# Data Collection and Preprocessing Phase

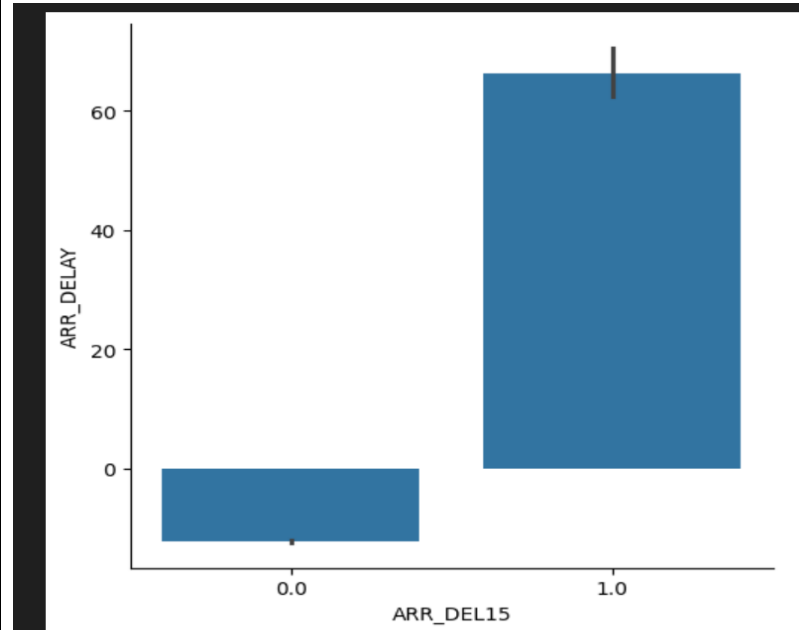| | |
|---|---|
| Date | 15 July 2024 |
| Team ID | 740080 |
| Project Title | Flight Delay Prediction using Machine Learning. |
| Maximum Marks | 6 Marks |

## Data Exploration and Preprocessing Template

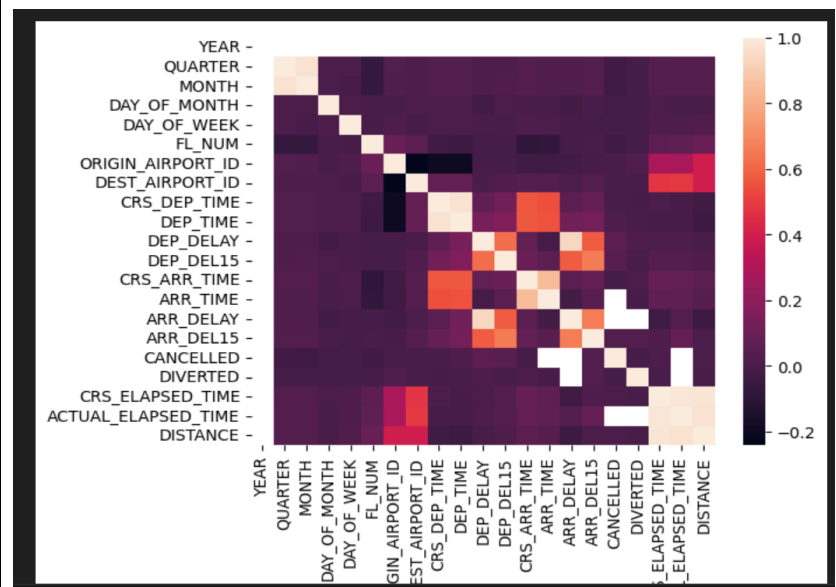Dataset variables will be statistically analyzed to identify patterns and outliers, with Python employed for preprocessing tasks like normalization and feature engineering. Data cleaning will address missing values and outliers, ensuring quality for subsequent analysis and modeling, and forming a strong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

| Section | Description |
|---|---|
| Data Overview |  |
| Univariate Analysis |  |

| | |
|---|---|
| Bivariate Analysis |  |
| Multivariate Analysis |  |
| Outliers and Anomalies | |

**Data Preprocessing Code Screenshots**

| | Loading Data |
|---|---|

```
flights=pd.read_csv("flightdata.csv")
flights
✓ 0.1s
```

| | YEAR | QUARTER | MONTH | DAY_OF_MONTH | DAY_OF_WEEK | UNIQUE_CARRIER | TAIL_NUM | FL_NUM | ORIGIN_AIRPORT_ID | ORIGIN | ... | DEP_DEL15 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 2016 | 1 | 1 | 1 | 5 | DL | N836DN | 1399 | 10397 | ATL | ... | 0.0 |
| 1 | 2016 | 1 | 1 | 1 | 5 | DL | N964DN | 1476 | 11433 | DTW | ... | 0.0 |
| 2 | 2016 | 1 | 1 | 1 | 5 | DL | N813DN | 1597 | 10397 | ATL | ... | 0.0 |
| 3 | 2016 | 1 | 1 | 1 | 5 | DL | N587NW | 1768 | 14747 | SEA | ... | 0.0 |
| 4 | 2016 | 1 | 1 | 1 | 5 | DL | N836DN | 1823 | 14747 | SEA | ... | 0.0 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 11226 | 2016 | 4 | 12 | 30 | 5 | DL | N940DL | 1715 | 11433 | DTW | ... | 0.0 |
| 11227 | 2016 | 4 | 12 | 30 | 5 | DL | N836DN | 1770 | 14747 | SEA | ... | 1.0 |
| 11228 | 2016 | 4 | 12 | 30 | 5 | DL | N583NW | 1823 | 11433 | DTW | ... | 0.0 |
| 11229 | 2016 | 4 | 12 | 30 | 5 | DL | N554NW | 1901 | 10397 | ATL | ... | 0.0 |
| 11230 | 2016 | 4 | 12 | 30 | 5 | DL | N842DN | 2005 | 10397 | ATL | ... | 0.0 |

| S_ARR_TIME | ARR_TIME | ARR_DELAY | ARR_DEL15 | CANCELLED | DIVERTED | CRS_ELAPSED_TIME | ACTUAL_ELAPSED_TIME | DISTANCE |
|---|---|---|---|---|---|---|---|---|
| 2143 | 2102.0 | -41.0 | 0.0 | 0 | 0 | 338 | 295.0 | 2182 |
| 1435 | 1439.0 | 4.0 | 0.0 | 0 | 0 | 110 | 115.0 | 528 |
| 1215 | 1142.0 | -33.0 | 0.0 | 0 | 0 | 335 | 300.0 | 2182 |
| 1335 | 1345.0 | 10.0 | 0.0 | 0 | 0 | 196 | 205.0 | 1399 |
| 607 | 615.0 | 8.0 | 0.0 | 0 | 0 | 247 | 259.0 | 1927 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 1223 | 1148.0 | -35.0 | 0.0 | 0 | 0 | 138 | 105.0 | 594 |
| 2046 | 2100.0 | 14.0 | 0.0 | 0 | 0 | 201 | 181.0 | 1399 |
| 2210 | 2154.0 | -16.0 | 0.0 | 0 | 0 | 311 | 295.0 | 1927 |
| 1806 | 1801.0 | -5.0 | 0.0 | 0 | 0 | 336 | 332.0 | 2182 |
| 925 | 913.0 | -12.0 | 0.0 | 0 | 0 | 120 | 110.0 | 594 |

| | Handling Missing Data |
|---|---|

```
flights=flights.fillna({'ARR_DEL15':1})
flights=flights.fillna({'dep_del15':0})
flights.iloc[177:185]
✓ 0.0s
```

| | 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 flights.iterrows():
    flights.loc[index,'CRS_ARR_TIME']=math.floor(row['CRS_ARR_TIME']/100)
flights.head()
✓ 3.2s
```

| | 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 |

| Data Transformation |  |
|---|---|
| Feature Engineering | Attached the codes in the final submission. |
| Save Processed Data | _ |