



Project Cover Sheet

Assignment Title:	Python Final Term Project		
Assignment No:	01	Date of Submission:	8 August 2022
Course Title:	Programming in Python		
Course Code:	CSC4162	Section:	A
Semester:	Summer	2021-22	Course Teacher: AKINUL ISLAM JONY

Declaration and Statement of Authorship:

- I/we hold a copy of this Assignment/Case-Study, which can be produced if the original is lost/damaged.
- This Assignment/Case-Study is my/our original work and no part of it has been copied from any other student's work or from any other source except where due acknowledgement is made.
- No part of this Assignment/Case-Study has been written for me/us by any other person except where such collaboration has been authorized by the concerned teacher and is clearly acknowledged in the assignment.
- I/we have not previously submitted or currently submitting this work for any other course/unit.
- This work may be reproduced, communicated, compared and archived for the purpose of detecting plagiarism.
- I/we give permission for a copy of my/our marked work to be retained by the Faculty for review and comparison, including review by external examiners.
- I/we understand that Plagiarism is the presentation of the work, idea or creation of another person as though it is your own. It is a form of cheating and is a very serious academic offence that may lead to expulsion from the University. Plagiarized material can be drawn from, and presented in, written, graphic and visual form, including electronic data, and oral presentations. Plagiarism occurs when the origin of them material used is not appropriately cited.
- I/we also understand that enabling plagiarism is the act of assisting or allowing another person to plagiarize or to copy my/our work.

* Student(s) must complete all details except the faculty use part.

** Please submit all assignments to your course teacher or the office of the concerned teacher.

Group Name/No.:

No	Name	ID	Program	Signature
1	Mithun Chanda Shuvo	19-39592-1	BSc [CSE]	
2	Salah Ahommod Emu	19-39603-1	BSc [CSE]	
3	Md. Asif Mostofa	19-41266-3	BSc [CSE]	
4	Nusrat Jahan Mou	19-39583-1	BSc [CSE]	
5			Choose an item.	
6			Choose an item.	
7			Choose an item.	
8			Choose an item.	
9			Choose an item.	
10			Choose an item.	

Faculty use only

FACULTY COMMENTS	Marks Obtained	
	Total Marks	

Section 1: Project Overview

The goal of this study is to forecast a person's likelihood of having a stroke based on personal information and figure out if they are related to the case in any manner. To make the prediction, we will use a dataset with over 4000 cases to train our model. There are many symptoms of strokes. All these symptoms resemble different diseases like it occurs in aging persons, so it becomes a challenging task to get a correct diagnosis. But as time passes, a lot of research data and patients records of hospitals are available. There are many open sources for accessing the patient's records and research can be conducted so that various computer technologies could be used for doing the correct diagnosis of the patients and detecting this disease to stop it from becoming fatal. Nowadays it is well known that machine learning and artificial intelligence are playing a huge role in the medical industry. We can use different machine learning models to diagnose the disease and classify or predict the results. A complete data analysis can easily be done using machine learning models. Models can be trained for knowledge predictions and medical records can be transformed and analyzed more deeply for better predictions.

The steps we followed in this project:

1. Data Collection
2. Data Cleaning
3. Exploratory data analysis
4. Feature selection
5. Machine Learning Model develop
6. Result analysis

Section 2: Dataset Overview

The valid URL of the dataset that we made use of for our project is given below:

<https://www.kaggle.com/datasets/fedesoriano/stroke-prediction-dataset>

The dataset that we have used in our project is called ‘Stroke Prediction Dataset’. This dataset holds columns named as ‘id’, ‘gender’, ‘age’, ‘hypertension’, ‘heart_disease’, ‘ever_married’, ‘work_type’, ‘residence_type’, ‘avg_glucose_level’, ‘bmi’, ‘smoking_status’, ‘stroke’.

This dataset will primarily figure out whether a person is susceptible to stroke based on personal information such as id, gender, age, whether the person is married, his working kind, and his residence type. The dataset also includes the individual's medical history, such as hypertension, heart disease, average glucose level, BMI, smoking status, and stroke. Every row in the collection holds information about the patient.

Section 3: Data Preprocessing and Exploratory data analysis

- data preprocessing/cleaning steps

Import all the Library and load dataset

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
from sklearn import datasets # for using built-in datasets
from sklearn import metrics # for checking the model accuracy
df = pd.read_csv('D:\\Dataset\\archive\\full_data.csv')
df
```

	gender	age	hypertension	heart_disease	ever_married	work_type	Residence_type	avg_glucose_level	bmi	smoking_status	stroke
0	Male	67.0	0	1	Yes	Private	Urban	228.69	36.6	formerly smoked	1
1	Male	80.0	0	1	Yes	Private	Rural	105.92	32.5	never smoked	1
2	Female	49.0	0	0	Yes	Private	Urban	171.23	34.4	smokes	1
3	Female	79.0	1	0	Yes	Self-employed	Rural	174.12	24.0	never smoked	1
4	Male	81.0	0	0	Yes	Private	Urban	186.21	29.0	formerly smoked	1
...
4976	Male	41.0	0	0	No	Private	Rural	70.15	29.8	formerly smoked	0
4977	Male	40.0	0	0	Yes	Private	Urban	191.15	31.1	smokes	0
4978	Female	45.0	1	0	Yes	Govt_job	Rural	95.02	31.8	smokes	0
4979	Male	40.0	0	0	Yes	Private	Rural	83.94	30.0	smokes	0
4980	Female	80.0	1	0	Yes	Private	Urban	83.75	29.1	never smoked	0

4981 rows × 11 columns

Jupyter interface showing data preprocessing steps:

Data PreProcessing

In [56]: # Return number of missing values for each column
df.isnull().sum()

Out[56]:

gender	0
age	0
hypertension	0
heart_disease	0
ever_married	0
Residence_type	0
avg_glucose_level	0
bmi	0
smoking_status	0
stroke	0
dtype:	int64

In [57]: # Check what percentage of each column's data is missing
df.isnull().sum()/len(df)

Out[57]:

gender	0.0
age	0.0
hypertension	0.0
heart_disease	0.0
ever_married	0.0
Residence_type	0.0
avg_glucose_level	0.0
bmi	0.0
smoking_status	0.0

In [58]:

```
from sklearn.preprocessing import LabelEncoder  
le = LabelEncoder()  
  
dfle = df  
dfle.gender = le.fit_transform(dfle.gender)  
dfle.ever_married = le.fit_transform(dfle.ever_married)  
dfle.Residence_type = le.fit_transform(dfle.Residence_type)  
dfle.smoking_status = le.fit_transform(dfle.smoking_status)  
dfle = dfle  
dfle.head()
```

Out[58]:

	gender	age	hypertension	heart_disease	ever_married	Residence_type	avg_glucose_level	bmi	smoking_status	stroke
0	1	67.0	0	1	1	1	228.69	36.6	1	1
1	1	80.0	0	1	1	0	105.92	32.5	2	1
2	0	49.0	0	0	1	1	171.23	34.4	3	1
3	0	79.0	1	0	1	0	174.12	24.0	2	1
4	1	81.0	0	0	1	1	186.21	29.0	1	1

In [59]: # Return basic statistics summary of the dataset
print(df.describe())

count

gender	4981.000000
age	4981.000000
hypertension	4981.000000
heart_disease	4981.000000
ever_married	4981.000000

mean

gender	0.416382
age	43.419859
hypertension	0.096165
heart_disease	0.055210
ever_married	0.658502

std

gender	0.493008
age	22.662755
hypertension	0.294848
heart_disease	0.228412
ever_married	0.474260

min

gender	0.000000
age	0.000000
hypertension	0.000000
heart_disease	0.000000
ever_married	0.000000

25%

gender	0.000000
age	25.000000
hypertension	0.000000
heart_disease	0.000000
ever_married	0.000000

50%

gender	0.000000
age	45.000000
hypertension	0.000000
heart_disease	0.000000
ever_married	1.000000

75%

gender	1.000000
age	61.000000
hypertension	0.000000
heart_disease	0.000000
ever_married	1.000000

max

gender	1.000000
age	82.000000
hypertension	1.000000
heart_disease	1.000000
ever_married	1.000000

Residence_type

count	4981.000000
mean	0.508332
std	0.499981
min	0.000000
25%	0.000000
50%	0.000000
75%	1.000000
max	1.000000

avg_glucose_level

count	4981.000000
mean	105.943562
std	45.075373
min	55.120000
25%	77.230000
50%	91.850000
75%	113.860000
max	271.740000

bmi

count	4981.000000
mean	28.498173
std	6.790464
min	14.000000
25%	23.700000
50%	28.100000
75%	32.600000
max	48.900000

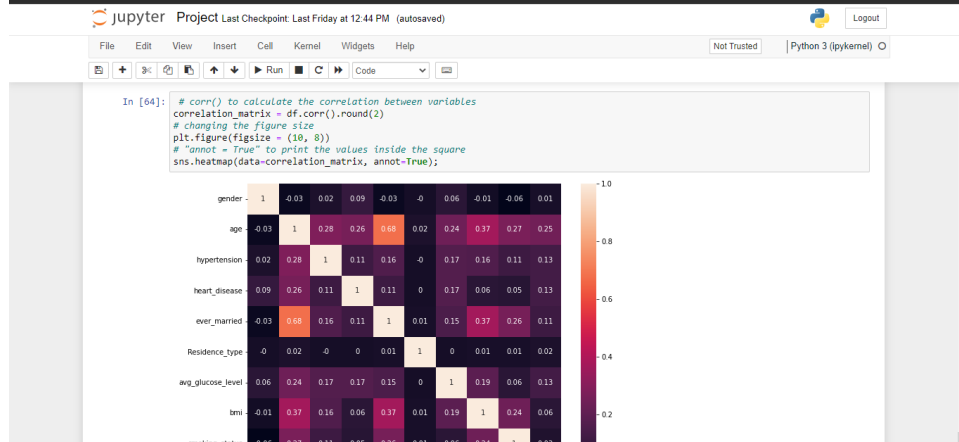
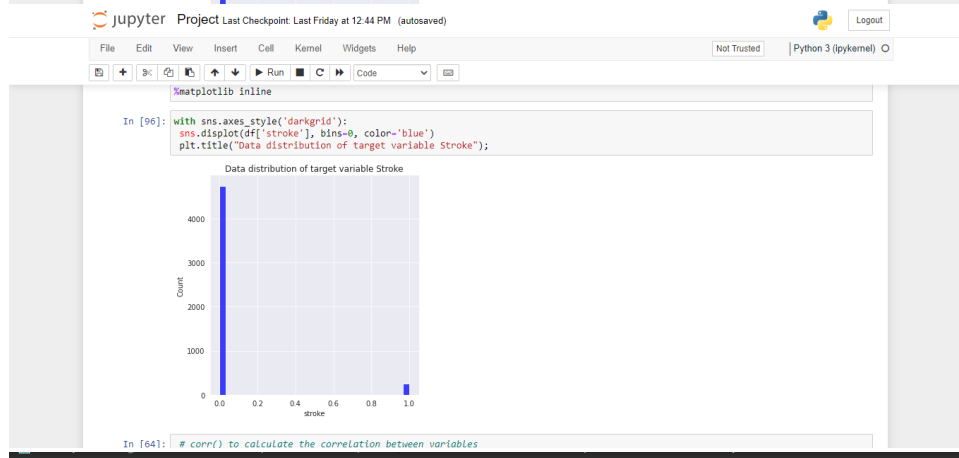
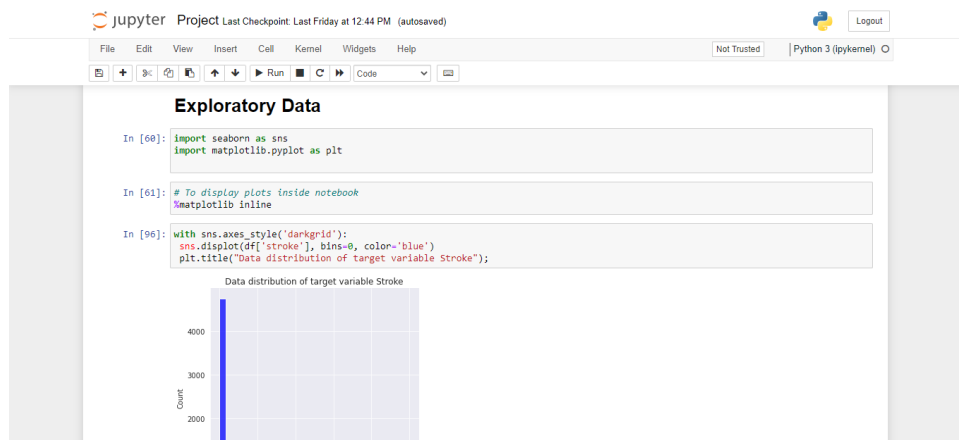
smoking_status

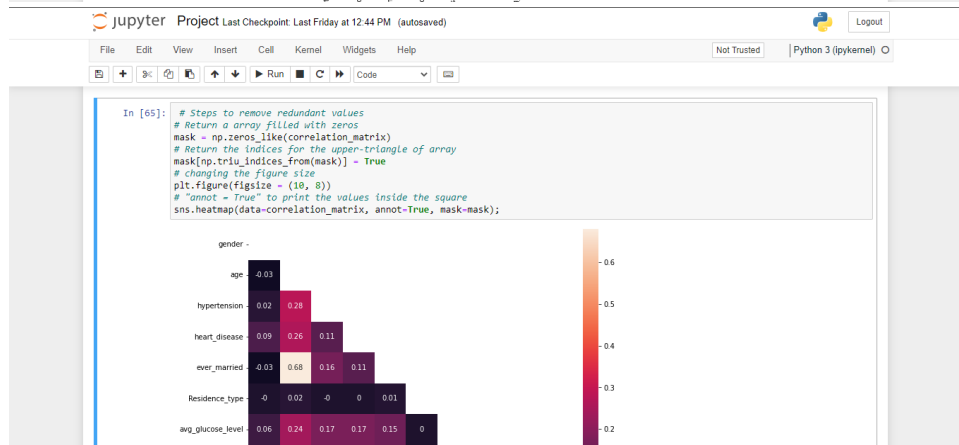
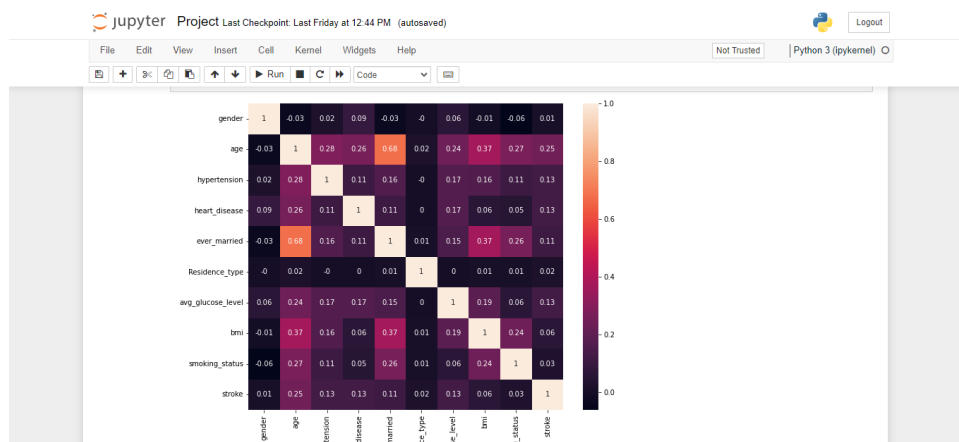
count	4981.000000
mean	1.379442
std	1.072180
min	0.000000
25%	0.000000
50%	0.000000
75%	2.000000
max	3.000000

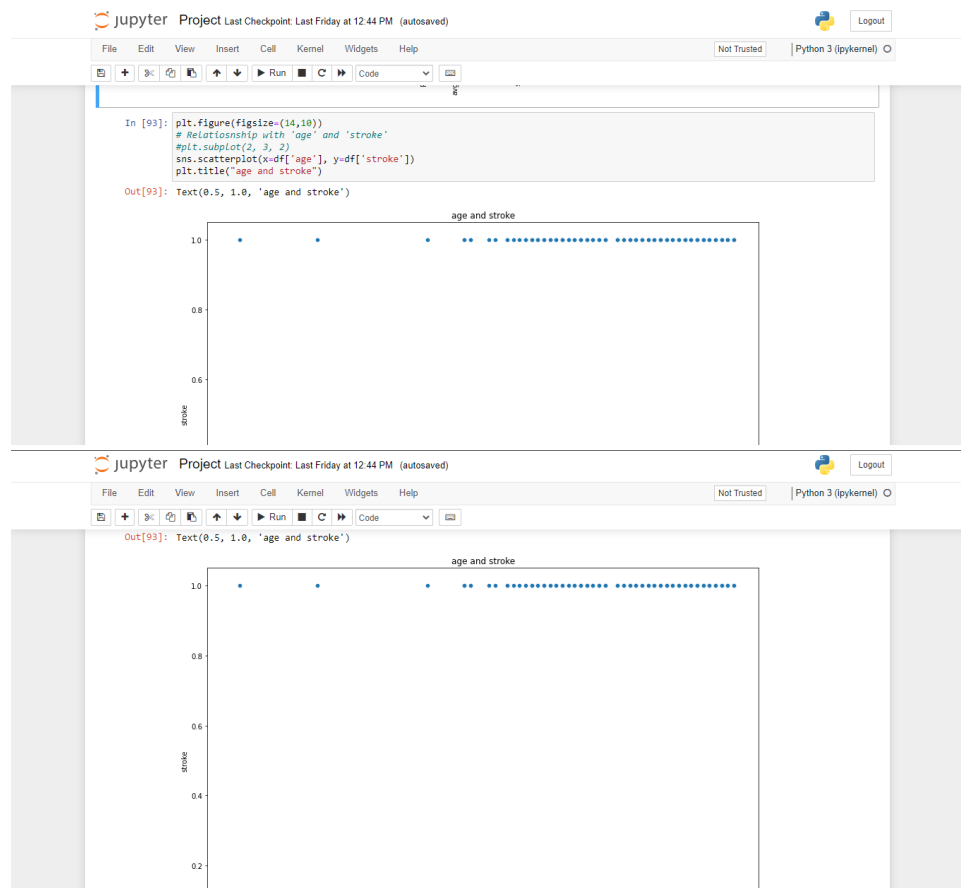
stroke

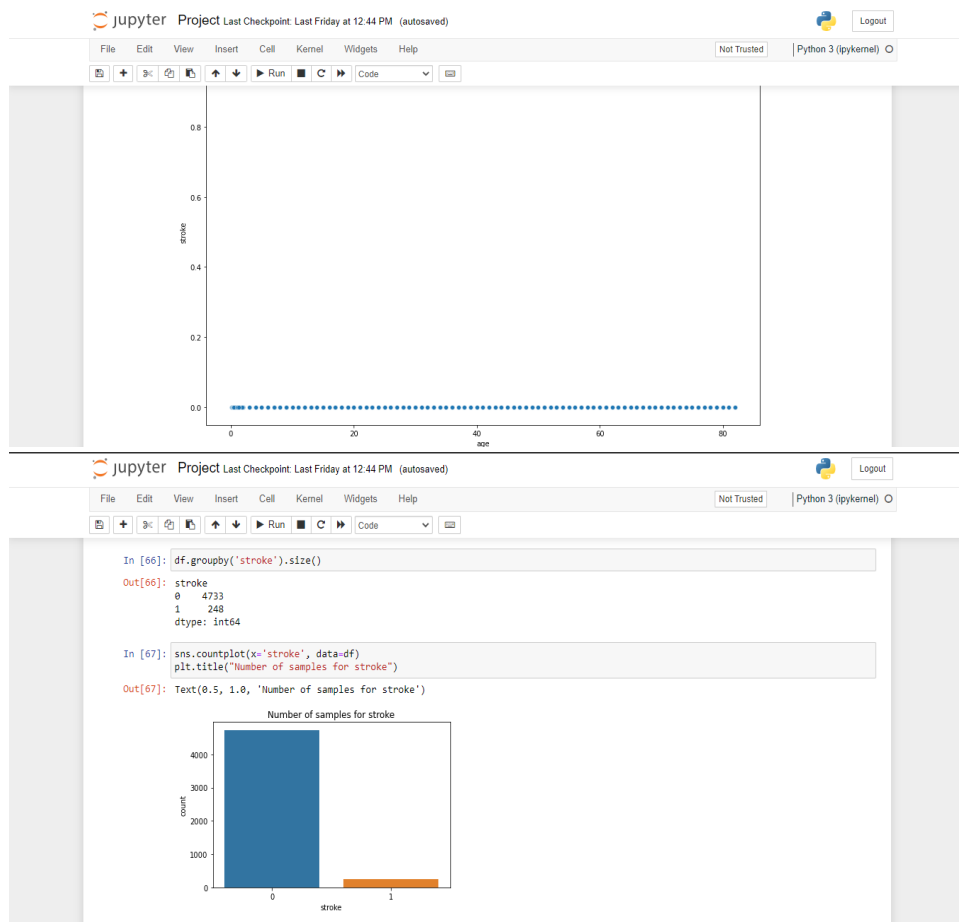
count	4981.000000
mean	0.049789
std	0.217531
min	0.000000
25%	0.000000
50%	0.000000

- Exploratory data analysis









section 4: Model Development

jupyter Project Last Checkpoint: Last Friday at 12:44 PM (autosaved) Logout

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

```
In [68]: X = df.drop('stroke',axis='columns')
y = df.stroke
```

Split the dataset

```
In [69]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.3, random_state = 16)
print("X_train shape: ", X_train.shape)
print("X_test shape: ", X_test.shape)
print("y_train shape: ", y_train.shape)
print("y_test shape: ", y_test.shape)

X_train shape: (3486, 9)
X_test shape: (1495, 9)
y_train shape: (3486,)
y_test shape: (1495,)
```

jupyter Project Last Checkpoint: Last Friday at 12:44 PM (autosaved) Logout

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

Create Model: Support Vector Machine (SVM)

```
In [72]: # importing the necessary package to use the classification algorithm
from sklearn import svm #for Support Vector Machine (SVM) Algorithm
model_svm = svm.SVC() #select the algorithm
model_svm.fit(X_train, y_train) #train the model with the training dataset
y_prediction_svm = model_svm.predict(X_test) # pass the testing data to the trained model
# checking the accuracy of the algorithm.
# by comparing predicted output by the model and the actual output
score_svm = metrics.accuracy_score(y_prediction_svm, y_test).round(4)
print("-----")
print('The accuracy of the SVM is: {}'.format(score_svm))
print("-----")
# save the accuracy score
score = set()
score.add(('SVM', score_svm))

-----
The accuracy of the SVM is: 0.9525
-----
```

Decision Tree

```
In [73]: # importing the necessary package to use the classification algorithm
from sklearn.tree import DecisionTreeClassifier #for using Decision Tree Algorithm
```

jupyter Project Last Checkpoint: Last Friday at 12:44 PM (autosaved) Logout

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

Decision Tree

```
In [73]: # importing the necessary package to use the classification algorithm
from sklearn.tree import DecisionTreeClassifier #for using Decision Tree Algorithm
model_dt = DecisionTreeClassifier(random_state=4)
model_dt.fit(X_train, y_train) #train the model with the training dataset
y_prediction_dt = model_dt.predict(X_test) #pass the testing data to the trained model
# checking the accuracy of the algorithm.
# by comparing predicted output by the model and the actual output
score_dt = metrics.accuracy_score(y_prediction_dt, y_test).round(4)
print("-----")
print('The accuracy of the DT is: {}'.format(score_dt))
print("-----")
# save the accuracy score
score.add(('DT', score_dt))

-----
The accuracy of the DT is: 0.905
-----
```

K Nearest Neighbours (KNN)

```
In [74]: # importing the necessary package to use the classification algorithm
from sklearn.neighbors import KNeighborsClassifier # for K nearest neighbours
from sklearn.linear_model import LogisticRegression # for Logistic Regression algorithm
model_knn = KNeighborsClassifier(n_neighbors=3) # 3 neighbours for putting the new data into a class
```

jupyter Project Last Checkpoint: Last Friday at 12:44 PM (autosaved) Logout

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

K Nearest Neighbours (KNN)

```
In [74]: # importing the necessary package to use the classification algorithm
from sklearn.neighbors import KNeighborsClassifier # for K nearest neighbours
# from sklearn.linear_model import LogisticRegression # for Logistic Regression algorithm
model_knn = KNeighborsClassifier(n_neighbors=3) # 3 neighbours for putting the new data into a class
model_knn.fit(X_train, y_train) # train the model with the training dataset
y_prediction_knn = model_knn.predict(X_test) # pass the testing data to the trained model
# checking the accuracy of the algorithm.
# by comparing predicted output by the model and the actual output
score_knn = metrics.accuracy_score(y_prediction_knn, y_test).round(4)
print("-----")
print('The accuracy of the KNN is: {}'.format(score_knn))
print("-----")
# save the accuracy score
score.add('KNN', score_knn)

-----
The accuracy of the KNN is: 0.9425
-----
```

Logistic Regression

```
In [75]: # importing the necessary package to use the classification algorithm
```

jupyter Project Last Checkpoint: Last Friday at 12:44 PM (autosaved) Logout

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

Logistic Regression

```
In [75]: # importing the necessary package to use the classification algorithm
from sklearn.linear_model import LogisticRegression # for Logistic Regression algorithm
model_lr = LogisticRegression()
model_lr.fit(X_train, y_train) # train the model with the training dataset
y_prediction_lr = model_lr.predict(X_test) # pass the testing data to the trained model
# checking the accuracy of the algorithm.
# by comparing predicted output by the model and the actual output
score_lr = metrics.accuracy_score(y_prediction_lr, y_test).round(4)
print("-----")
print('The accuracy of the LR is: {}'.format(score_lr))
print("-----")
# save the accuracy score
score.add('LR', score_lr)

-----
The accuracy of the LR is: 0.9518
-----
```

The image displays two screenshots of a Jupyter Notebook interface. The top screenshot shows a code cell titled "Naive Bayes" where a Gaussian Naive Bayes model is trained on training data and evaluated on test data. The output shows an accuracy of 0.8749. The bottom screenshot shows a subsequent code cell that prints the accuracy scores for five different models: DT (0.985), LR (0.9518), SVM (0.9525), KNN (0.9425), and NB (0.8749).

```
n_iter_1 = _check_optimize_result(

Naive Bayes

In [76]: # importing the necessary package to use the classification algorithm
from sklearn.naive_bayes import GaussianNB
model_nb = GaussianNB()
model_nb.fit(X_train, y_train) #train the model with the training dataset
y_prediction_nb = model_nb.predict(X_test) #pass the testing data to the trained model
# checking the accuracy of the algorithm.
# by comparing predicted output by the model and the actual output
score_nb = metrics.accuracy_score(y_prediction_nb, y_test).round(4)
print("-----")
print('The accuracy of the NB is: {}'.format(score_nb))
print("-----")
# save the accuracy score
score.add(('NB', score_nb))

-----
The accuracy of the NB is: 0.8749
-----

In [79]: print("The accuracy scores of different Models:")
print("-----")
for s in score:
    print(s)

The accuracy scores of different Models:
-----
('DT', 0.985)
('LR', 0.9518)
('SVM', 0.9525)
('KNN', 0.9425)
('NB', 0.8749)

In [ ]:
```

Section 5: Discussion & Conclusion

Here in our project, we have used 5 models.

- Naive Bayes
- KNN
- Decision Tree
- Logistic Regression
- SVM

We can see from the comparison among the models that

By training the model with training datasets,

- The accuracy score of Naive Bayes (NB) is 0.8749.

- The accuracy score of KNN is 0.9425
- The accuracy score of Decision tree (DT) is 0.905
- The accuracy score of Logistic Regression (LR) is 0.9518
- The accuracy score of SVM is 0.9525

Here we compared predicted output of all the models with their actual output. So, the accuracy score of SVM model was the highest which is 0.9525 and it is most efficient. The least accuracy score is found from the Naive Bayes which is 0.8749.