111)
# creating the boxplot
bp = ax.boxplot(data_to_plot);
```

In [12]:

```
sns.boxplot(data1['SepalWidthCm'])
```

C:\Users\rushil\AppData\Local\Programs\Python\Python310\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.  
warnings.warn(  
Out[12]: <AxesSubplot:xlabel='SepalWidthCm'>

In [13]:

```
print(np.where(data1['SepalWidthCm']>4.0))  
(array([15, 32, 33], dtype=int64),)
```

In [ ]:

# Data Science and Big Data Analytics

Group : B2

Lab Assignment : 12

A (4) Leadership and Personality D... | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A9 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A9.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel)

jupyter DSBDA\_Lab\_A9 Last Checkpoint: 3 hours ago (autosaved)

Logout

## Data Visualisation II

```
In [2]: import pandas as pd  
import numpy as np  
import seaborn as sns  
import matplotlib.pyplot as plt
```

```
In [4]: titanic = sns.load_dataset('titanic')
```

```
In [5]: titanic
```

```
Out[5]:
```

	survived	pclass	sex	age	sibsp	parch	fare	embarked	class	who	adult_male	deck	embark_town	alive	alone
0	0	3	male	22.0	1	0	7.2500	S	Third	man	True	Nan	Southampton	no	False
1	1	1	female	38.0	1	0	71.2833	C	First	woman	False	C	Cherbourg	yes	False
2	1	3	female	26.0	0	0	7.9250	S	Third	woman	False	Nan	Southampton	yes	True
3	1	1	female	35.0	1	0	53.1000	S	First	woman	False	C	Southampton	yes	False
4	0	3	male	35.0	0	0	8.0500	S	Third	man	True	Nan	Southampton	no	True
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...
886	0	2	male	27.0	0	0	13.0000	S	Second	man	True	Nan	Southampton	no	True
887	1	1	female	19.0	0	0	30.0000	S	First	woman	False	B	Southampton	yes	True
888	0	3	female	NaN	1	2	23.4500	S	Third	woman	False	Nan	Southampton	no	False
889	1	1	male	26.0	0	0	30.0000	C	First	man	True	C	Cherbourg	yes	True
890	0	3	male	32.0	0	0	7.7500	Q	Third	man	True	Nan	Queenstown	no	True

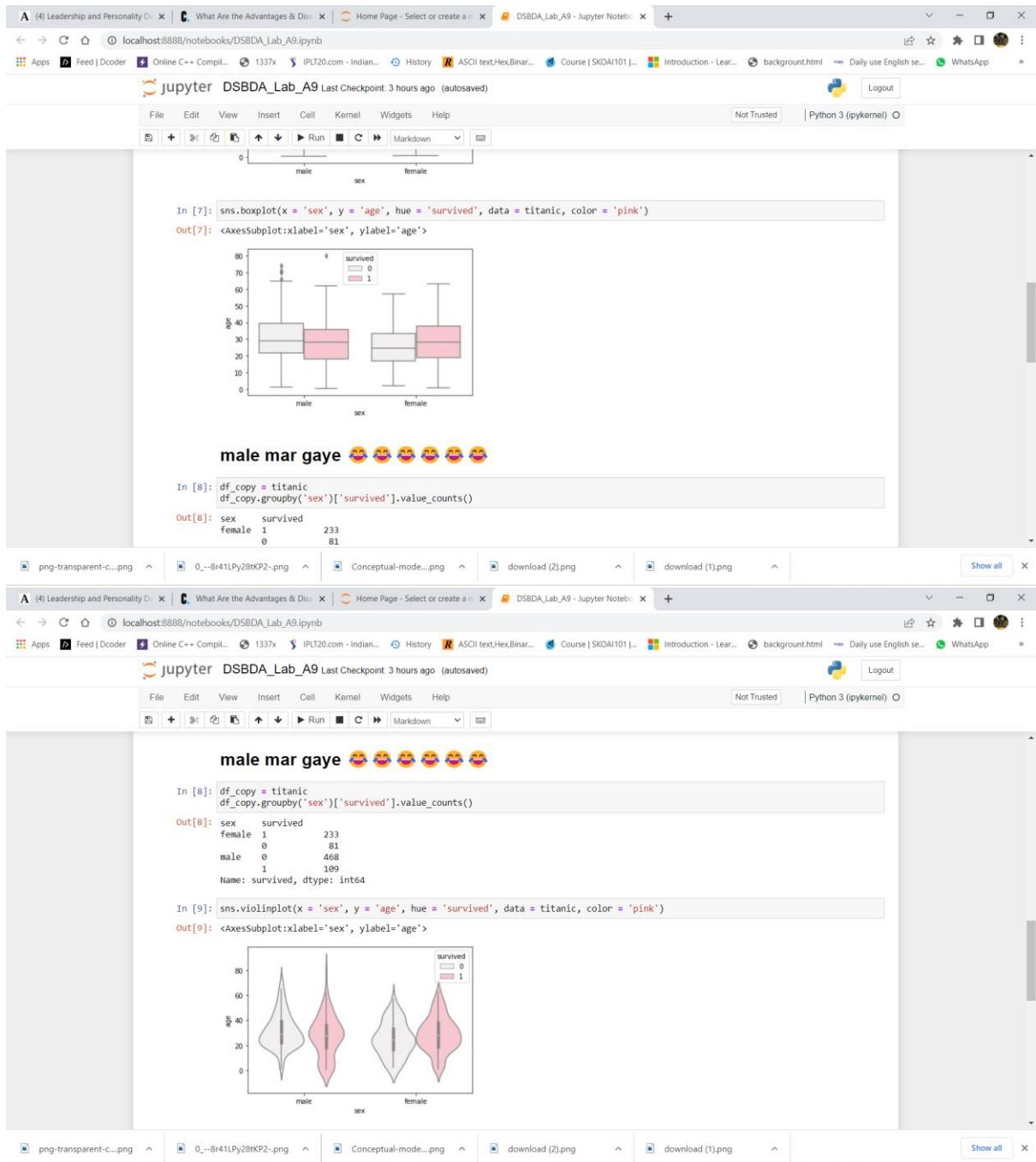
891 rows × 15 columns

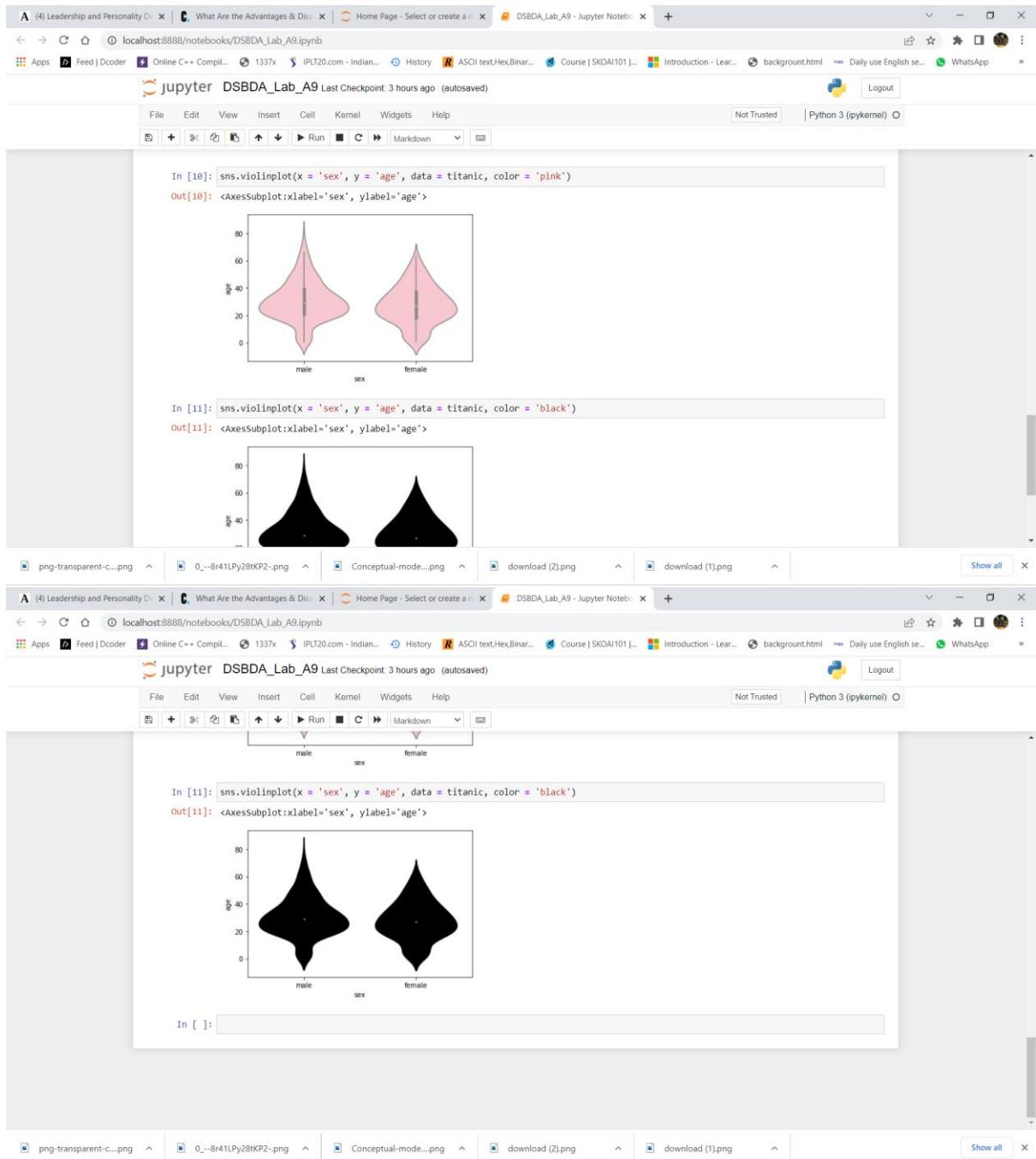
```
In [6]: sns.boxplot(x = 'sex', y = 'age', data = titanic, color = 'pink')
```

```
Out[6]: <AxesSubplot:xlabel='sex', ylabel='age'>
```

```
In [7]: sns.boxplot(x = 'sex', y = 'age', hue = 'survived', data = titanic, color = 'pink')
```

```
Out[7]: <AxesSubplot:xlabel='sex', ylabel='age'>
```





**Data Visualization III**

```
In [2]: import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

In [3]: data1 = pd.read_csv("Iris.csv")
data1.head()

Out[3]:   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 [4]: print(data1.columns)
Index(['Id', 'SepalLengthCm', 'SepalWidthCm', 'PetalLengthCm', 'PetalWidthCm',
       'Species'],
      dtype='object')

In [5]: #anotherway
column = list(data1)
print(column)

In [5]: #anotherway
column = list(data1)
print(column)

In [6]: data1.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 [7]: data1.dtypes

Out[7]: Id          int64
SepallengthCm  float64
SepalwidthCm   float64
PetallengthCm  float64
PetalwidthCm   float64
Species        object
dtype: object

In [8]: data1.hist()
```

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A10 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A10.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [8]: `data1.hist()`

Out[8]: array([[<AxesSubplot:title={'center':'Id'}>,<AxesSubplot:title={'center':'SepallengthCm'}>,<AxesSubplot:title={'center':'SepalwidthCm'}>,<AxesSubplot:title={'center':'PetallengthCm'}>,<AxesSubplot:title={'center':'PetalwidthCm'}>],<AxesSubplot:title={'center': 'PetalwidthCm'}>],<AxesSubplot:>]],  
dtype=object)

In [9]: `fig, axes = plt.subplots(2, 2, figsize=(16, 8))`

axes[0,0].set\_title("Distribution of First Column")  
axes[0,0].hist(data1["SepallengthCm"]);  
axes[0,1].set\_title("Distribution of Second Column")  
axes[0,1].hist(data1["SepalwidthCm"]);

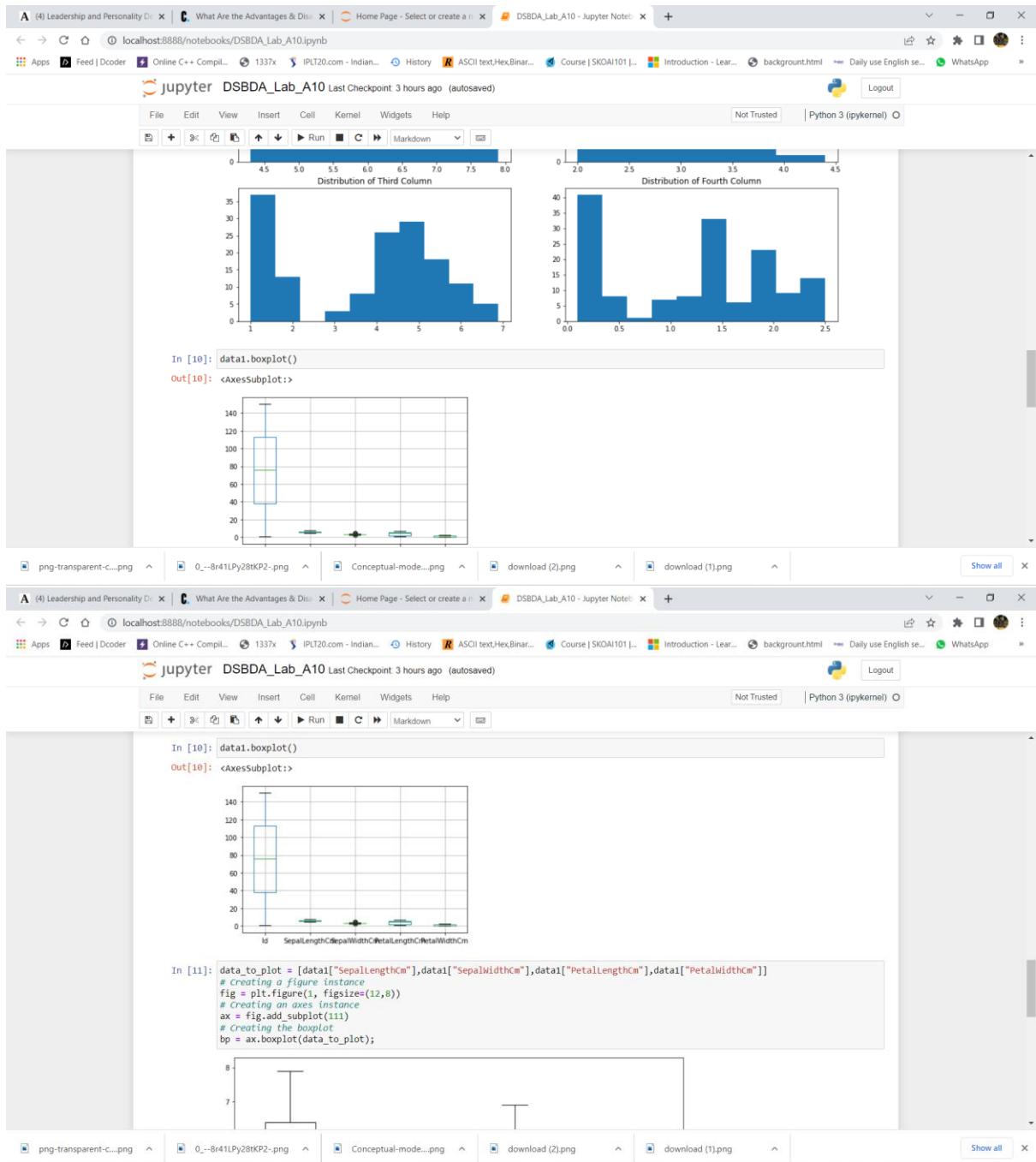
A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Lab\_A10 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Lab\_A10.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [9]: `fig, axes = plt.subplots(2, 2, figsize=(16, 8))`

axes[0,0].set\_title("Distribution of First Column")  
axes[0,0].hist(data1["SepallengthCm"]);  
axes[0,1].set\_title("Distribution of Second Column")  
axes[0,1].hist(data1["SepalwidthCm"]);  
axes[1,0].set\_title("Distribution of Third Column")  
axes[1,0].hist(data1["PetallengthCm"]);  
axes[1,1].set\_title("Distribution of Fourth Column")  
axes[1,1].hist(data1["PetalwidthCm"]);



In [11]:

```
data_to_plot = [data1["SepalLengthCm"],data1["SepalWidthCm"],data1["PetalLengthCm"],data1["PetalWidthCm"]]
# Creating a figure instance
fig = plt.figure(1, figsize=(12,8))
# Creating an axes instance
ax = fig.add_subplot(111)
# creating the boxplot
bp = ax.boxplot(data_to_plot);
```

In [12]:

```
sns.boxplot(data1['SepalWidthCm'])
```

C:\Users\rushil\AppData\Local\Programs\Python\Python310\lib\site-packages\seaborn\\_decorators.py:36: FutureWarning: Pass the following variable as a keyword arg: x. From version 0.12, the only valid positional argument will be `data`, and passing other arguments without an explicit keyword will result in an error or misinterpretation.  
warnings.warn(  
Out[12]:

`<AxesSubplot:xlabel='SepalWidthCm'>`

In [13]:

```
print(np.where(data1['SepalWidthCm']>4.0))
```

```
(array([15, 32, 33], dtype=int64),)
```

In [ ]:

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Assg\_B2 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Assg\_B2.ipynb

File Edit View Insert Cell Kernel Widgets Help

In [4]: `#import the dataset`  
`import numpy as np`  
`import pandas as pd`

In [5]: `data = pd.read_csv('weather.csv')`

In [6]: `data.info()`

#	Column	Non-Null Count	Dtype
0	MinTemp	366	non-null float64
1	MaxTemp	366	non-null float64
2	Rainfall	366	non-null float64
3	Evaporation	366	non-null float64
4	Sunshine	363	non-null float64
5	WindGustDir	363	non-null object
6	WindGustSpeed	364	non-null float64
7	WindDir9am	335	non-null object
8	WindDir3pm	365	non-null object
9	WindSpeed9am	359	non-null float64
10	WindSpeed3pm	366	non-null int64
11	Humidity9am	366	non-null int64
12	Humidity3pm	366	non-null int64
13	Pressure9am	366	non-null float64
14	Pressure3pm	366	non-null float64
15	Cloud9am	366	non-null int64
16	Cloud3pm	366	non-null int64
17	Temp9am	366	non-null float64
18	Temp3pm	366	non-null float64
19	RainToday	366	non-null object
20	RISK_MM	366	non-null float64
21	RainTomorrow	366	non-null object

dtypes: float64(12), int64(5), object(5)  
memory usage: 63.0+ KB

In [7]: `data.head(10)`

Out[7]:

	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDir9am	WindDir3pm	WindSpeed9am	...
0	8.0	24.3	0.0	3.4	6.3	NW	30.0	SW	NW	6.0	...
1	14.0	26.9	3.6	4.4	9.7	ENE	39.0	E	W	4.0	...
2	13.7	23.4	3.6	5.8	3.3	NW	85.0	N	NNE	6.0	...
3	13.3	15.5	39.8	7.2	9.1	NW	54.0	NNW	W	30.0	...
4	7.6	16.1	2.8	5.6	10.6	SSE	50.0	SSE	ESE	20.0	...
5	6.2	16.9	0.0	5.8	8.2	SE	44.0	SE	E	20.0	...

Windows Update  
Your device will restart to update outside of active hours (estimate: 6 min)  
Leave it on and plugged in. Open Settings to adjust your active hours and whether you want reminders.  
Restart now OK

Show all

A (4) Leadership and Personality D | C. What Are the Advantages & Disadvantages of Machine Learning? | Home Page - Select or create a new notebook | DSBDA\_Asgg\_B2 - Jupyter Notebook

localhost:8888/notebooks/DSBDA\_Asgg\_B2.ipynb

File Edit View Insert Cell Kernel Widgets Help

Not Trusted Python 3 (ipykernel)

jupyter DSBDA\_Asgg\_B2 Last Checkpoint: Yesterday at 12:31 AM (autosaved)

Logout

In [7]: `data.head(10)`

Out[7]:

	MinTemp	MaxTemp	Rainfall	Evaporation	Sunshine	WindGustDir	WindGustSpeed	WindDir9am	WindDir3pm	WindSpeed9am	Humidity3pm	Pressure9am	
0	8.0	24.3	0.0	3.4	6.3	NW	30.0	SW	NW	6.0	...	29	1019
1	14.0	26.9	3.6	4.4	9.7	ENE	39.0	E	W	4.0	...	36	1012
2	13.7	23.4	3.6	5.8	3.3	NW	85.0	N	NNE	6.0	...	69	1009
3	13.3	15.5	39.8	7.2	9.1	NW	54.0	VNW	W	30.0	...	56	1005
4	7.6	16.1	2.8	5.6	10.6	SSE	50.0	SSE	ESE	20.0	...	49	1018
5	6.2	16.9	0.0	5.8	8.2	SE	44.0	SE	E	20.0	...	57	1023
6	8.1	18.2	0.2	4.2	8.4	SE	43.0	SE	ESE	19.0	...	47	1024
7	8.3	17.0	0.0	5.6	4.6	E	41.0	SE	E	11.0	...	57	1026
8	8.8	19.5	0.0	4.0	4.1	S	48.0	E	ENE	19.0	...	48	1026
9	8.4	22.8	16.2	5.4	7.7	E	31.0	S	ESE	7.0	...	32	1024

10 rows × 22 columns

In [8]: `print("describe: ")  
print(data.describe())`

Windows Update

Your device will restart to update outside of active hours (estimate: 6 min)  
Leave it on and plugged in. Open Settings to adjust your active hours and whether you want reminders.

Restart now OK

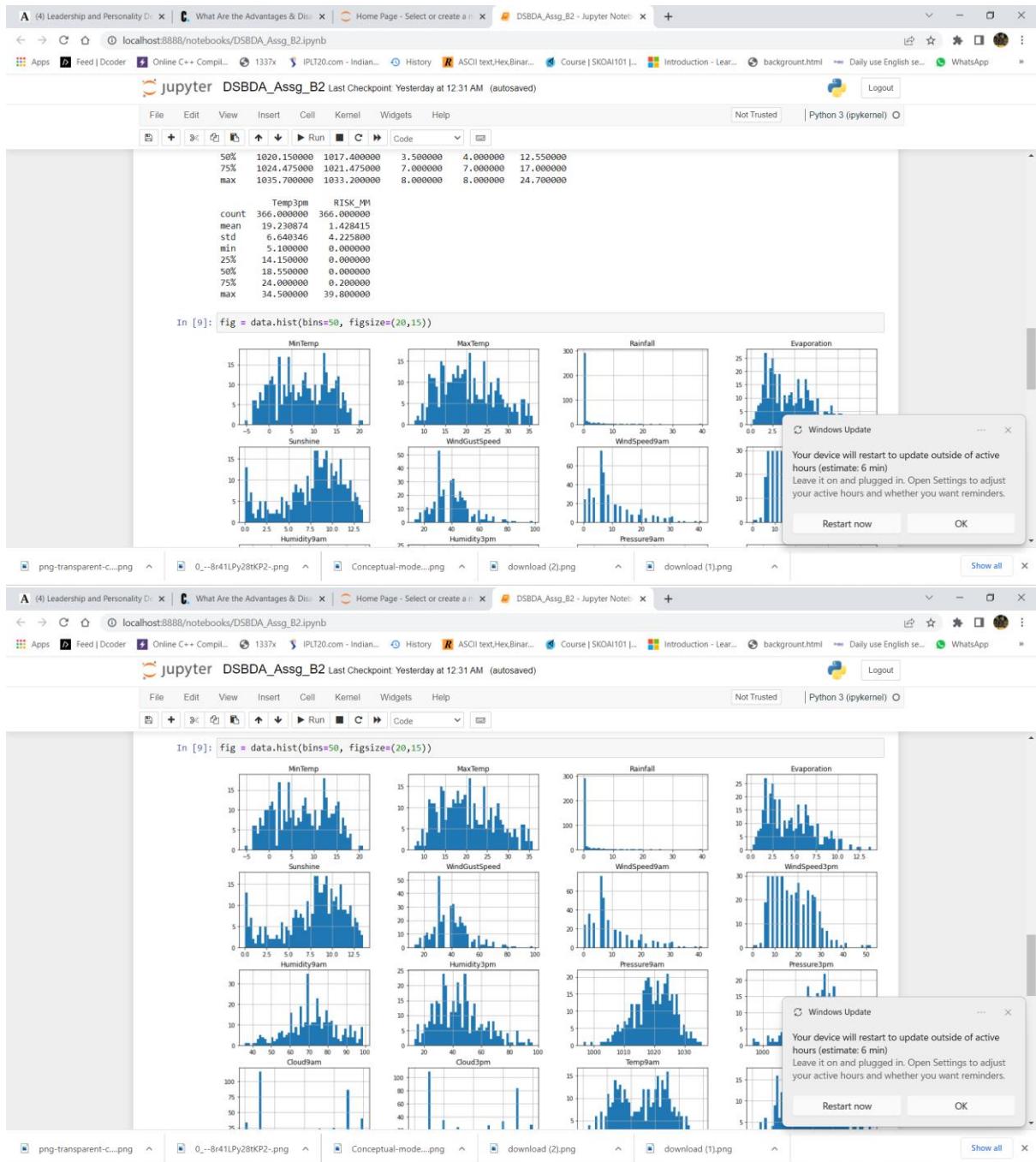
Show all

In [8]: `print("describe: ")  
print(data.describe())`

Windows Update

Your device will restart to update outside of active hours (estimate: 6 min)  
Leave it on and plugged in. Open Settings to adjust your active hours and whether you want reminders.

Restart now OK



Screenshot of a Jupyter Notebook session titled "DSBDA\_Assg\_B2 - Jupyter Notebook".

The notebook contains the following code:

```
In [10]: data.dropna(inplace=True)
data.drop_duplicates(inplace=True)
print("Success")
Success
```

Below the code, there are four histograms showing the distribution of variables: Cloud9am, Cloud3pm, Temp9am, and Temp3pm.

In [14]:

```
import matplotlib.pyplot as plt
import seaborn as sns
plt.figure(figsize=(20,10))
sns.heatmap(data.corr(), annot=True, cmap='cubehelix_r')
plt.show()
```

A heatmap visualization of the correlation matrix for the dataset. The columns and rows are labeled with variables: MinTemp, MaxTemp, Rainfall, Evaporation, Sunshine, WindGustSpeed, WindSpeed9am, WindSpeed3pm, Humidity9am, Humidity3pm, Pressure9am, Pressure3pm, Cloud9am, Cloud3pm, Temp9am, Temp3pm, and Risk\_MM. The color scale ranges from -1.0 (dark purple) to 1.0 (dark red).

A Windows Update dialog box is visible in the top right corner, indicating a scheduled restart.

In [1]:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
```

In [6]:

```
df = pd.read_csv('data_science.csv')
df.head()
```

```
C:\Users\Acer\anaconda3\lib\site-packages\IPython\core\interactiveshell.py:3165: DtypeWarning: Columns (9) have mixed types. Specify dtype option on import or set low_memory=False
.
    has_raised = await self.run_ast_nodes(code_ast.body, cell_name,
```

Out[6]:

	<b>id</b>	<b>conversation_id</b>	<b>created_at</b>	<b>date</b>	<b>time</b>	<b>timezone</b>	<b>user_id</b>	<b>username</b>
0	1406400408545804288	1406400396264943616		2021-06-20 05:26:01	2021-06-20 05:26:01	05:26:01	530	1113747629282930688 ballouxfrancois
1	1406390341176016897	1406390341176016897		2021-06-20 04:46:01	2021-06-20 04:46:01	04:46:01	530	788898706586275840 tdatascience
2	1406386311481774083	1406386311481774083		2021-06-20 04:30:00	2021-06-20 04:30:00	04:30:00	530	19402238 sciencenews
3	1406383545153638402	1406383545153638402		2021-06-20 04:19:01	2021-06-20 04:19:01	04:19:01	530	788898706586275840 tdatascience
4	1406358632648818689	1406358632648818689		2021-06-20 02:40:01	2021-06-20 02:40:01	02:40:01	530	788898706586275840 tdatascience

5 rows × 36 columns

In [7]:

```
df.shape
```

Out[7]:

```
(241386, 36)
```

In [8]:

```
df.isnull().sum()
```

Out[8]:

id	0
conversation_id	0
created_at	0
date	0
time	0
timezone	0
user_id	0
username	0
name	0

```
place          241032
tweet           0
language         0
mentions         0
urls             0
photos            0
replies_count    0
retweets_count   0
likes_count      0
hashtags          0
cashtags          0
link              0
retweet            0
quote_url        231065
video              0
thumbnail        131048
near              241386
geo               241386
source             241386
user_rt_id       241386
user_rt           241386
retweet_id        241386
reply_to          0
retweet_date      241386
translate          241386
trans_src          241386
trans_dest         241386
dtype: int64
```

In [9]:

```
#data cleaning
```

In [10]:

```
columns_to_drop = list(df.columns[27:])
df.drop(columns = columns_to_drop, axis = 1, inplace = True)
```

In [11]:

```
df.isnull().sum()
```

Out[11]:

```
id                  0
conversation_id     0
created_at          0
date                0
time                0
timezone            0
user_id              0
username             0
name                0
place          241032
tweet           0
language           0
mentions           0
urls              0
photos             0
replies_count      0
retweets_count     0
likes_count         0
hashtags           0
cashtags           0
link              0
retweet            0
quote_url        231065
video              0
thumbnail        131048
near              241386
geo               241386
dtype: int64
```

In [12]:

```
df.drop(columns = ['place', 'quote_url', 'thumbnail', 'near', 'geo'], axis = 1, inplace = True)
```

In [13]:

```
df.isnull().sum()
```

Out[13]:

```
id          0  
conversation_id 0  
created_at    0  
date          0  
time          0  
timezone      0  
user_id       0  
username      0  
name          0  
tweet          0  
language       0  
mentions       0  
urls          0  
photos         0  
replies_count 0  
retweets_count 0  
likes_count    0  
hashtags       0  
cashtags       0  
link           0  
retweet         0  
video          0  
dtype: int64
```

In [14]:

```
pd.set_option('display.max_columns', None)  
df.head()
```

Out[14]:

	<b>id</b>	<b>conversation_id</b>	<b>created_at</b>	<b>date</b>	<b>time</b>	<b>timezone</b>	<b>user_id</b>	<b>username</b>
0	1406400408545804288	1406400396264943616		2021-06-20 05:26:01 IST	2021-06-20 05:26:01	530	1113747629282930688	ballouxfrancois
1	1406390341176016897	1406390341176016897		2021-06-20 04:46:01 IST	2021-06-20 04:46:01	530	788898706586275840	tdatascience
2	1406386311481774083	1406386311481774083		2021-06-20 04:30:00 IST	2021-06-20 04:30:00	530	19402238	sciencenews
3	1406383545153638402	1406383545153638402		2021-06-20 04:19:01 IST	2021-06-20 04:19:01	530	788898706586275840	tdatascience
4	1406358632648818689	1406358632648818689		2021-06-20 02:40:01 IST	2021-06-20 02:40:01	530	788898706586275840	tdatascience

In [15]:

```
#Dropping cashtags column
```

```
df.drop(columns = ['cashtags'], axis = 1, inplace= True)
```

In [16]:

```
df.head()
```

Out[16]:

	<b>id</b>	<b>conversation_id</b>	<b>created_at</b>	<b>date</b>	<b>time</b>	<b>timezone</b>	<b>user_id</b>	<b>username</b>
0	1406400408545804288	1406400396264943616	2021-06-20 05:26:01 IST	2021-06-20	05:26:01		530	1113747629282930688 ballouxfrancois
1	1406390341176016897	1406390341176016897	2021-06-20 04:46:01 IST	2021-06-20	04:46:01		530	788898706586275840 tdatasience
2	1406386311481774083	1406386311481774083	2021-06-20 04:30:00 IST	2021-06-20	04:30:00		530	19402238 sciencenews
3	1406383545153638402	1406383545153638402	2021-06-20 04:19:01 IST	2021-06-20	04:19:01		530	788898706586275840 tdatasience
4	1406358632648818689	1406358632648818689	2021-06-20 02:40:01 IST	2021-06-20	02:40:01		530	788898706586275840 tdatasience

In [17]:

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 241386 entries, 0 to 241385
Data columns (total 21 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   id               241386 non-null   int64  
 1   conversation_id  241386 non-null   int64  
 2   created_at       241386 non-null   object  
 3   date             241386 non-null   object  
 4   time             241386 non-null   object  
 5   timezone         241386 non-null   int64  
 6   user_id          241386 non-null   int64  
 7   username         241386 non-null   object  
 8   name              241386 non-null   object  
 9   tweet             241386 non-null   object  
 10  language          241386 non-null   object  
 11  mentions          241386 non-null   object  
 12  urls              241386 non-null   object  
 13  photos             241386 non-null   object  
 14  replies_count    241386 non-null   int64  
 15  retweets_count   241386 non-null   int64  
 16  likes_count       241386 non-null   int64  
 17  hashtags          241386 non-null   object  
 18  link              241386 non-null   object  
 19  retweet            241386 non-null   bool    
 20  video              241386 non-null   int64  
dtypes: bool(1), int64(8), object(12)
memory usage: 37.1+ MB
```

In [18]:

```
#Removing IST from created_at column
df['created_at'] = df['created_at'].apply(lambda tweet: tweet.replace('IST', '') .strip())
```

```
)
```

```
In [19]:
```

```
df['created_at'] = pd.to_datetime(df['created_at'])
df['date'] = pd.to_datetime(df['date'])
df['time'] = pd.to_datetime(df['time'])
```

```
In [20]:
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 241386 entries, 0 to 241385
Data columns (total 21 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   id               241386 non-null   int64  
 1   conversation_id  241386 non-null   int64  
 2   created_at       241386 non-null   datetime64[ns]
 3   date             241386 non-null   datetime64[ns]
 4   time             241386 non-null   datetime64[ns]
 5   timezone         241386 non-null   int64  
 6   user_id          241386 non-null   int64  
 7   username         241386 non-null   object  
 8   name              241386 non-null   object  
 9   tweet             241386 non-null   object  
 10  language          241386 non-null   object  
 11  mentions          241386 non-null   object  
 12  urls              241386 non-null   object  
 13  photos             241386 non-null   object  
 14  replies_count    241386 non-null   int64  
 15  retweets_count   241386 non-null   int64  
 16  likes_count       241386 non-null   int64  
 17  hashtags          241386 non-null   object  
 18  link              241386 non-null   object  
 19  retweet            241386 non-null   bool    
 20  video              241386 non-null   int64  
dtypes: bool(1), datetime64[ns](3), int64(8), object(9)
memory usage: 37.1+ MB
```

```
In [21]:
```

```
df = df.astype({'username':'category', 'name':'category', 'tweet':'category', 'language':'category',
               'link':'category', 'urls':'category'})
```

```
In [22]:
```

```
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 241386 entries, 0 to 241385
Data columns (total 21 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   id               241386 non-null   int64  
 1   conversation_id  241386 non-null   int64  
 2   created_at       241386 non-null   datetime64[ns]
 3   date             241386 non-null   datetime64[ns]
 4   time             241386 non-null   datetime64[ns]
 5   timezone         241386 non-null   int64  
 6   user_id          241386 non-null   int64  
 7   username         241386 non-null   category 
 8   name              241386 non-null   category 
 9   tweet             241386 non-null   category 
 10  language          241386 non-null   category 
 11  mentions          241386 non-null   object  
 12  urls              241386 non-null   category 
 13  photos             241386 non-null   object  
 14  replies_count    241386 non-null   int64  
 15  retweets_count   241386 non-null   int64  
 16  likes_count       241386 non-null   int64  
 17  hashtags          241386 non-null   object  
 18  link              241386 non-null   object  
 19  retweet            241386 non-null   bool    
 20  video              241386 non-null   int64  
dtypes: bool(1), category(10), int64(8)
```

```
+--|--|--|--|--+  
17 hashtags          241386 non-null object  
18 link              241386 non-null category  
19 retweet            241386 non-null bool  
20 video              241386 non-null int64  
dtypes: bool(1), category(6), datetime64[ns](3), int64(8), object(3)  
memory usage: 56.1+ MB
```

In [23]:

```
#splitting dataframe into 4 parts
```

In [24]:

```
parts = np.array_split(df, 4)
```

In [25]:

```
parts[0].head()
```

Out[25]:

	<b>id</b>	<b>conversation_id</b>	<b>created_at</b>	<b>date</b>	<b>time</b>	<b>timezone</b>	<b>user_id</b>	<b>username</b>
0	1406400408545804288	1406400396264943616	2021-06-20 05:26:01	2021-06-20	2022-05-04 05:26:01	530	111374762928930688	ballouxfrancois
1	1406390341176016897	1406390341176016897	2021-06-20 04:46:01	2021-06-20	2022-05-04 04:46:01	530	788898706586275840	tdatascience
2	1406386311481774083	1406386311481774083	2021-06-20 04:30:00	2021-06-20	2022-05-04 04:30:00	530	19402238	sciencenews
3	1406383545153638402	1406383545153638402	2021-06-20 04:19:01	2021-06-20	2022-05-04 04:19:01	530	788898706586275840	tdatascience
4	1406358632648818689	1406358632648818689	2021-06-20 02:40:01	2021-06-20	2022-05-04 02:40:01	530	788898706586275840	tdatascience

In [26]:

```
#Converting to lowercase  
def to_lowercase(text):  
    return text.lower()
```

In [27]:

```
for i in range(4):  
    parts[i]['tweet'] = parts[i]['tweet'].apply(lambda tweet: to_lowercase(tweet))
```

In [28]:

```
for i in range(4):  
    print(parts[i]['tweet'].head())
```

```
0    what can be done? - never blindly trust an ab...  
1    "we need a paradigm shift from model-centric t...  
2    using high-resolution satellite data and compu...  
3    .@stephenson_data shares four steps that will ...  
4    "curricula is inherently brittle in a world wh...  
Name: tweet, dtype: object  
60347    the world is digital – learn about the power o...  
60348    build your data brilliance in 2020 – explore y...  
60349    the end is near, but there is still time to ea...  
60350    how to make the transition from analyst to #da...  
60351    that time when the australian worked with clim...  
Name: tweet, dtype: object  
120601    checking out "22 tips for better data science"
```

```
120694     checking out 22 tips for better data science ...
120695 10 questions every data decision-maker should ...
120696 tired of data silos? join our webinar on may 2...
120697 data science is helping people in india with t...
120698 data science master's students visit accenture...
Name: tweet, dtype: object
181040 join the #datascience movement. immerse yourse...
181041 don't forget to vote for @rapidminer in @kdnug...
181042 want to learn more about our #datascience and ...
181043 @andrewchen @mikeschiemer @dberkowitz @rwang0 ...
181044 data are data and our field is exceptionally g...
Name: tweet, dtype: object
```

In [29]:

```
#removing special character ,html tags and numbers
```

In [30]:

```
from bs4 import BeautifulSoup
import re
```

In [31]:

```
def clean_text(text):
    soup = BeautifulSoup(text, 'html.parser')
    cleaned_text = soup.get_text()
    cleaned_text = re.sub('[^A-Za-z0-9]+', ' ', cleaned_text)
    return cleaned_text
```

In [32]:

```
for i in range(4):
    parts[i]['tweet'] = parts[i]['tweet'].apply(lambda tweet: clean_text(tweet))
```

In [33]:

```
#tokenization
```

In [41]:

```
from nltk.tokenize import word_tokenize
```

In [42]:

```
def tokenization(text):
    tokens = word_tokenize(text)
    return tokens
```

In [ ]:

```
for i in range(4):
    parts[i]['tweet'] = parts[i]['tweet'].apply(lambda tweet: tokenization(tweet))
```

In [ ]:

```
#removing stopwords
```

In [44]:

```
from nltk.corpus import stopwords
```

In [ ]:

```
for i in range(4):
    stop_words = set(stopwords.words('english'))
    parts[i]['tweet'] = [[word for word in tweet if word not in stop_words] for tweet in parts[i]['tweet']]
```

In [46]:

```
for i in range(4):
    print(parts[i]['tweet'].head())

0    what can be done never blindly trust an abstra...
1    we need a paradigm shift from model centric t...
2    using high resolution satellite data and compu...
3    stephenson data shares four steps that will h...
4    curricula is inherently brittle in a world wh...
Name: tweet, dtype: object
60347    the world is digital learn about the power of ...
60348    build your data brilliance in 2020 explore you...
60349    the end is near but there is still time to ear...
60350    how to make the transition from analyst to dat...
60351    that time when the australian worked with clim...
Name: tweet, dtype: object
120694    checking out 22 tips for better data science d...
120695    10 questions every data decision maker should ...
120696    tired of data silos join our webinar on may 24...
120697    data science is helping people in india with t...
120698    data science master s students visit accenture...
Name: tweet, dtype: object
181040    join the datascience movement immerse yourself...
181041    don t forget to vote for rapidminer in kdnugge...
181042    want to learn more about our datascience and w...
181043    andrewchen mikeschiemer dberkowitz rwang0 our...
181044    data are data and our field is exceptionally g...
Name: tweet, dtype: object
```

In [47]:

```
#Stemming and Lemmatization
```

In [48]:

```
from nltk.stem import PorterStemmer
porter = PorterStemmer()
```

In [49]:

```
for i in range(4):
    parts[i]['tweet'] = [[porter.stem(word) for word in tweet] for tweet in parts[i]['t
weet']]
```

In [50]:

```
#Sentiment Analysis
```

In [ ]:

```
from textblob import TextBlob
```

In [ ]:

```
def get_sentiment(tweet):
    text = ' '.join(word for word in tweet)
    analysis = TextBlob(text)
    if analysis.sentiment.polarity > 0:
        return 'positive'
    elif analysis.sentiment.polarity == 0:
        return 'neutral'
    else:
        return 'negative'
```

In [ ]:

```
for i in range(4):
    parts[i]['sentiment'] = parts[i]['tweet'].apply(lambda tweet : get_sentiment(tweet))
```

In [ ]:

```
for i in range(4):
    print(parts[i][['tweet', 'sentiment']].head())
```

In [52]:

```
#joining all in one dataframe
```

In [53]:

```
df_sentiment = parts[0].append([parts[1], parts[2], parts[3]])
```

In [54]:

```
df_sentiment.head()
```

Out[54]:

	<b>id</b>	<b>conversation_id</b>	<b>created_at</b>	<b>date</b>	<b>time</b>	<b>timezone</b>	<b>user_id</b>	<b>username</b>
0	1406400408545804288	1406400396264943616	2021-06-20 05:26:01	2021-06-20	2022-05-04 05:26:01		530	1113747629282930688 ballouxfrancois
1	1406390341176016897	1406390341176016897	2021-06-20 04:46:01	2021-06-20	2022-05-04 04:46:01		530	788898706586275840 tdatasience
2	1406386311481774083	1406386311481774083	2021-06-20 04:30:00	2021-06-20	2022-05-04 04:30:00		530	19402238 sciencenews
3	1406383545153638402	1406383545153638402	2021-06-20 04:19:01	2021-06-20	2022-05-04 04:19:01		530	788898706586275840 tdatasience
4	1406358632648818689	1406358632648818689	2021-06-20 02:40:01	2021-06-20	2022-05-04 02:40:01		530	788898706586275840 tdatasience

In [ ]: