

PROJECT DEVELOPMENT PHASE

SPRINT 2 – CODE AND TESTCASE

Date	10 November 2022
Team ID	PNT2022TMID18451
Project	Flight delay prediction using Machine learning
Marks	8 Marks

We have performed the uploading the Dataset and performed the Data Pre-processing and also we have split the dataset into train data and Test dataset in this Sprint development phase.

Jupyter notebook :

Screenshots :

The screenshot shows a Jupyter Notebook running on a local host. The browser address bar indicates the URL is localhost:8888/notebooks/PNT2022TMID18451-Sprint%202.ipynb. The Jupyter interface includes a menu bar (File, Edit, View, Insert, Cell, Kernel, Widgets, Help) and a toolbar with icons for file operations, running, and saving. The notebook title is 'PNT2022TMID18451-Sprint 2' with a 'Last Checkpoint: 2 minutes ago (autosaved)' status. The code in the first cell imports various Python libraries: sys, numpy (np), pandas (pd), seaborn (sns), pickle, sklearn, and imblearn. It also imports specific modules from sklearn for preprocessing (LabelEncoder, OneHotEncoder), model selection (train_test_split), and metrics (DecisionTreeClassifier, accuracy_score). The second cell shows the code to read a CSV file named 'flightdata.csv' into a pandas DataFrame. The output of the second cell is displayed as a table with columns: YEAR, QUARTER, MONTH, DAY_OF_MONTH, DAY_OF_WEEK, UNIQUE_CARRIER, TAIL_NUM, FL_NUM, ORIGIN_AIRPORT_ID, ORIGIN, CRS_ARR_TIME. The first five rows of data are shown.

```
In [2]: import sys
import numpy as np
import pandas as pd
import seaborn as sns
import pickle
import sklearn
# import imblearn
%matplotlib inline
from sklearn.preprocessing import LabelEncoder
from sklearn.preprocessing import OneHotEncoder
from sklearn.model_selection import train_test_split
from sklearn.tree import DecisionTreeClassifier
from sklearn.metrics import accuracy_score
import sklearn.metrics as metrics

In [3]: data=pd.read_csv('flightdata.csv')
data
```

	YEAR	QUARTER	MONTH	DAY_OF_MONTH	DAY_OF_WEEK	UNIQUE_CARRIER	TAIL_NUM	FL_NUM	ORIGIN_AIRPORT_ID	ORIGIN	CRS_ARR_TIME
0	2016	1	1	1	5	DL	N836DN	1399	10397	ATL	2143
1	2016	1	1	1	5	DL	N964DN	1476	11433	DTW	1435
2	2016	1	1	1	5	DL	N813DN	1597	10397	ATL	1215
3	2016	1	1	1	5	DL	N587NW	1768	14747	SEA	1335
4	2016	1	1	1	5	DL	N836DN	1823	14747	SEA	607

Browser tabs: (4) WhatsApp, Home Page - Select or create a n..., PNT2022TMD18451-Sprint 2 - J... | Address bar: localhost:8888/notebooks/PNT2022TMD18451-Sprint%202.ipynb | Welcome, SHAIK S... | Gmail | YouTube | HireMee Assessment | WhatsApp | GDB online Debug... | Apply-LCP - IBM

jupyter PNT2022TMD18451-Sprint 2 Last Checkpoint: 3 minutes ago (autosaved) Logout

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

11226 2016 4 12 30 5 DL N940DL 1715 11433 DTW ... 1223
11227 2016 4 12 30 5 DL N836DN 1770 14747 SEA ... 2046
11228 2016 4 12 30 5 DL N583NW 1823 11433 DTW ... 2210
11229 2016 4 12 30 5 DL N554NW 1901 10397 ATL ... 1800
11230 2016 4 12 30 5 DL N843DN 2005 10397 ATL ... 925

11231 rows x 26 columns

```
In [4]: data.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11231 entries, 0 to 11230
Data columns (total 26 columns):
 #   Column              Non-Null Count  Dtype
---  -
 0   YEAR                11231 non-null  int64
 1   QUARTER             11231 non-null  int64
 2   MONTH              11231 non-null  int64
 3   DAY_OF_MONTH        11231 non-null  int64
 4   DAY_OF_WEEK         11231 non-null  int64
 5   UNIQUE_CARRIER     11231 non-null  object
 6   TAIL_NUM            11231 non-null  object
 7   FL_NUM              11231 non-null  int64
 8   ORIGIN_AIRPORT_ID   11231 non-null  int64
 9   ORIGIN              11231 non-null  object
```

Browser tabs: (4) WhatsApp, Home Page - Select or create a n..., PNT2022TMD18451-Sprint 2 - J... | Address bar: localhost:8888/notebooks/PNT2022TMD18451-Sprint%202.ipynb | Welcome, SHAIK S... | Gmail | YouTube | HireMee Assessment | WhatsApp | GDB online Debug... | Apply-LCP - IBM

jupyter PNT2022TMD18451-Sprint 2 Last Checkpoint: 3 minutes ago (autosaved) Logout

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

```
0   YEAR                11231 non-null  int64
 1   QUARTER             11231 non-null  int64
 2   MONTH              11231 non-null  int64
 3   DAY_OF_MONTH        11231 non-null  int64
 4   DAY_OF_WEEK         11231 non-null  int64
 5   UNIQUE_CARRIER     11231 non-null  object
 6   TAIL_NUM            11231 non-null  object
 7   FL_NUM              11231 non-null  int64
 8   ORIGIN_AIRPORT_ID   11231 non-null  int64
 9   ORIGIN              11231 non-null  object
10  DEST_AIRPORT_ID     11231 non-null  int64
11  DEST                11231 non-null  object
12  CRS_DEP_TIME        11231 non-null  int64
13  DEP_TIME            11124 non-null  float64
14  DEP_DELAY           11124 non-null  float64
15  DEP_DELAY           11124 non-null  float64
16  CRS_ARR_TIME        11231 non-null  int64
17  ARR_TIME            11116 non-null  float64
18  ARR_DELAY           11043 non-null  float64
19  ARR_DELAY           11043 non-null  float64
20  CANCELLED           11231 non-null  float64
21  DIVERTED            11231 non-null  float64
22  CRS_ELAPSED_TIME    11231 non-null  float64
23  ACTUAL_ELAPSED_TIME 11043 non-null  float64
24  DISTANCE            11231 non-null  float64
25  Unnamed: 25         0 non-null      float64
dtypes: float64(12), int64(10), object(4)
memory usage: 2.2+ MB
```

localhost:8888/notebooks/PNT2022TMID18451-Sprint2-202.ipynb

Welcome, SHAIK S... | Gmail | YouTube | HireMee Assessment | WhatsApp | GDB online Debug... | Apply-LCP - IBM

Jupyter PNT2022TMID18451-Sprint 2 Last Checkpoint: 3 minutes ago (autosaved) Logout

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

```
In [5]: data.describe()
```

```
Out[5]:
```

	YEAR	QUARTER	MONTH	DAY_OF_MONTH	DAY_OF_WEEK	FL_NUM	ORIGIN_AIRPORT_ID	DEST_AIRPORT_ID	CRS_DEP_TIME	DEP_
count	11231.0	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11124.00
mean	2016.0	2.544475	6.628973	15.790758	3.960199	1334.325617	12334.516695	12302.274508	1320.798326	1327.18
std	0.0	1.090701	3.354678	8.782056	1.995257	811.875227	1595.026510	1601.988550	490.737845	500.30
min	2016.0	1.000000	1.000000	1.000000	1.000000	7.000000	10397.000000	10397.000000	10.000000	1.00
25%	2016.0	2.000000	4.000000	8.000000	2.000000	624.000000	10397.000000	10397.000000	905.000000	905.00
50%	2016.0	3.000000	7.000000	16.000000	4.000000	1267.000000	12478.000000	12478.000000	1320.000000	1324.00
75%	2016.0	3.000000	9.000000	23.000000	6.000000	2032.000000	13487.000000	13487.000000	1735.000000	1739.00
max	2016.0	4.000000	12.000000	31.000000	7.000000	2853.000000	14747.000000	14747.000000	2359.000000	2400.00

8 rows x 22 columns

```
In [6]: data.isnull().sum()
```

```
Out[6]:
```

YEAR	0
QUARTER	0
MONTH	0
DAY_OF_MONTH	0
DAY_OF_WEEK	0
UNIQUE_CARRIER	0
...	...

localhost:8888/notebooks/PNT2022TMID18451-Sprint2-202.ipynb

Welcome, SHAIK S... | Gmail | YouTube | HireMee Assessment | WhatsApp | GDB online Debug... | Apply-LCP - IBM

Jupyter PNT2022TMID18451-Sprint 2 Last Checkpoint: 3 minutes ago (autosaved) Logout

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

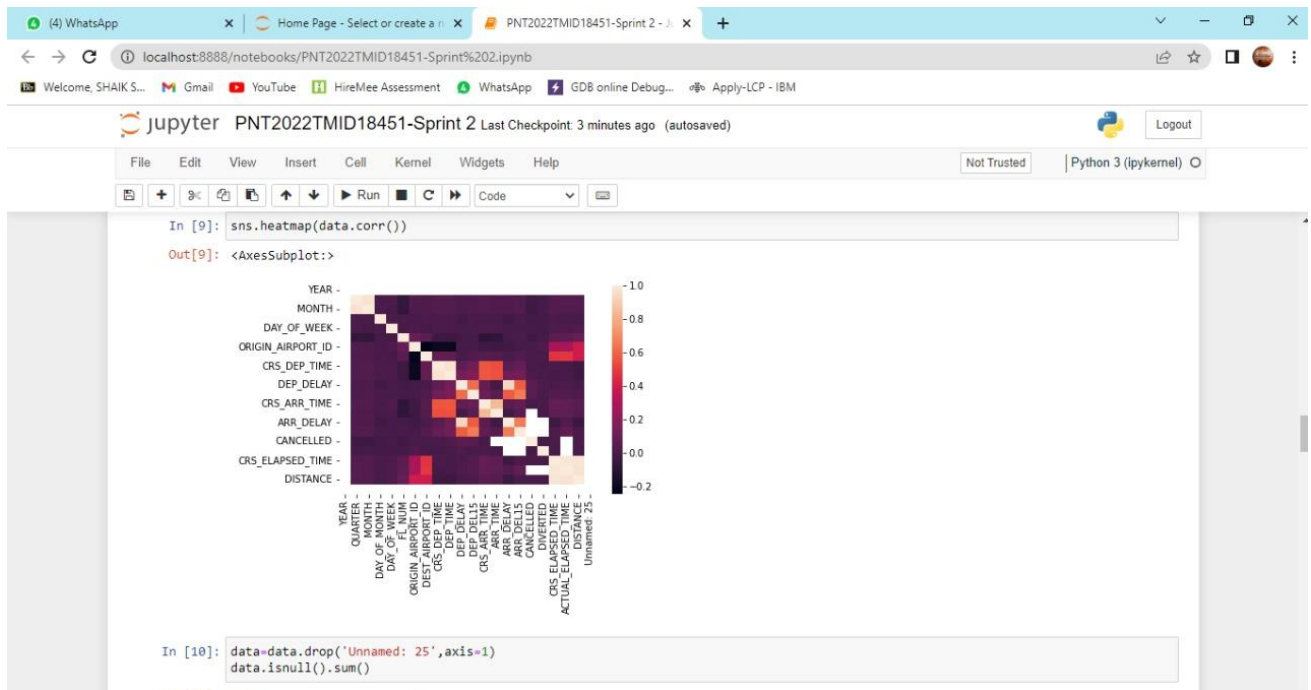
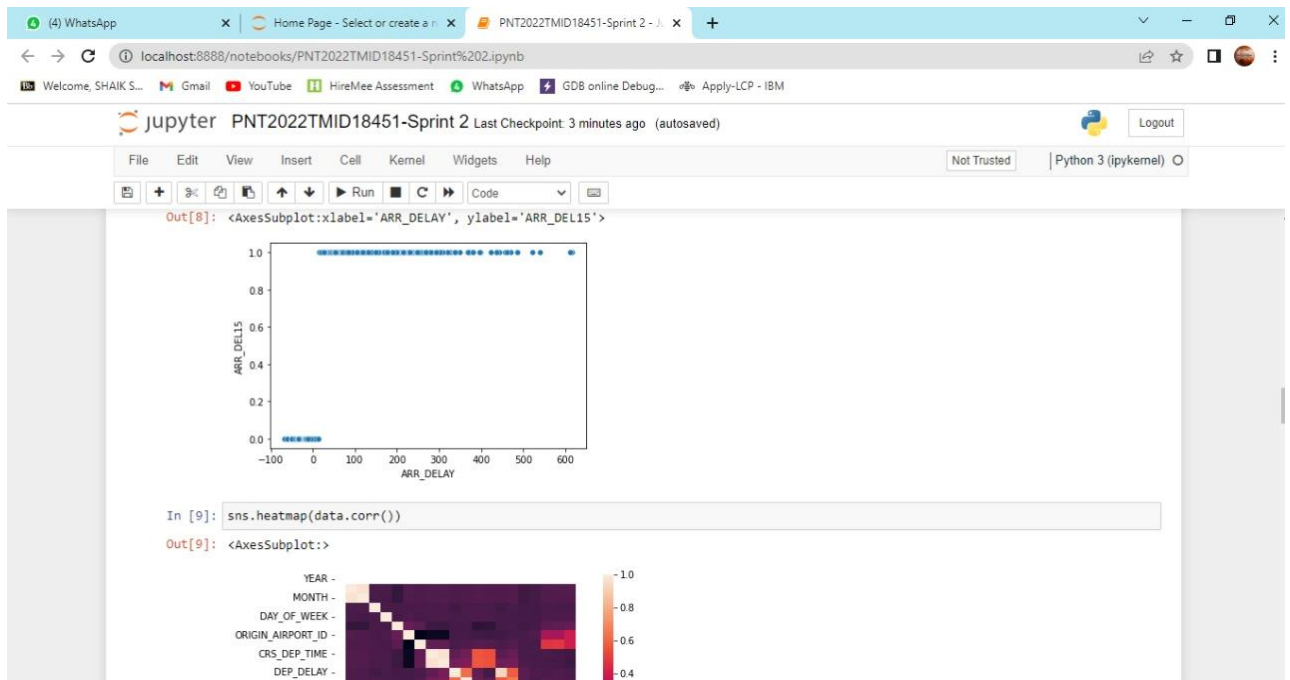
```
DEST_AIRPORT_ID 0
DEST             0
CRS_DEP_TIME     0
DEP_TIME        107
DEP_DELAY        107
DEP_DEL15        107
CRS_ARR_TIME     0
ARR_TIME        115
ARR_DELAY        188
ARR_DEL15        188
CANCELLED        0
DIVERTED         0
CRS_ELAPSED_TIME 0
ACTUAL_ELAPSED_TIME 188
DISTANCE         0
Unnamed: 25      11231
dtype: int64
```

```
In [7]: data['DEST'].unique()
```

```
Out[7]: array(['SEA', 'MSP', 'DTW', 'ATL', 'JFK'], dtype=object)
```

```
In [8]: sns.scatterplot(x='ARR_DELAY', y='ARR_DEL15', data=data)
```

```
Out[8]: <AxesSubplot:xlabel='ARR_DELAY', ylabel='ARR_DEL15'>
```



Browser tabs: (4) WhatsApp, Home Page - Select or create a n, PNT2022TMD18451-Sprint 2 - J. x

Address bar: localhost:8888/notebooks/PNT2022TMD18451-Sprint%202.ipynb

Navigation: Welcome, SHAIK S..., Gmail, YouTube, HireMee Assessment, WhatsApp, GDB online Debug..., Apply-LCP - IBM

Jupyter interface: PNT2022TMD18451-Sprint 2 Last Checkpoint: 3 minutes ago (autosaved) Python 3 (ipykernel)

File Edit View Insert Cell Kernel Widgets Help

Out[10]:

```
YEAR 0
QUARTER 0
MONTH 0
DAY_OF_MONTH 0
DAY_OF_WEEK 0
UNIQUE_CARRIER 0
TAIL_NUM 0
FL_NUM 0
ORIGIN_AIRPORT_ID 0
ORIGIN 0
DEST_AIRPORT_ID 0
DEST 0
CRS_DEP_TIME 0
DEP_TIME 107
DEP_DELAY 107
DEP_DEL15 107
CRS_ARR_TIME 0
ARR_TIME 115
ARR_DELAY 188
ARR_DEL15 188
CANCELLED 0
DIVERTED 0
CRS_ELAPSED_TIME 0
ACTUAL_ELAPSED_TIME 188
DISTANCE 0
dtype: int64
```

In [11]:

```
data=data[["FL_NUM","MONTH","DAY_OF_MONTH","DAY_OF_WEEK","ORIGIN","DEST","CRS_ARR_TIME","DEP_DEL15","ARR_DEL15"]]
data.isnull().sum()
```

Browser tabs: (4) WhatsApp, Home Page - Select or create a n, PNT2022TMD18451-Sprint 2 - J. x

Address bar: localhost:8888/notebooks/PNT2022TMD18451-Sprint%202.ipynb

Navigation: Welcome, SHAIK S..., Gmail, YouTube, HireMee Assessment, WhatsApp, GDB online Debug..., Apply-LCP - IBM

Jupyter interface: PNT2022TMD18451-Sprint 2 Last Checkpoint: 3 minutes ago (autosaved) Python 3 (ipykernel)

File Edit View Insert Cell Kernel Widgets Help

Out[11]:

```
FL_NUM 0
MONTH 0
DAY_OF_MONTH 0
DAY_OF_WEEK 0
ORIGIN 0
DEST 0
CRS_ARR_TIME 0
DEP_DEL15 107
ARR_DEL15 188
dtype: int64
```

In [12]:

```
data=data.fillna({'ARR_DEL15':1})
data=data.fillna({'DEP_DEL15':0})
data.iloc[177:185]
```

Out[12]:

	FL_NUM	MONTH	DAY_OF_MONTH	DAY_OF_WEEK	ORIGIN	DEST	CRS_ARR_TIME	DEP_DEL15	ARR_DEL15
177	2834	1	9	6	MSP	SEA	852	0.0	1.0
178	2839	1	9	6	DTW	JFK	1724	0.0	0.0
179	86	1	10	7	MSP	DTW	1632	0.0	1.0
180	87	1	10	7	DTW	MSP	1649	1.0	0.0
181	423	1	10	7	JFK	ATL	1600	0.0	0.0
182	440	1	10	7	JFK	ATL	849	0.0	0.0
183	485	1	10	7	JFK	SEA	1945	1.0	0.0
184	557	1	10	7	MSP	DTW	912	0.0	1.0

Localhost: 8888/notebooks/PNT2022TMID18451-Sprint2.ipynb

Python 3 (ipykernel)

In [20]: `t`

Out[20]: `array([[0., 0., 0., 0., 1.],
 [0., 0., 0., 1., 0.],
 [0., 0., 0., 0., 1.],
 ...,
 [0., 0., 0., 0., 1.],
 [0., 0., 0., 0., 1.],
 [0., 1., 0., 0., 0.]])`

In [21]: `x=np.delete(x,[4,5],axis=1)`
`x.shape`

Out[21]: `(11231, 6)`

In [22]: `x=np.concatenate((t,z,x),axis=1)`
`x.shape`

Out[22]: `(11231, 16)`

In [23]: `data=pd.get_dummies(data,columns=['ORIGIN','DEST'])`
`data.head()`

Out[23]:

	FL_NUM	MONTH	DAY_OF_MONTH	DAY_OF_WEEK	CRS_ARR_TIME	DEP_DEL15	ARR_DEL15	ORIGIN_0	ORIGIN_1	ORIGIN_2	ORIGIN_3	ORIGIN_4	DEST
0	1399	1	1	5	21	0.0	0.0	1	0	0	0	0	0
1	1476	1	1	5	14	0.0	0.0	0	1	0	0	0	0

Localhost: 8888/notebooks/PNT2022TMID18451-Sprint2.ipynb

Python 3 (ipykernel)

In [24]: `y=data.iloc[:,5:6].values`

In [25]: `from sklearn.model_selection import train_test_split`
`x_train,x_test,y_train,y_test = train_test_split(x,y,test_size=0.2,random_state=0)`
`x_test.shape`

Out[25]: `(2247, 16)`

In [26]: `x_test.shape`

Out[26]: `(2247, 16)`

In [27]: `x_train.shape`

Out[27]: `(8984, 16)`

In [28]: `y_test.shape`

Out[28]: `(2247, 1)`

In [29]: `y_train.shape`

Out[29]: `(8984, 1)`

