

# PROJECT DEVELOPMENT PHASE

## SPRINT 2 – CODE AND TESTCASE

<b>Date</b>	10 November 2022
<b>Team ID</b>	PNT2022TMID10242
<b>Project</b>	Flight delay prediction using Machine learning
<b>Marks</b>	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 interface with the following content:

### Sprint 1 - Code and Testcases

Date -- 08 November 2022 Team Id -- PNT2022TMID10242 Project -- Flight delay prediction using Machine Learning Marks -- 8 marks

```
In [1]: import sys
import numpy as np
import pandas as pd
import seaborn as sns
import pickle
import sklearn
%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 [2]: data=pd.read_csv("flightdata.csv")
data
```

Out[2]:

	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	...	2145
1	2016	1	1	1	5	DL	N964DN	1476	11433	DTW	...	1438
2	2016	1	1	1	5	DL	N813DN	1597	10397	ATL	...	1211
3	2016	1	1	1	5	DL	N587NW	1768	14747	SEA	...	1338
4	2016	1	1	1	5	DL	N836DN	1823	14747	SEA	...	601

Documents/IBM/IBM Project/ x PNT2022TMD10242 - sprint1 - x +

localhost:8888/notebooks/Documents/IBM/IBM%20Project/PNT2022TMD10242%20-%20sprint1.ipynb

jupyter PNT2022TMD10242 - sprint1 Last Checkpoint: 39 minutes ago (autosaved) Logout

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

Run

3	2016	1	1	1	5	DL	N587NW	1768	14747	SEA	...	133f
4	2016	1	1	1	5	DL	N836DN	1823	14747	SEA	...	60f
...	...	...	...	...	...	...	...	...	...	...	...	...
11226	2016	4	12	30	5	DL	N940DL	1715	11433	DTW	...	122f
11227	2016	4	12	30	5	DL	N836DN	1770	14747	SEA	...	204f
11228	2016	4	12	30	5	DL	N583NW	1823	11433	DTW	...	221f
11229	2016	4	12	30	5	DL	N554NW	1901	10397	ATL	...	180f
11230	2016	4	12	30	5	DL	N843DN	2005	10397	ATL	...	92f

11231 rows x 26 columns

In [3]: 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
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
```

Documents/IBM/IBM Project/ x PNT2022TMD10242 - sprint1 - x +

localhost:8888/notebooks/Documents/IBM/IBM%20Project/PNT2022TMD10242%20-%20sprint1.ipynb

jupyter PNT2022TMD10242 - sprint1 Last Checkpoint: 40 minutes ago (autosaved) Logout

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

Run

In [4]: data.describe()

Out[4]:

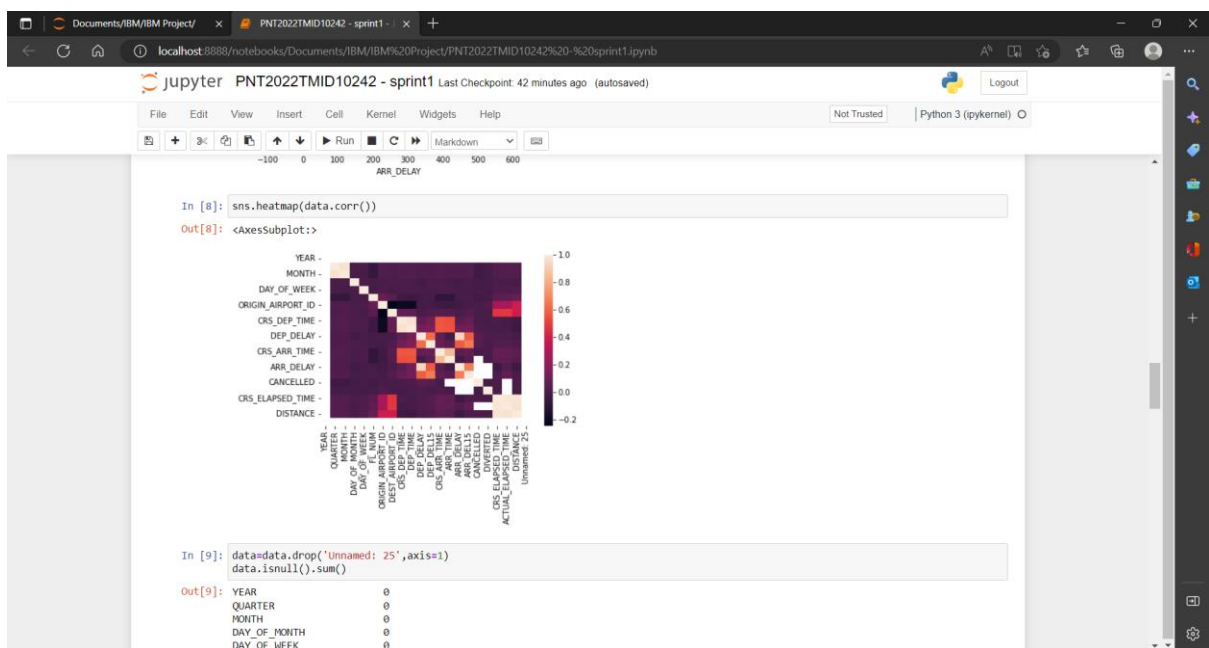
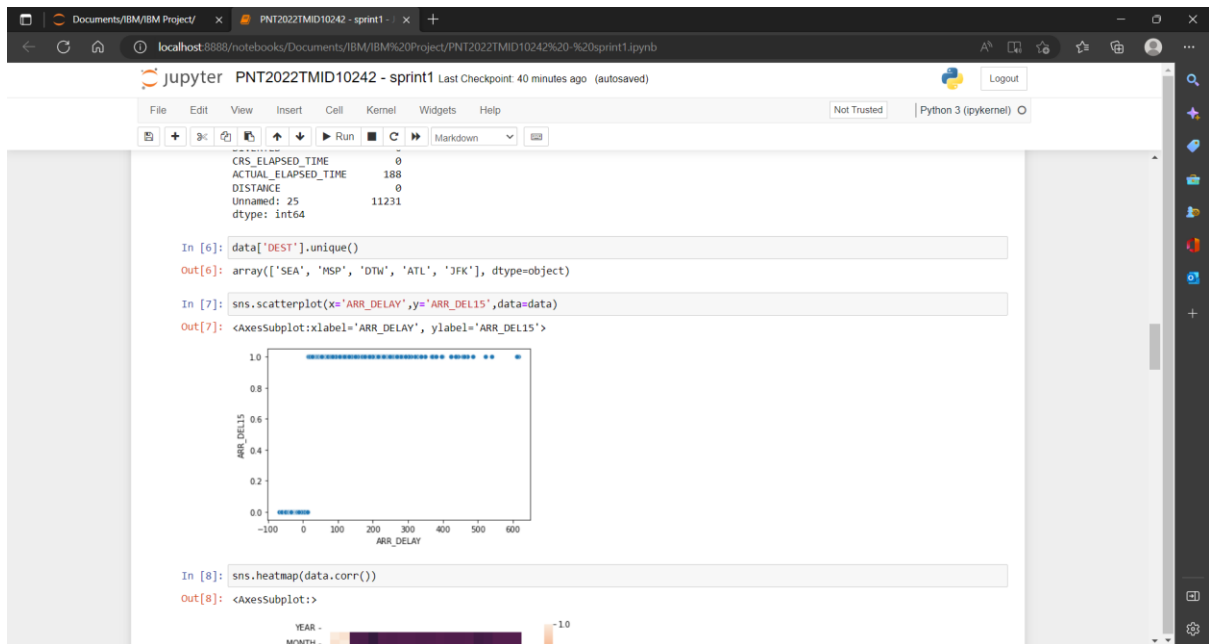
	YEAR	QUARTER	MONTH	DAY_OF_MONTH	DAY_OF_WEEK	FL_NUM	ORIGIN_AIRPORT_ID	DEST_AIRPORT_ID	CRS_DEP_TIME	DEP_
count	11231.0	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11231.000000	11124.00
mean	2016.0	2.544475	6.628973	15.790758	3.960199	1334.325617	12334.516695	12302.274508	1320.798326	1327.16
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 [5]: data.isnull().sum()

Out[5]:

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



```
Documents/IBM/IBM Project/ x PNT2022TMD10242 - sprint1 - x +
localhost:8888/notebooks/Documents/IBM/IBM%20Project/PNT2022TMD10242%20-%20sprint1.ipynb

jupyter PNT2022TMD10242 - sprint1 Last Checkpoint: 42 minutes ago (autosaved)
Logout

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

In [9]: data=data.drop('Unnamed: 25',axis=1)
data.isnull().sum()

Out[9]: 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
CANCELED 0
DIVERTED 0
CRS_ELAPSED_TIME 0
ACTUAL_ELAPSED_TIME 188
DISTANCE 0
dtype: int64
```

```
Documents/IBM/IBM Project/ x PNT2022TMD10242 - sprint1 - x +
localhost:8888/notebooks/Documents/IBM/IBM%20Project/PNT2022TMD10242%20-%20sprint1.ipynb

jupyter PNT2022TMD10242 - sprint1 Last Checkpoint: 42 minutes ago (autosaved)
Logout

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

ACTUAL_ELAPSED_TIME 188
DISTANCE 0
dtype: int64

In [10]: data=data[["FL_NUM","MONTH","DAY_OF_MONTH","DAY_OF_WEEK","ORIGIN","DEST","CRS_ARR_TIME","DEP_DEL15","ARR_DEL15"]]
data.isnull().sum()

Out[10]: 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 [11]: data=data.fillna({'ARR_DEL15':1})
data=data.fillna({'DEP_DEL15':0})
data.iloc[177:185]

Out[11]:
```

	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

Documents/IBM/IBM Project/ x PNT2022TMD10242 - sprint1 - x +

localhost:8888/notebooks/Documents/IBM/IBM%20Project/PNT2022TMD10242%20-%20sprint1.ipynb

jupyter PNT2022TMD10242 - sprint1 Last Checkpoint: 42 minutes ago (autosaved) Logout

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

```
In [12]: import math
for index,row in data.iterrows():
    data.loc[index,'CRS_ARR_TIME'] = math.floor(row['CRS_ARR_TIME'] / 100)
data.head()
```

Out[12]:

	FL_NUM	MONTH	DAY_OF_MONTH	DAY_OF_WEEK	ORIGIN	DEST	CRS_ARR_TIME	DEP_DEL15	ARR_DEL15
0	1399	1	1	5	ATL	SEA	21	0.0	0.0
1	1476	1	1	5	DTW	MSP	14	0.0	0.0
2	1597	1	1	5	ATL	SEA	12	0.0	0.0
3	1768	1	1	5	SEA	MSP	13	0.0	0.0
4	1823	1	1	5	SEA	DTW	6	0.0	0.0

```
In [13]: from sklearn.preprocessing import LabelEncoder
le=LabelEncoder()
data['DEST']=le.fit_transform(data['DEST'])
data['ORIGIN'] = le.fit_transform(data['ORIGIN'])
```

```
In [14]: data.head()
```

Out[14]:

	FL_NUM	MONTH	DAY_OF_MONTH	DAY_OF_WEEK	ORIGIN	DEST	CRS_ARR_TIME	DEP_DEL15	ARR_DEL15
0	1399	1	1	5	0	4	21	0.0	0.0
1	1476	1	1	5	1	3	14	0.0	0.0
2	1597	1	1	5	0	4	12	0.0	0.0
3	1768	1	1	5	4	3	13	0.0	0.0
4	1823	1	1	5	4	1	6	0.0	0.0

```
In [15]: x=data.iloc[:,0:8].values
```

Documents/IBM/IBM Project/ x PNT2022TMD10242 - sprint1 - x +

localhost:8888/notebooks/Documents/IBM/IBM%20Project/PNT2022TMD10242%20-%20sprint1.ipynb

jupyter PNT2022TMD10242 - sprint1 Last Checkpoint: 43 minutes ago (autosaved) Logout

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

```
4 1823 1 1 5 4 1 6 0.0 0.0
```

```
In [15]: x=data.iloc[:,0:8].values
y=data.iloc[:,8:9].values
x.shape
```

Out[15]: (11231, 8)

```
In [16]: y
```

Out[16]: array([[0.],
[0.],
[0.],
...,
[0.],
[0.],
[0.]])

```
In [17]: from sklearn.preprocessing import OneHotEncoder
oh=OneHotEncoder()
z=oh.fit_transform(data.iloc[:,4:5]).toarray()
t=oh.fit_transform(data.iloc[:,5:6]).toarray()
```

```
In [18]: z
```

Out[18]: array([[1., 0., 0., 0., 0.],
[0., 1., 0., 0., 0.],
[1., 0., 0., 0., 0.],
...,
[0., 1., 0., 0., 0.],
[1., 0., 0., 0., 0.],
[1., 0., 0., 0., 0.]])

```
In [19]: t
```

Documents/IBM/IBM Project/ x PNT2022TMD10242 - sprint1 - x +

localhost:8888/notebooks/Documents/IBM/IBM%20Project/PNT2022TMD10242%20-%20sprint1.ipynb

jupyter PNT2022TMD10242 - sprint1 Last Checkpoint: 43 minutes ago (autosaved) Logout

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

[1., 0., 0., 0., 0.]]

```
In [19]: t
Out[19]: 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 [20]: x=np.delete(x,[4,5],axis=1)
x.shape
Out[20]: (11231, 6)

In [21]: x=np.concatenate((t,z,x),axis=1)
x.shape
Out[21]: (11231, 16)

In [22]: data=pd.get_dummies(data,columns=['ORIGIN','DEST'])
data.head()
Out[22]:
```

	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	DES
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
2	1597	1	1	5	12	0.0	0.0	1	0	0	0	0	0
3	1768	1	1	5	13	0.0	0.0	0	0	0	0	0	1
4	1823	1	1	5	6	0.0	0.0	0	0	0	0	0	1

Documents/IBM/IBM Project/ x PNT2022TMD10242 - sprint1 - x +

localhost:8888/notebooks/Documents/IBM/IBM%20Project/PNT2022TMD10242%20-%20sprint1.ipynb

jupyter PNT2022TMD10242 - sprint1 Last Checkpoint: 43 minutes ago (autosaved) Logout

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

```
2 1597 1 1 5 12 0.0 0.0 1 0 0 0 0 0
3 1768 1 1 5 13 0.0 0.0 0 0 0 0 0 1
4 1823 1 1 5 6 0.0 0.0 0 0 0 0 0 1
```

```
In [23]: y=data.iloc[:,5:6].values

In [24]: 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)

In [25]: x_test.shape
Out[25]: (2247, 16)

In [26]: x_train.shape
Out[26]: (8984, 16)

In [27]: y_test.shape
Out[27]: (2247, 1)

In [28]: y_train.shape
Out[28]: (8984, 1)
```