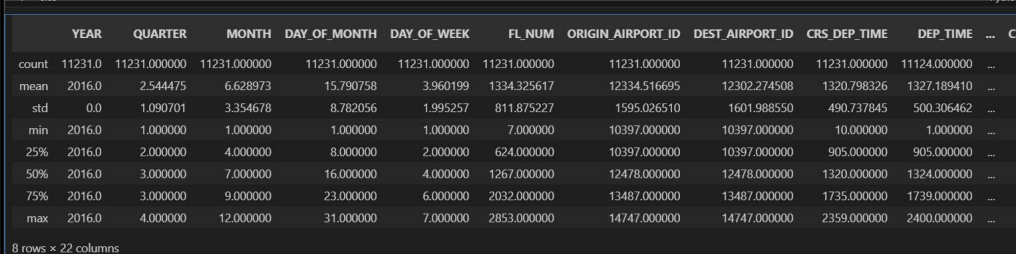


## Data Collection and Preprocessing Phase

Date	18 June 2024
Team ID	739634
Project Title	Flight Delays Prediction Using Machine Learning
Maximum Marks	6 Marks

### Data Exploration and Preprocessing Report

Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong foundation for insights and predictions.

Section	Description
Data Overview	<p><u>Dimension:</u> 11231 rows × 26 Columns</p> <p><u>Descriptive statistics:</u></p> <p>dataset.describe()</p>  <p>8 rows × 22 columns</p>

dataset.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
14  DEP_DELAY                         11124 non-null  float64
15  DEP_DEL15                         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_DEL15                         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
```

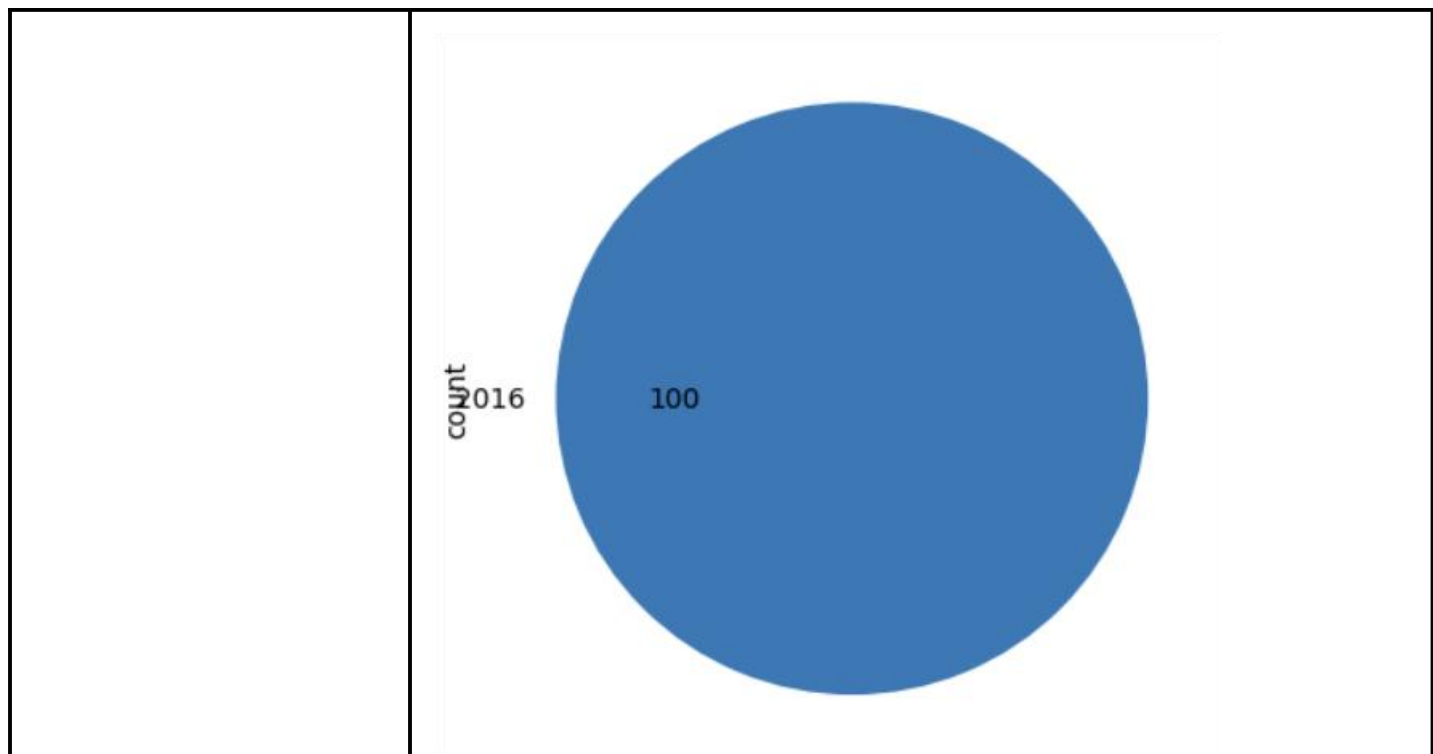
```
dataset.isnull().sum()
```

```
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
Unnamed: 25     11231
dtype: int64
```

```
dataset['DEST'].unique()
```

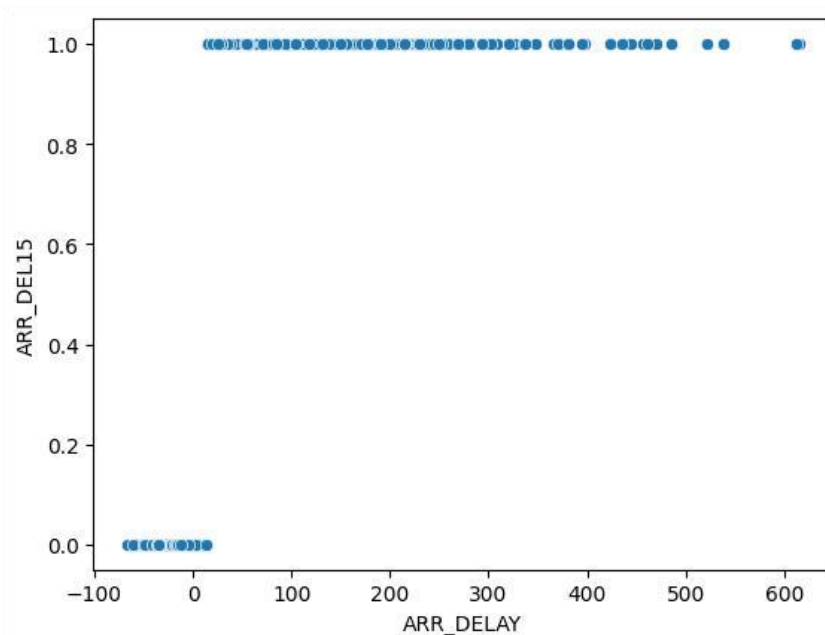
```
array(['SEA', 'MSP', 'DTW', 'ATL', 'JFK'], dtype=object)
```

Univariate Analysis

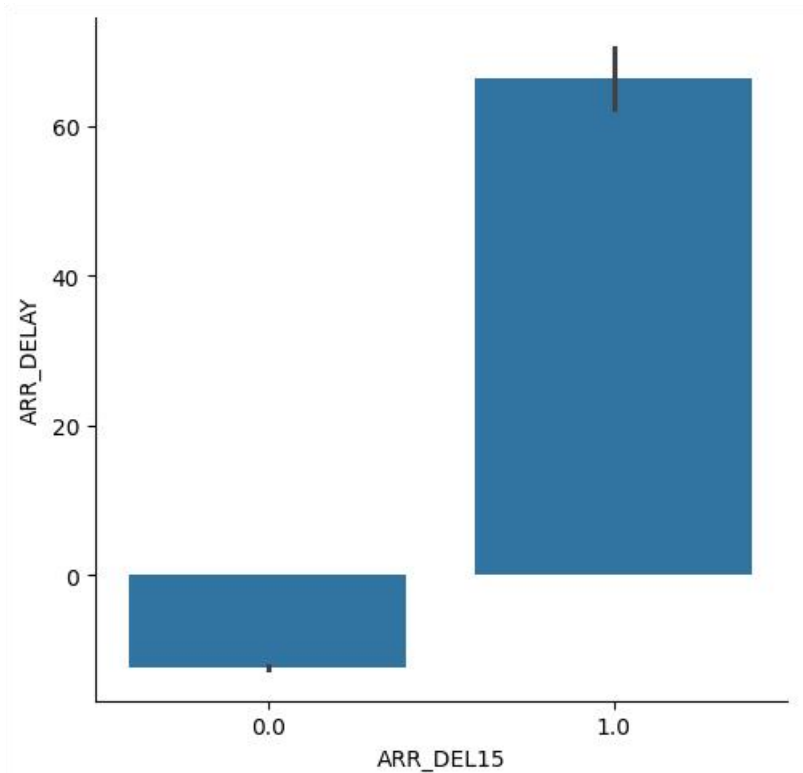


## Bivariate Analysis

### Scatterplot:



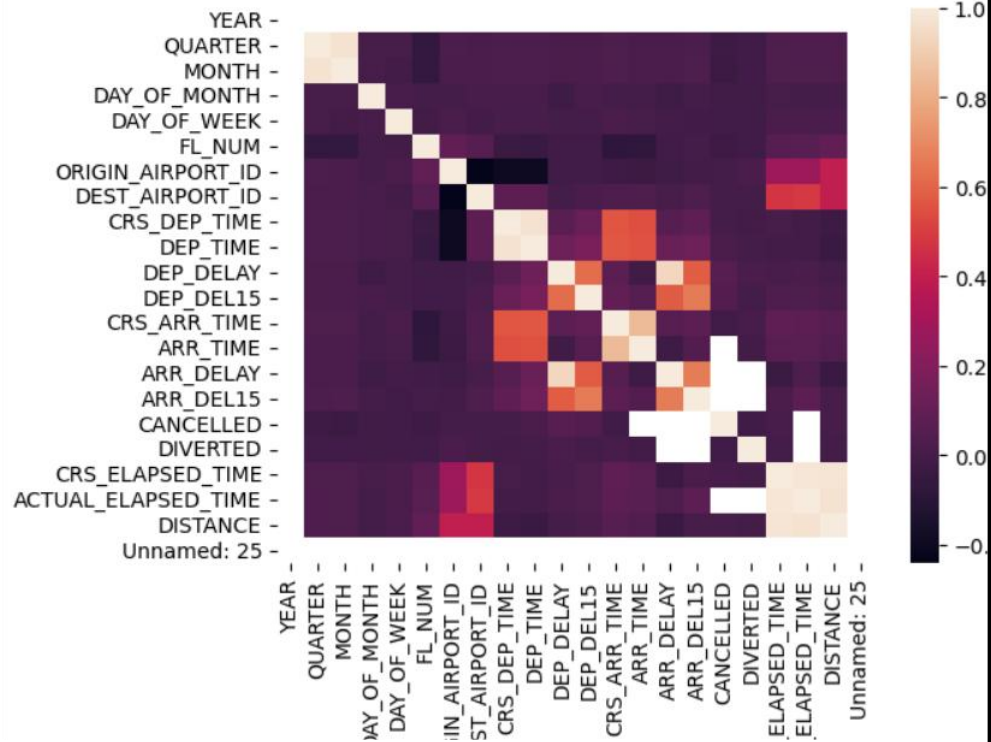
### CATPLOT:



### Heatmap:

Multivariate

## HEAT MAP:



Outliers and Anomalies

-

## Data Preprocessing Code Screenshots

Loading Data

```
dataset=pd.read_csv("flightdata.csv")
```

Pyth

```
dataset.head()
```

Pyth

YEAR	QUARTER	MONTH	DAY OF MONTH	DAY OF WEEK	UNIQUE CARRIER	TAIL NUM	FL NUM	ORIGIN AIRPORT ID	ORIGIN	...	CRS ARR TIME	ARR TIME
2016	1	1	1	5	DL	N836DN	1399	10397	ATL	...	2143	2102.0
2016	1	1	1	5	DL	N964DN	1476	11433	DTW	...	1435	1439.0
2016	1	1	1	5	DL	N813DN	1597	10397	ATL	...	1215	1142.0
2016	1	1	1	5	DL	N587NW	1768	14747	SEA	...	1335	1345.0
2016	1	1	1	5	DL	N836DN	1823	14747	SEA	...	607	615.0

Handling Missing Data

```
import pandas as pd
dataset=dataset.drop("Unnamed: 25",axis=1)
dataset.isnull().sum()
```

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

```
dataset=dataset.fillna({'ARR_DEL15':1})
dataset=dataset.fillna({'dep_del15':0})
dataset.iloc[177:185]
```

	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	NaN	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

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

	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

	<pre>from sklearn.preprocessing import LabelEncoder le=LabelEncoder() dataset['ORIGIN']=le.fit_transform(dataset['ORIGIN']) dataset['DEST']=le.fit_transform(dataset['DEST']) dataset.head()</pre> <table><thead><tr><th></th><th>FL_NUM</th><th>MONTH</th><th>DAY_OF_MONTH</th><th>DAY_OF_WEEK</th><th>ORIGIN</th><th>DEST</th><th>CRS_ARR_TIME</th><th>DEP_DEL15</th><th>ARR_DEL15</th></tr></thead><tbody><tr><td>0</td><td>1399</td><td>1</td><td>1</td><td>5</td><td>0</td><td>4</td><td>21</td><td>0.0</td><td>0.0</td></tr><tr><td>1</td><td>1476</td><td>1</td><td>1</td><td>5</td><td>1</td><td>3</td><td>14</td><td>0.0</td><td>0.0</td></tr><tr><td>2</td><td>1597</td><td>1</td><td>1</td><td>5</td><td>0</td><td>4</td><td>12</td><td>0.0</td><td>0.0</td></tr><tr><td>3</td><td>1768</td><td>1</td><td>1</td><td>5</td><td>4</td><td>3</td><td>13</td><td>0.0</td><td>0.0</td></tr><tr><td>4</td><td>1823</td><td>1</td><td>1</td><td>5</td><td>4</td><td>1</td><td>6</td><td>0.0</td><td>0.0</td></tr></tbody></table>		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
	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																																																				
Data Transformation	<pre>from sklearn.preprocessing import OneHotEncoder oh=OneHotEncoder() z=oh.fit_transform(dataset.iloc[:,4:5]).toarray() t=oh.fit_transform(dataset.iloc[:,5:6]).toarray()</pre> <pre>dataset = dataset.dropna()</pre> <p>✓ 0.0s</p> <pre>x=dataset.iloc[:,0:8].values y=dataset.iloc[:,8:9].values</pre> <pre>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)</pre>																																																												
Feature Engineering	Attached the codes in final submission.																																																												
Save Processed Data	-																																																												