

## Assignment 01

Name : Kiran K. Patil

ID:211070904

Sem : 06

Course : Machine Learning

### **Kaggle :**

Kaggle is a platform where data scientists and machine learning developers can compete against each other to solve complex data science problems. Companies and organizations also use Kaggle to host competitions in order to crowdsource solutions to problems they are facing.

The platform provides a variety of tools for participants, including a cloud-based workbench for developing and running code, as well as a large dataset repository. Kaggle also hosts a community forum where participants can ask questions, share ideas and collaborate on projects. Additionally, Kaggle offers a number of educational resources, including tutorials, courses, and a blog to help users learn data science and machine learning. It is also one of the most popular platform for data scientist to showcase their skills and find job opportunities



### **Pandas :**

Pandas is a open-source library for the Python programming language for data manipulation and analysis. It provides data structures such as dataframes and series that allow for easy manipulation and analysis of large datasets. Pandas is built on top of the NumPy library and is widely used in the data science community for its powerful data manipulation capabilities. Some of the main features of Pandas include.

Pandas is a powerful library with a wide range of functions for data manipulation and analysis. Here is a list of some of the most commonly used functions in Pandas, along with a brief description of what they do:

#### Pandas Functions :

- `read_csv()`: reads a CSV file and returns a dataframe
- `to_csv()`: writes a dataframe to a CSV file
- `head()`: returns the first n rows of a dataframe (default is 5)
- `tail()`: returns the last n rows of a dataframe (default is 5)
- `shape`: returns the dimensions of a dataframe (rows, columns)
- `info()`: returns information about a dataframe, including the data types of each column and the number of non-null values
- `describe()`: returns basic statistics for each numeric column in a dataframe
- `columns`: returns the column labels of a dataframe
- `index`: returns the index (row labels) of a dataframe
- `value_counts()`: returns the frequency counts for each unique value in a column
- `sort_values()`: sorts a dataframe by one or more columns
- `groupby()`: groups a dataframe by one or more columns and applies a function to each group
- `merge()`: merges two dataframes on one or more columns
- `pivot_table()`: creates a pivot table from a dataframe
- `melt()`: "melts" a dataframe and returns a new, reshaped dataframe
- `stack()`: "stacks" the columns of a dataframe and returns a new, reshaped dataframe

- `unstack()`: "unstacks" the rows of a dataframe and returns a new, reshaped dataframe
- `crosstab()`: creates a cross-tabulation (frequency table) of two or more factors

## **Seaborn :**

Seaborn is a Python data visualization library based on Matplotlib. It is built on top of Matplotlib and allows for easy creation of beautiful, informative, and highly-customizable statistical graphics. Some of the main features of Seaborn include:

Here is a list of some commonly used functions in Seaborn:

- `sns.lineplot()`: creates a line plot
- `sns.barplot()`: creates a bar plot
- `sns.scatterplot()`: creates a scatter plot
- `sns.histplot()`: creates a histogram
- `sns.boxplot()`: creates a box plot
- `sns.violinplot()`: creates a violin plot
- `sns.catplot()`: creates a categorical plot
- `sns.pairplot()`: creates a pair plot
- `sns.jointplot()`: creates a joint plot
- `sns.heatmap()`: creates a heat map
- `sns.regplot()`: creates a regression plot
- `sns.kdeplot()`: creates a kernel density estimate plot
- `sns.lmplot()`: creates a scatter plot with linear regression line
- `sns.countplot()`: creates a bar plot of counts
- `sns.despine()`: removes the top and right spines from the plot
- `sns.set_style()`: sets the background theme of the plot
- `sns.set_context()`: sets the context of the plot (paper, notebook, talk, poster)

## **Matplotlib :**

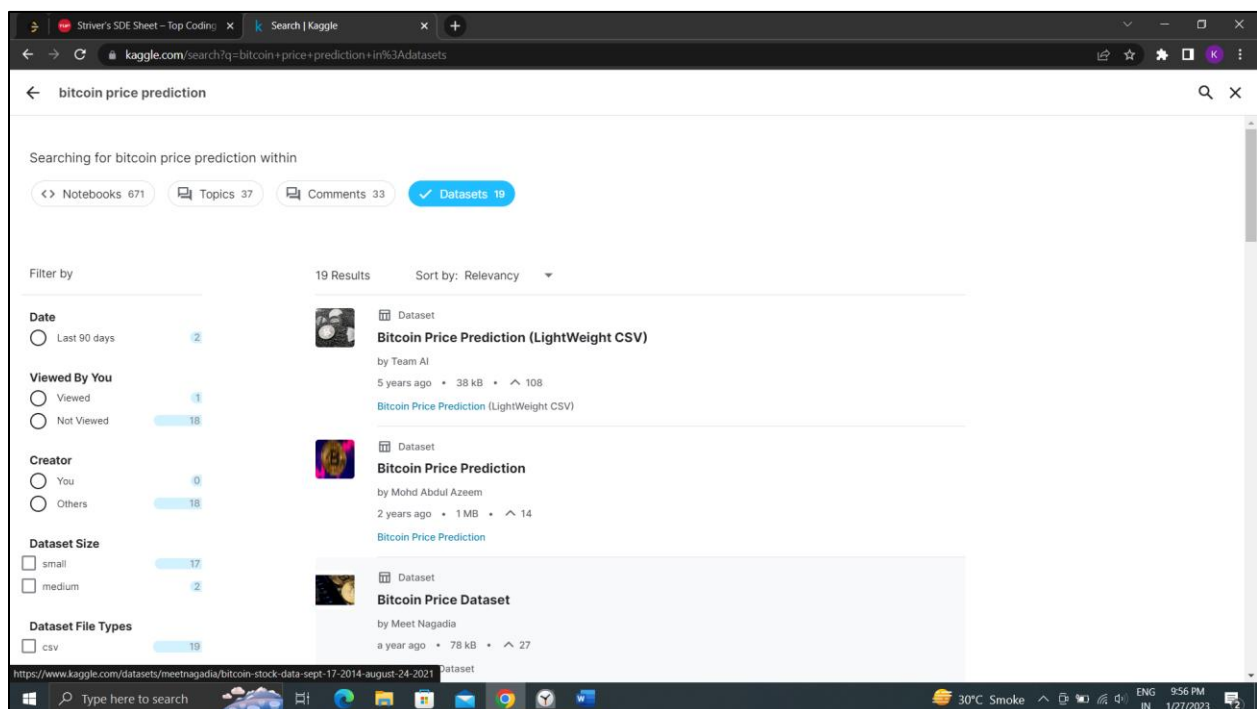
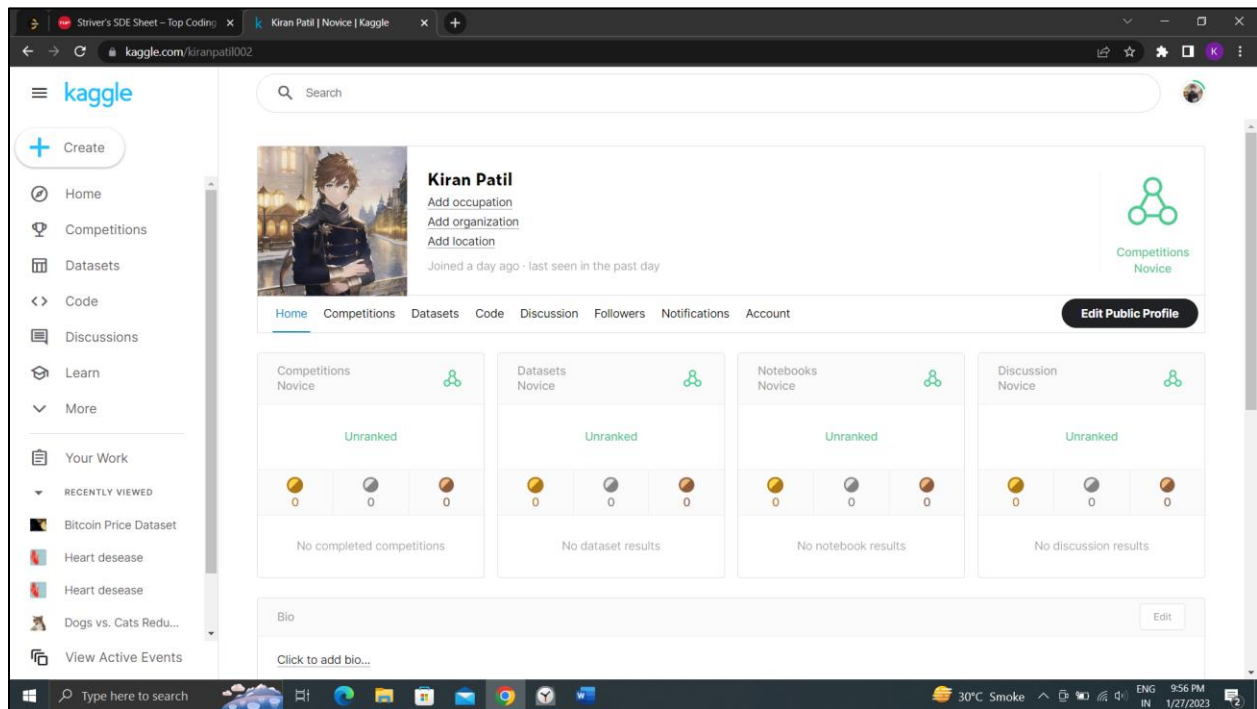
**Matplotlib** is a plotting library for the Python programming language and its numerical mathematics extension NumPy

Matplotlib is a powerful data visualization library that provides a wide variety of customizable plots, but it can be complex to use for creating more advanced visualizations. Seaborn is built on top of Matplotlib and provides a simpler, more convenient interface for creating many common types of plots, especially those used in statistical data visualization.

Here is a list of some commonly used functions in Matplotlib:

- `plt.plot()`: creates a line plot
- `plt.scatter()`: creates a scatter plot
- `plt.bar()`: creates a bar plot
- `plt.hist()`: creates a histogram
- `plt.pie()`: creates a pie chart
- `plt.boxplot()`: creates a box plot
- `plt.stem()`: creates a stem plot
- `plt.quiver()`: creates a vector field plot
- `plt.semilogx()`: creates a semilogarithmic plot of x-axis
- `plt.semilogy()`: creates a semilogarithmic plot of y-axis
- `plt.loglog()`: creates a log-log plot
- `plt.polar()`: creates a polar plot
- `plt.contour()`: creates a contour plot
- `plt.imshow()`: displays an image
- `plt.colorbar()`: adds a colorbar to a plot
- `plt.xlabel()`: adds a label to the x-axis
- `plt.ylabel()`: adds a label to the y-axis
- `plt.title()`: adds a title to the plot
- `plt.legend()`: adds a legend to the plot
- `plt.show()`: displays the plot

Setting up an account and download the dataset from Kaggle :



Striver's SDE Sheet - Top Coding | Bitcoin Price Dataset | Kaggle

kaggle.com/datasets/meetnagadia/bitcoin-stock-data-sept-17-2014-august-24-2021

### Bitcoin Price Dataset

Data Card Code (5) Discussion (0)

BTC-USD.csv (221.19 kB)

Detail Compact Column

7 of 7 columns

About this file

This is the dataset about Bitcoin Price

Date	# Open	# High	# Low	# Close	# Adj
2014-09-17	465.864014	468.174011	452.421997	457.334015	457
2014-09-18	456.859985	456.859985	413.104004	424.440002	424
2014-09-19	424.102997	427.834991	384.532013	394.795990	394
2014-09-20	394.673004	423.295990	389.882996	408.903992	408
2014-09-21	408.084991	412.425995	393.181000	398.821014	398
2014-09-22	399.100006	406.915985	397.130005	402.152008	402

Version 4 (221.19 kB)

BTC-USD.csv

View Active Events

https://www.kaggle.com/datasets/meetnagadia/bitcoin-stock-data-sept-17-2014-august-24-2021/download?datasetVersionNumber=4

Striver's SDE Sheet - Top Coding | Bitcoin Price Dataset | Kaggle | VITI-ML-Lab-01.ipynb - Colaboratory | ML-Lab 1.ipynb - Colaboratory

colab.research.google.com/drive/1y2kqf8QvofdR-034K0AFVp0m0q35x?authuser=1

### ML-Lab 1.ipynb

File Edit View Insert Runtime Tools Help Last saved at 11:01 AM

+ Code + Text

```
!pip install -q kaggle
```

```
[ ] from google.colab import files
    files.upload()
```

Choose Files No file chosen Upload widget is only available when the cell has been executed in the current browser session. Please rerun this cell to enable.

```
['kaggle.json': b'{"username": "kiranpatil002", "key": "a2a5805c073a7803eaa00712e5245b2b"}']
```

```
[ ] mkdir ~/.kaggle
```

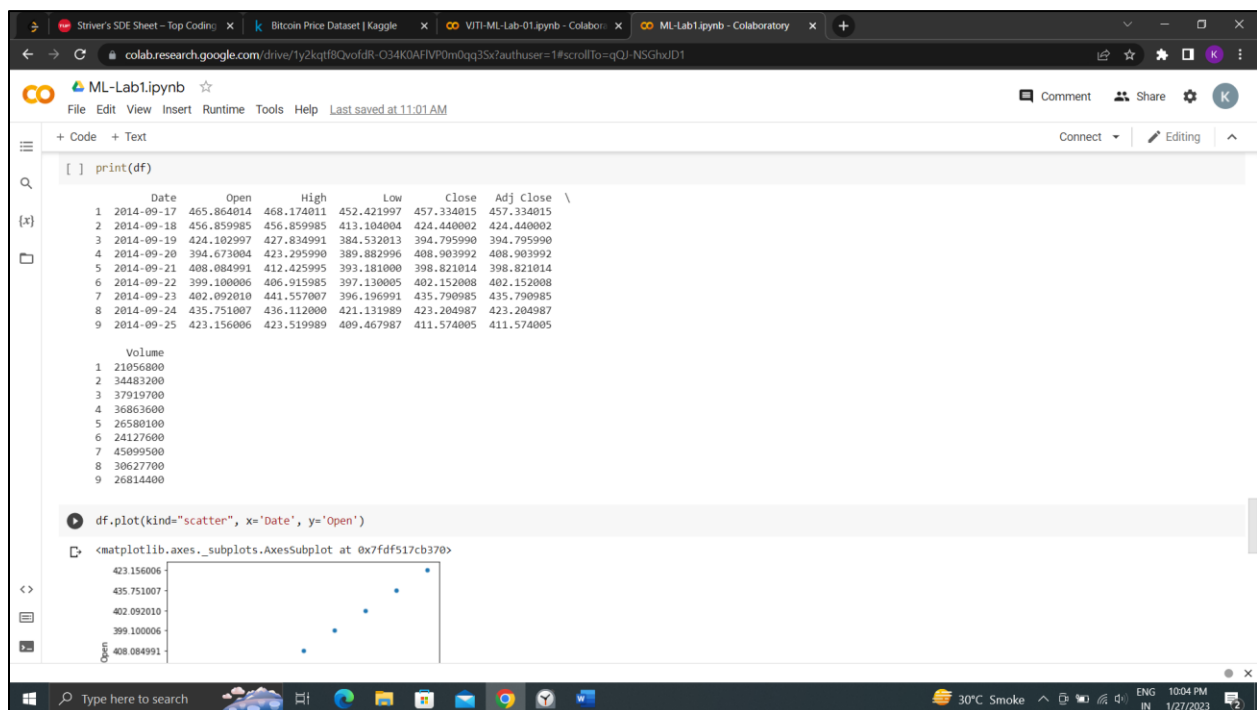
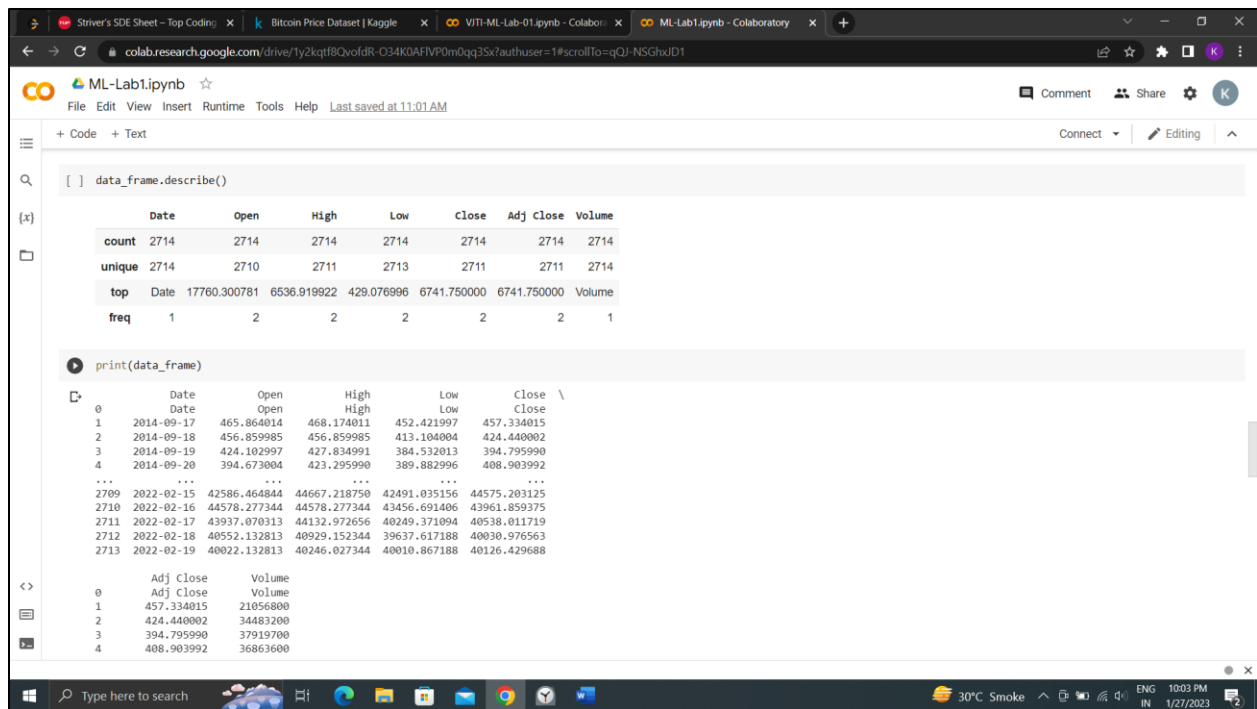
```
[ ] cp kaggle.json ~/.kaggle/
```

```
[ ] chmod 600 ~/.kaggle/kaggle.json
```

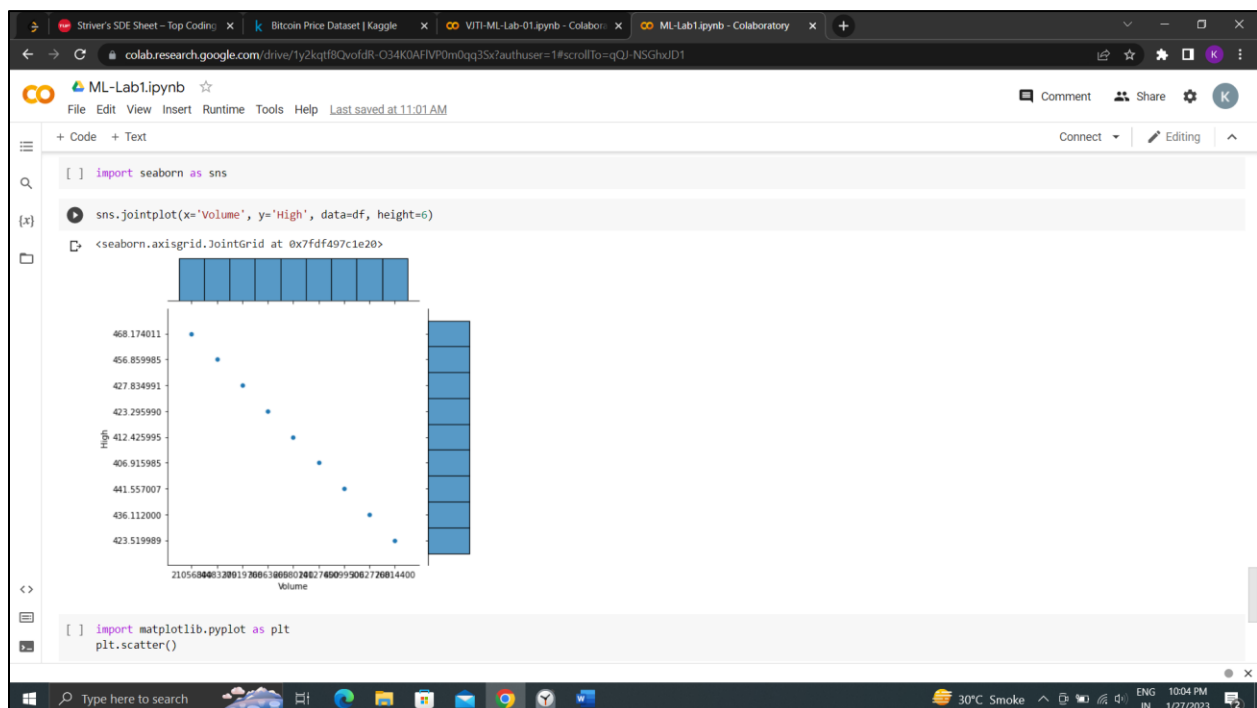
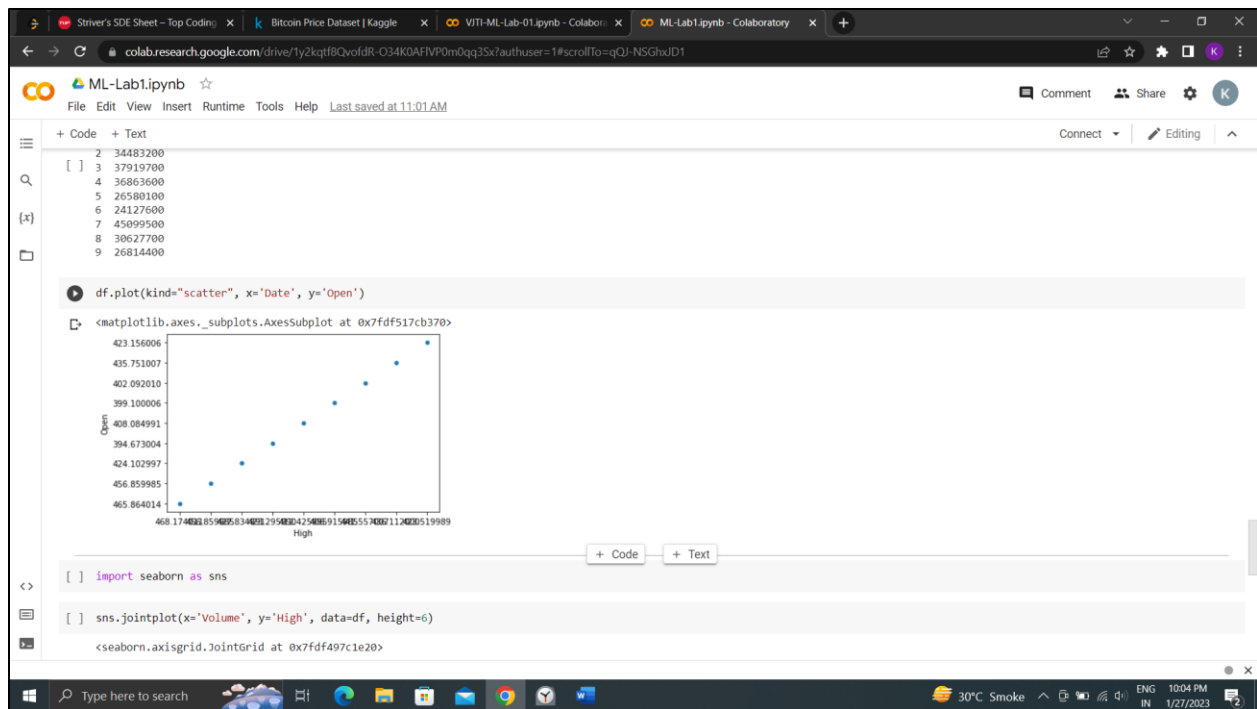
```
[ ] kaggle datasets list
```

ref	title	size	lastUpdated	downloadCount	voteCount	usabilityRating
ahsan81/hotel-reservations-classification-dataset	Hotel Reservations Dataset	480KB	2023-01-04 12:50:31	5041	175	1.0
senapatirajesh/netflix-tv-shows-and-movies	Latest Netflix TV shows and movies	1MB	2023-01-14 17:03:12	1490	46	0.9411765
johnny1994/divorce-rates-data-should-you-get-married	Divorce Rates Data: Should You Get Married?	22KB	2023-01-15 12:49:06	841	31	0.88235295
rakeshharv/spotify-top-10000-streamed-songs	Spotify Top 10000 Streamed Songs	280KB	2023-01-02 08:17:15	2196	66	1.0
diegorigephit/fifa-world-cup-2022-complete-dataset	Fifa World Cup 2022: Complete Dataset	7KB	2022-12-18 22:51:11	7425	227	1.0
themrityunjaypathak/indb-top-100-movies	IMDb Top 100 Movies	4KB	2023-01-11 17:15:09	1250	40	1.0
thedevasator/tesla-accident-fatalities-analysis-and-statistic	Tesla Deaths	18KB	2023-01-02 23:34:23	1130	30	0.9411765
rishikeshtkonapure/home-loan-approval	Home Loan Approval	13KB	2023-01-12 06:28:57	1406	36	1.0
ayushnith/starcraft-players-dataset	Starcraft Players Dataset - Gamers Analytics	205KB	2023-01-18 23:14:29	227	29	1.0
ruddygunawan/per-capita-income-by-county-2021-vs-education	Per Capita Income by County (2021) vs. Education	89KB	2022-12-28 14:37:42	1455	31	1.0

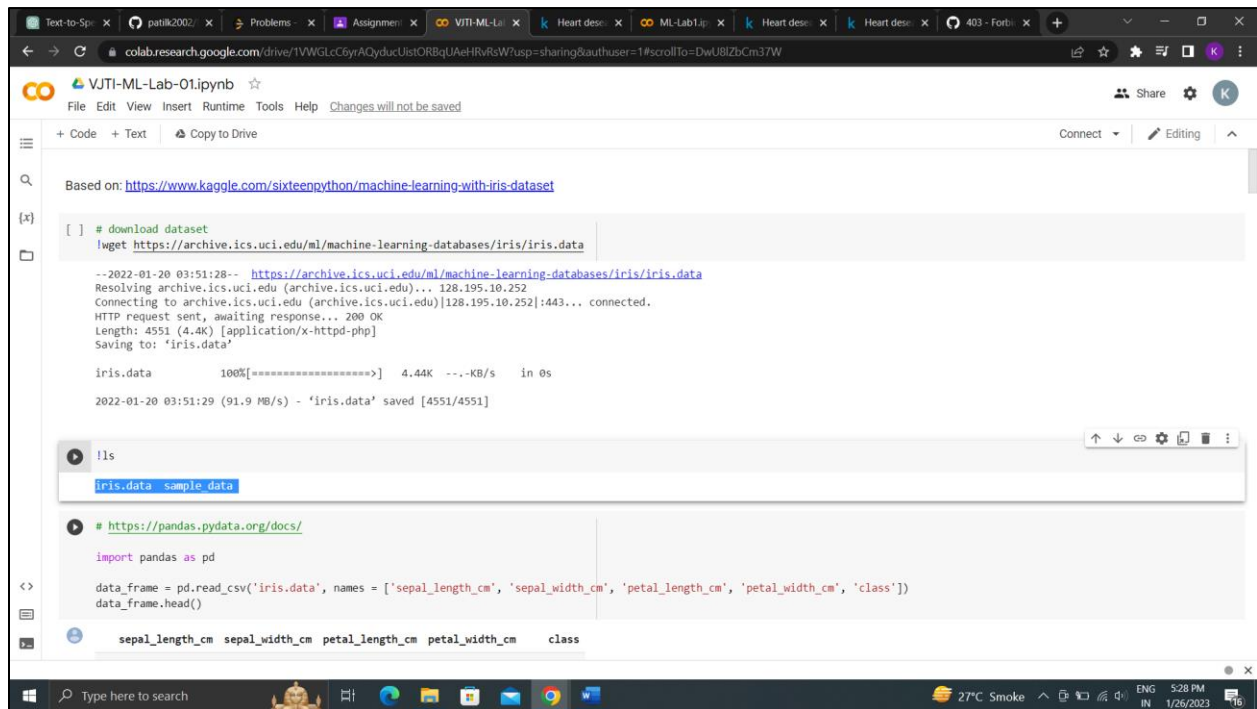








## Functions performed using Pandas :



The screenshot shows a Jupyter Notebook titled "VJTI-ML-Lab-01.ipynb" with the following code and output:

```
[ ] # download dataset
!wget https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data

--2022-01-20 03:51:28-- https://archive.ics.uci.edu/ml/machine-learning-databases/iris/iris.data
Resolving archive.ics.uci.edu (archive.ics.uci.edu)... 128.195.10.252
Connecting to archive.ics.uci.edu (archive.ics.uci.edu)|128.195.10.252|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: 4551 (4.4K) [application/x-httpd-php]
Saving to: 'iris.data'

iris.data      100%[=====] 4.44K  --.-KB/s  in 0s

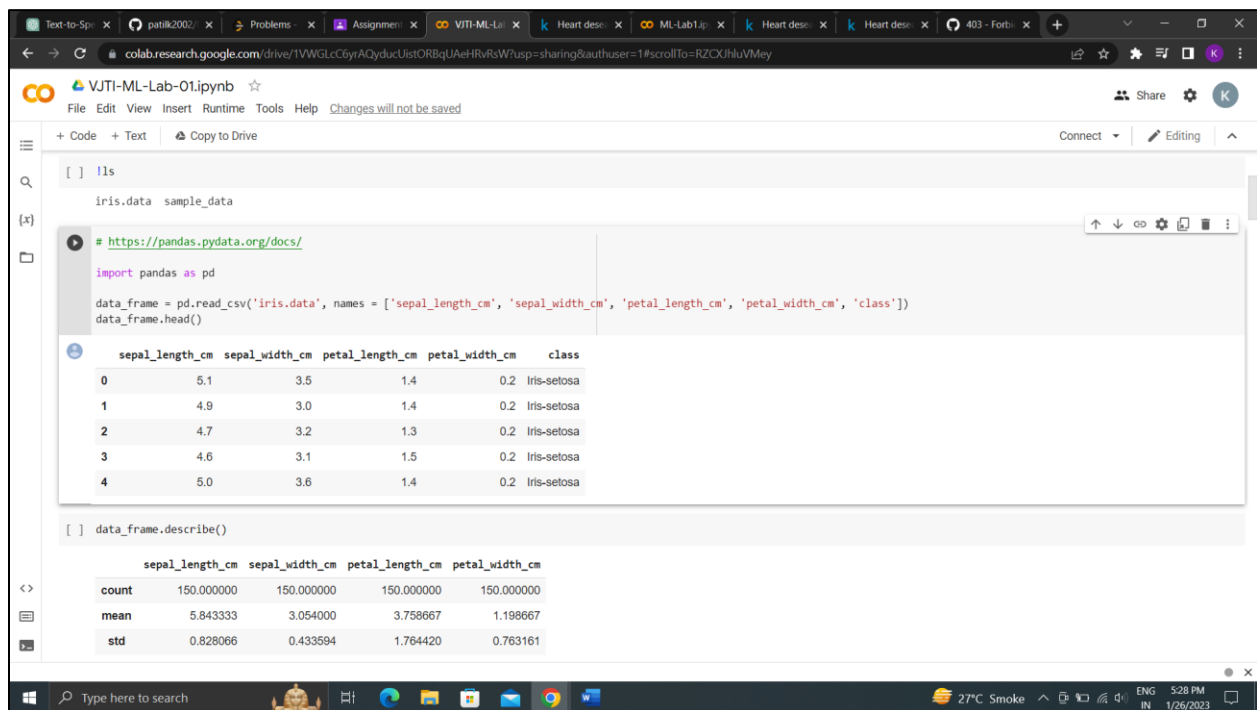
2022-01-20 03:51:29 (91.9 MB/s) - 'iris.data' saved [4551/4551]
```

```
[ ] !ls
iris.data sample_data
```

```
[ ] # https://pandas.pydata.org/docs/
import pandas as pd

data_frame = pd.read_csv('iris.data', names = ['sepal_length_cm', 'sepal_width_cm', 'petal_length_cm', 'petal_width_cm', 'class'])
data_frame.head()
```

	sepal_length_cm	sepal_width_cm	petal_length_cm	petal_width_cm	class
--	-----------------	----------------	-----------------	----------------	-------



The screenshot shows the continuation of the Jupyter Notebook with the following code and output:

```
[ ] !ls
iris.data sample_data
```

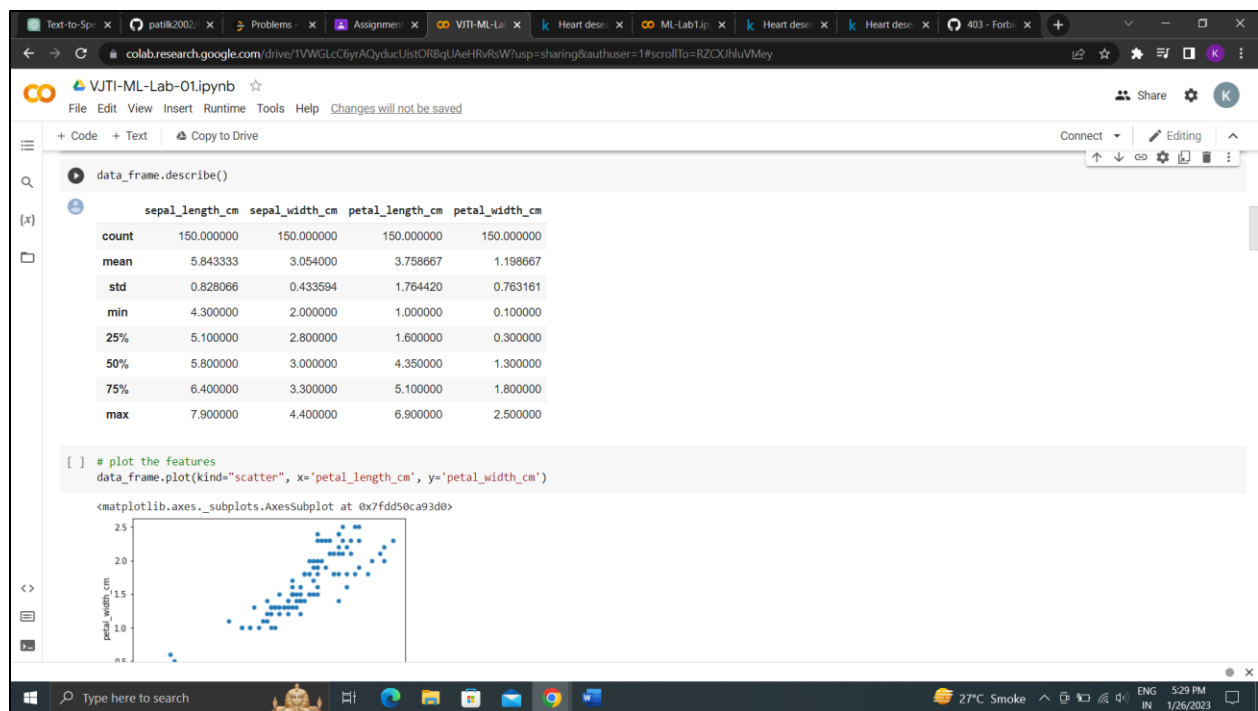
```
[ ] # https://pandas.pydata.org/docs/
import pandas as pd

data_frame = pd.read_csv('iris.data', names = ['sepal_length_cm', 'sepal_width_cm', 'petal_length_cm', 'petal_width_cm', 'class'])
data_frame.head()
```

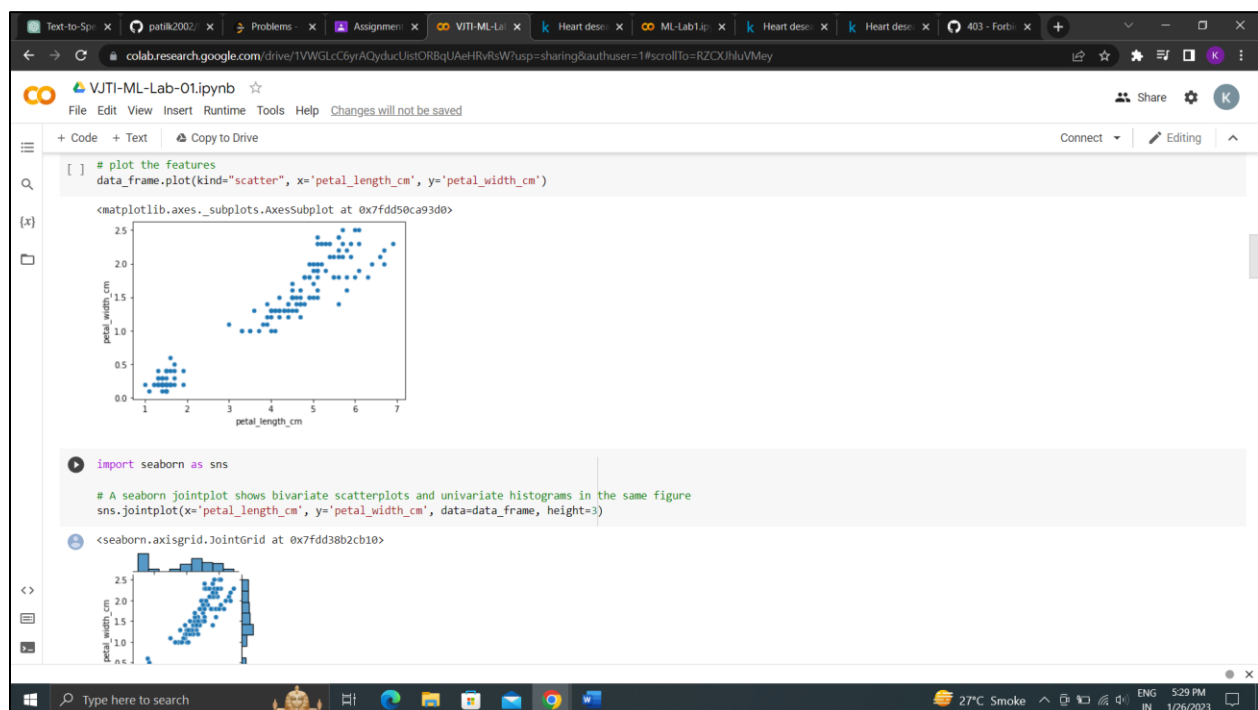
	sepal_length_cm	sepal_width_cm	petal_length_cm	petal_width_cm	class
0	5.1	3.5	1.4	0.2	Iris-setosa
1	4.9	3.0	1.4	0.2	Iris-setosa
2	4.7	3.2	1.3	0.2	Iris-setosa
3	4.6	3.1	1.5	0.2	Iris-setosa
4	5.0	3.6	1.4	0.2	Iris-setosa

```
[ ] data_frame.describe()
```

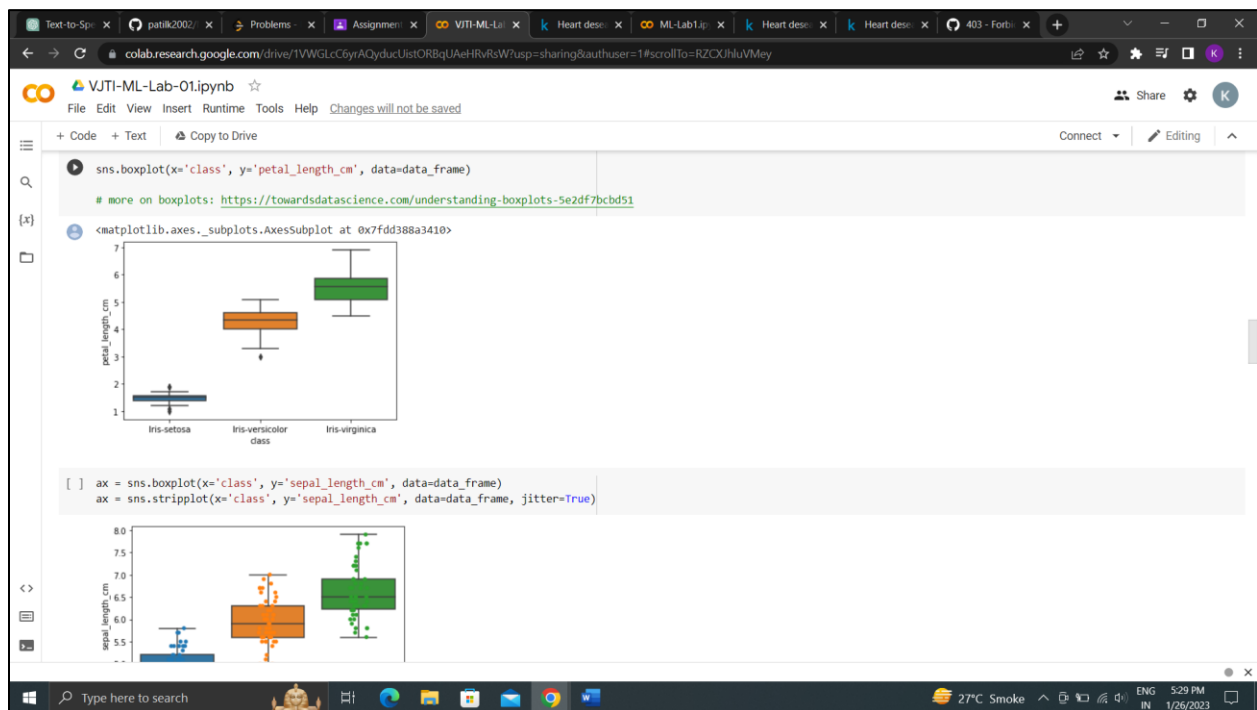
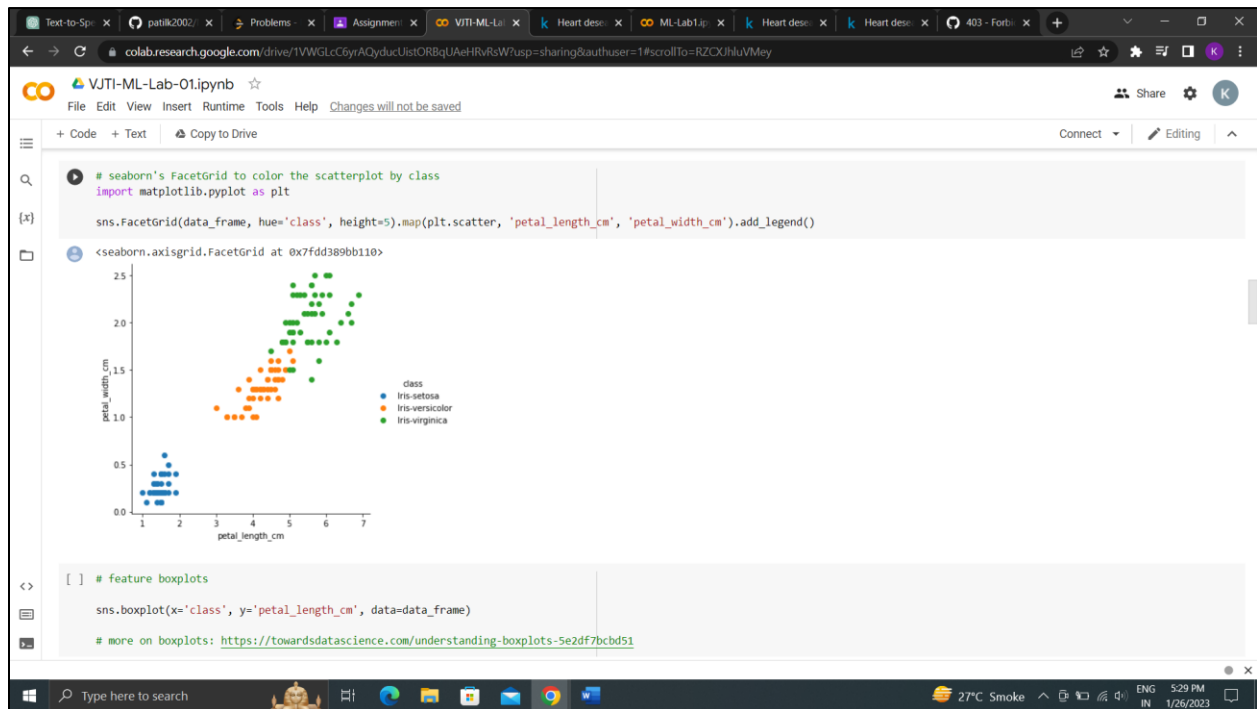
	sepal_length_cm	sepal_width_cm	petal_length_cm	petal_width_cm
count	150.000000	150.000000	150.000000	150.000000
mean	5.843333	3.054000	3.758667	1.198667
std	0.828066	0.433594	1.764420	0.763161

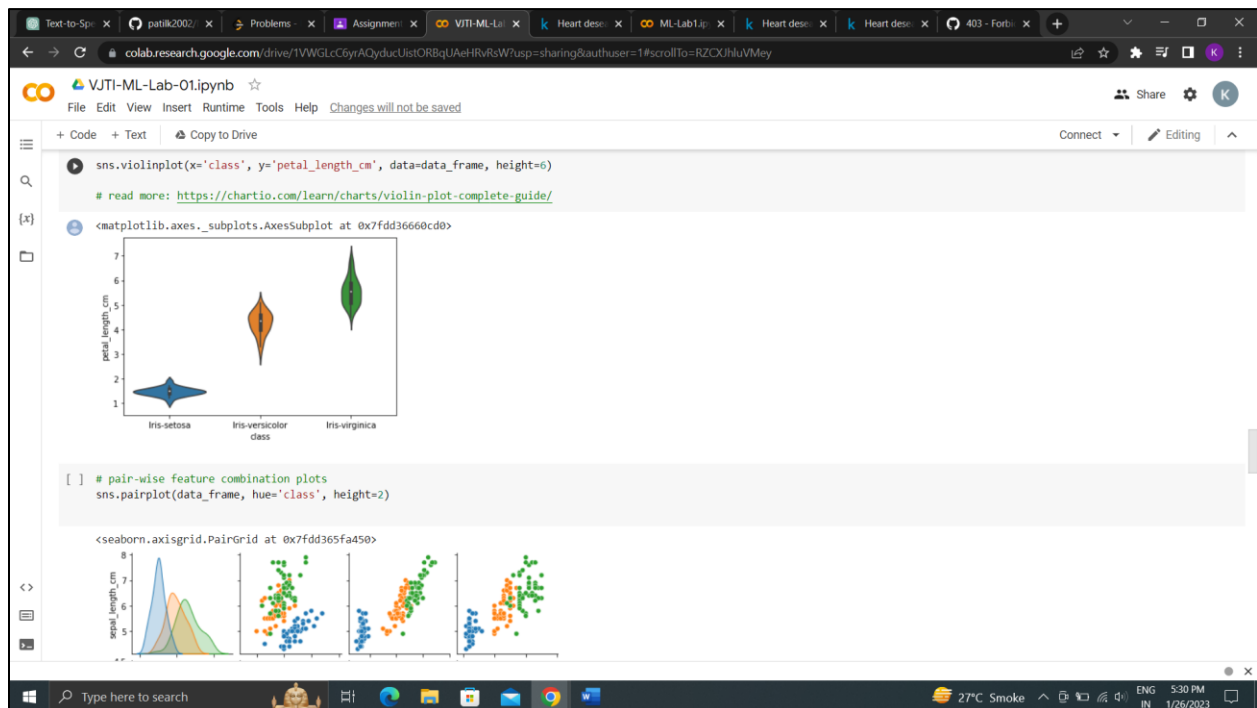
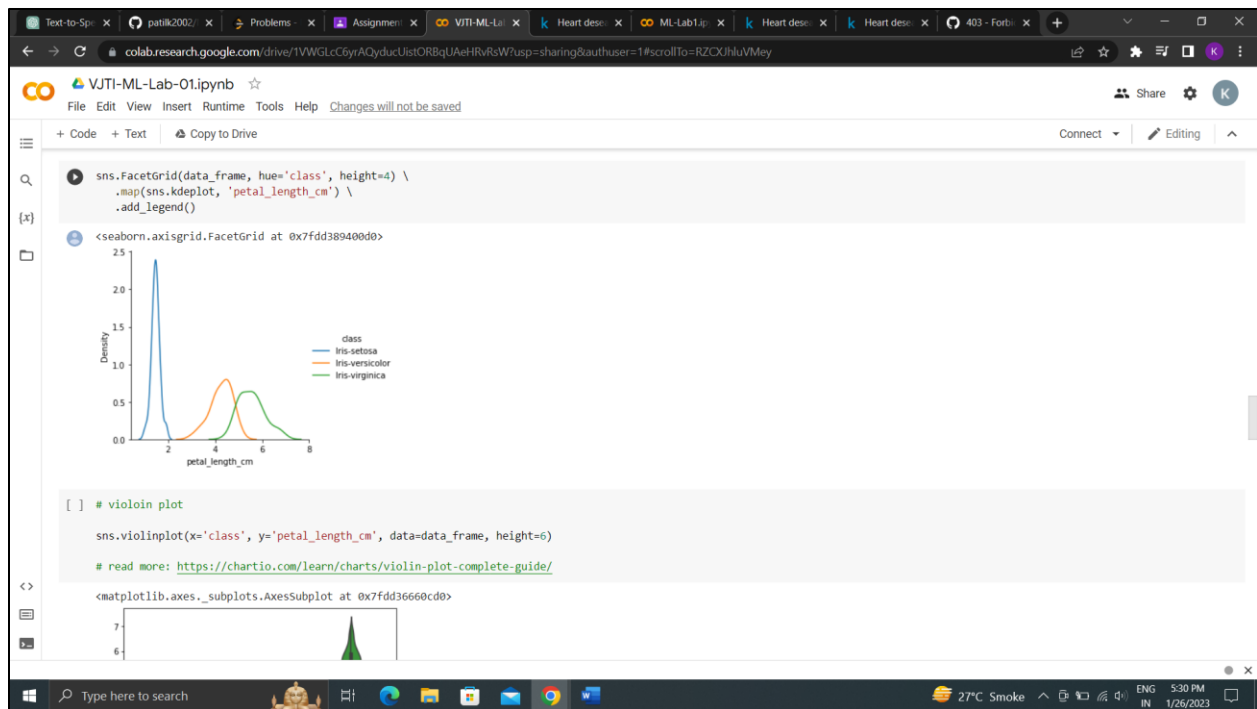


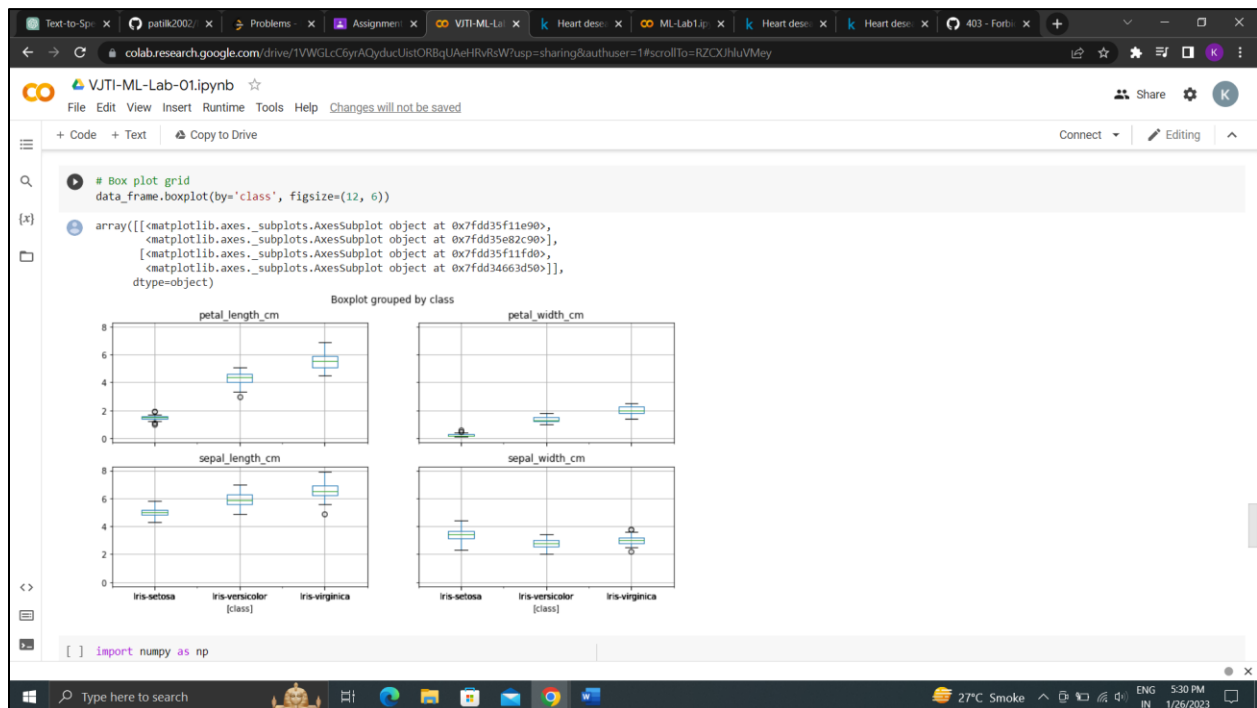
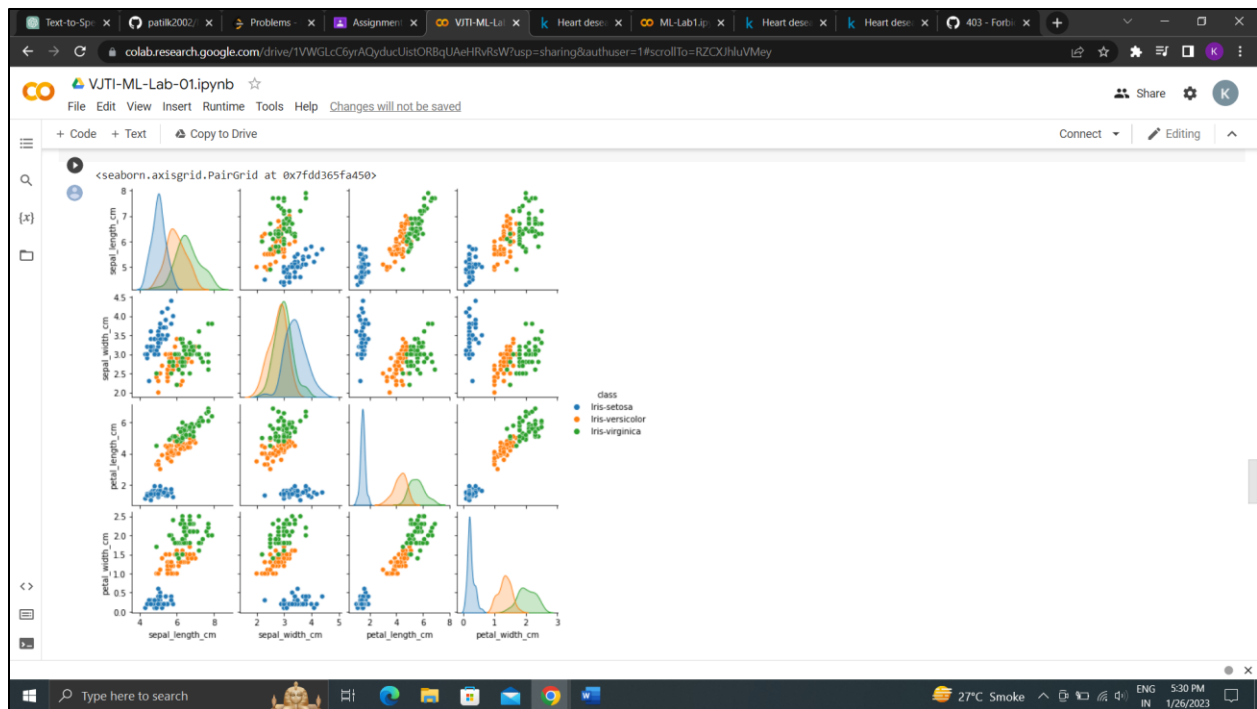
## Functions Performed Seaborn :



## Functions performed using matplotlib







```
Text-to-Sp... x patik2002... x Problems... x Assignment... x VJTI-ML-Lab... x ML-Lab1... x Heart dese... x Heart dese... x Heart dese... x 403 - Forb... x
colab.research.google.com/drive/1VWGLc6yrAQyducUistORbqUAeHRvR3W?usp=sharing&authuser=1#scrollTo=RZCXihluVMeY

VJTI-ML-Lab-01.ipynb
File Edit View Insert Runtime Tools Help Changes will not be saved
+ Code + Text Copy to Drive
Connect Editing
[ ] import numpy as np
import math

# Separating the data into dependent and independent variables
X = data_frame.iloc[:, :-1].values
y = data_frame.iloc[:, -1].values

# Create train and test split

def train_test_split(X, y, train=0.7):
    # first shuffle
    permute = np.random.permutation(len(X))
    X = X[permute]
    y = y[permute]

    X_train, y_train = X[0:math.floor(0.7*len(X))], y[0:math.floor(0.7*len(X))]
    X_test, y_test = X[math.floor(0.7*len(X)):], y[math.floor(0.7*len(X)):]

    return X_train, y_train, X_test, y_test

[ ] X_train, y_train, X_test, y_test = train_test_split(X, y, train=0.7)

# Logistic Regression

from sklearn.linear_model import LogisticRegression

# fit logistic regression model on training data
classifier = LogisticRegression()
classifier.fit(X_train, y_train)

# predict model on testing data
y_pred = classifier.predict(X_test)
```

```
Text-to-Sp... x patik2002... x Problems... x Assignment... x VJTI-ML-Lab... x ML-Lab1... x Heart dese... x Heart dese... x Heart dese... x 403 - Forb... x
colab.research.google.com/drive/1VWGLc6yrAQyducUistORbqUAeHRvR3W?usp=sharing&authuser=1#scrollTo=RZCXihluVMeY

VJTI-ML-Lab-01.ipynb
File Edit View Insert Runtime Tools Help Changes will not be saved
+ Code + Text Copy to Drive
Connect Editing
[ ] from sklearn.linear_model import LogisticRegression

# fit logistic regression model on training data
classifier = LogisticRegression()
classifier.fit(X_train, y_train)

# predict model on testing data
y_pred = classifier.predict(X_test)

from sklearn.metrics import confusion_matrix
from sklearn.metrics import classification_report

# Summary of the predictions made by the classifier
print(classification_report(y_test, y_pred, digits=3))
print(confusion_matrix(y_test, y_pred))

# Accuracy score
from sklearn.metrics import accuracy_score
print('accuracy is', accuracy_score(y_pred, y_test))

precision recall f1-score support
Iris-setosa 1.000 1.000 1.000 13
Iris-versicolor 0.941 1.000 0.970 16
Iris-virginica 1.000 0.938 0.968 16

accuracy 0.980 0.979 0.978 45
macro avg 0.979 0.978 0.978 45
weighted avg 0.979 0.978 0.978 45

[[13 0 0]
 [0 16 0]
 [0 1 15]]
accuracy is 0.9777777777777777
```

```
ML-Lab1.ipynb - Colaboratory
colab.research.google.com/drive/ly2kqf8QvofdR-O34K0AFVp0m0qg35x?authuser=1#scrollTo=xZnfU4l0liu

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

[1] !pip install -q kaggle

from google.colab import files
files.upload()

Choose Files | kaggle.json
  • kaggle.json(application/json) - 69 bytes, last modified: 1/26/2023 - 100% done
    Saving kaggle.json to kaggle.json
    {'kaggle.json': b'{"username": "kiranpatil1002", "key": "a2a5805c073a7803eaa00712e5245b2b"}'}
```

```
[3] mkdir ~/.kaggle

[4] !cp kaggle.json ~/.kaggle/

[5] !chmod 600 ~/.kaggle/kaggle.json

[6] !kaggle datasets list
```

ref	title	size	lastUpdated	downloadcount	voteCount	usabilityRating
ahsan81/hotel-reservations-classification-dataset	Hotel Reservations Dataset	480KB	2023-01-04 12:50:31	5041	175	1.0
senapatirajesh/netflix-tv-shows-and-movies	Latest Netflix TV shows and movies	1MB	2023-01-14 17:03:12	1490	46	0.9411765
johnny1994/divorce-rates-data-should-you-get-married	Divorce Rates Data: Should You Get Married?	22KB	2023-01-15 12:49:06	841	31	0.88235295
rakeshary/spotify-top-10000-streamed-songs	Spotify Top 10000 Streamed Songs	280KB	2023-01-02 08:17:15	2196	66	1.0
dieborigephit/fifa-world-cup-2022-complete-dataset	Fifa World Cup 2022: Complete Dataset	7KB	2022-12-18 22:51:11	7425	227	1.0
themrityunjaypathak/imdb-top-100-movies	IMDb Top 100 Movies	4KB	2023-01-11 17:15:09	1250	40	1.0
thedevasator/tesla-accident-fatalities-analysis-and-statistic	Tesla Deaths	18KB	2023-01-02 23:34:23	1130	30	0.9411765
rishikeshkonapure/home-loan-approval	Home Loan Approval	13KB	2023-01-12 06:28:57	1406	36	1.0
ayushnib/starcraft-players-dataset	Starcraft Players Dataset - Gamers Analytics	205KB	2023-01-18 23:14:29	227	29	1.0

0s completed at 4:54 PM

```
ML-Lab1.ipynb - Colaboratory
colab.research.google.com/drive/ly2kqf8QvofdR-O34K0AFVp0m0qg35x?authuser=1#scrollTo=xZnfU4l0liu

File Edit View Insert Runtime Tools Help All changes saved

+ Code + Text

[3] mkdir ~/.kaggle

[4] !cp kaggle.json ~/.kaggle/

[5] !chmod 600 ~/.kaggle/kaggle.json

!kaggle datasets list
```

ref	title	size	lastUpdated	downloadcount	voteCount	usabilityRating
ahsan81/hotel-reservations-classification-dataset	Hotel Reservations Dataset	480KB	2023-01-04 12:50:31	5041	175	1.0
senapatirajesh/netflix-tv-shows-and-movies	Latest Netflix TV shows and movies	1MB	2023-01-14 17:03:12	1490	46	0.9411765
johnny1994/divorce-rates-data-should-you-get-married	Divorce Rates Data: Should You Get Married?	22KB	2023-01-15 12:49:06	841	31	0.88235295
rakeshary/spotify-top-10000-streamed-songs	Spotify Top 10000 Streamed Songs	280KB	2023-01-02 08:17:15	2196	66	1.0
dieborigephit/fifa-world-cup-2022-complete-dataset	Fifa World Cup 2022: Complete Dataset	7KB	2022-12-18 22:51:11	7425	227	1.0
themrityunjaypathak/imdb-top-100-movies	IMDb Top 100 Movies	4KB	2023-01-11 17:15:09	1250	40	1.0
thedevasator/tesla-accident-fatalities-analysis-and-statistic	Tesla Deaths	18KB	2023-01-02 23:34:23	1130	30	0.9411765
rishikeshkonapure/home-loan-approval	Home Loan Approval	13KB	2023-01-12 06:28:57	1406	36	1.0
ayushnib/starcraft-players-dataset	Starcraft Players Dataset - Gamers Analytics	205KB	2023-01-18 23:14:29	227	29	1.0
ruddygunawan/per-capita-income-by-county-2021-vs-education	Per Capita Income by County (2021) vs. Education	89KB	2022-12-28 14:37:42	1455	31	1.0
nguyenthicamlai/population-2022	World Population (1955-2020)	121KB	2023-01-20 12:01:55	720	35	0.7647059
schmoyote/coffee-reviews-dataset	Coffee Reviews Dataset	569KB	2023-01-19 02:25:42	733	26	1.0
thedevasator/state-ut-wise-road-accidents-due-to-driver-violat	Road Accidents due to Driver Violations (India)	2KB	2023-01-06 14:29:20	818	33	1.0
thedevasator/global-fossil-co2-emissions-by-country-2002-2022	Emissions by country	2MB	2023-01-24 04:24:10	2250	72	1.0
sudhanshu17/mobilephone	MobilePhone's Dataset	876KB	2023-01-20 14:24:37	947	33	0.88235295
thedevasator/physical-strength-correlation-with-fear-related	Physical Strength & Fear-Related Personality	27KB	2023-01-24 04:33:10	461	23	1.0
omkargowda/suicide-rates-overview-1985-to-2021	Suicide Rates Overview (1985 to 2021)	539KB	2023-01-04 15:11:45	1531	64	1.0
gan2gan/1000-imdb-movies-20062016	1000 IMDb movies (2006-2016)	134KB	2023-01-17 12:32:25	603	30	1.0
arbazmohammad/world-airports-and-airlines-datasets	Worlds Airports and Airlines Datasets	760KB	2023-01-11 07:39:19	617	30	1.0
ahsan81/used-handheld-device-data	Used Phones & Tablets Pricing Dataset	77KB	2023-01-07 16:47:53	693	25	1.0

0s completed at 4:54 PM

Conclusion : From this experiment we learned to Download a dataset from Kaggle and Performed data exploration using Pandas, Seaborn, Matplotlib python libraries