

## **Final Term Project**

Jyothi Sevakula

Virginia Polytechnic Institute and State University

CS 5805 Machine Learning 1

Professor Dr. Reza Jafari

December 8th, 2024

## Table of Contents

<b>Table of Contents.....</b>	<b>2</b>
<b>Abstract.....</b>	<b>7</b>
<b>Introduction.....</b>	<b>8</b>
<b>Description of Dataset.....</b>	<b>9</b>
<b>Phase I: Feature Engineering and EDA.....</b>	<b>10</b>
Data Preprocessing:.....	10
Data Cleaning:.....	10
Data Duplicates:.....	10
Down Sampling:.....	10
Dimensionality Reduction:.....	12
Random Forest Analysis:.....	12
Principal Component Analysis:.....	13
Singular Value Decomposition:.....	14
Variance Inflation Factor:.....	16
Collinearity Analysis:.....	17
Observations:.....	17
Discretization and Binarization:.....	17
Variable Transformation:.....	18
Outlier Analysis and Removal:.....	18
Covariance Matrix:.....	20
Pearson Correlation Matrix:.....	21
Balanced or Imbalanced Data:.....	21
<b>Phase II: Regression Analysis.....</b>	<b>22</b>
T-test Analysis:.....	22
F-test Analysis:.....	23
Linear Regression Model and Prediction of Dependent Variable:.....	25
Confidence Interval Analysis:.....	26
Stepwise regression and adjusted R-square analysis:.....	27
Comparing Linear regression and Backward stepwise regression:.....	37
<b>Phase III: Classification Analysis.....</b>	<b>37</b>
Decision Tree:.....	37
Pre-pruning:.....	37
Post-pruning:.....	39
Logistic Regression:.....	41
K-nearest Neighbors:.....	42
SVM:.....	45
Naive Bayes:.....	47

Neural Network:.....	49
Multi-layered Perceptron:.....	49
Analysis of all classifiers of Phase III:.....	51
<b>Phase IV: Clustering and Association.....</b>	<b>53</b>
K-mean Algorithm:.....	53
DBSCAN Algorithm:.....	55
Apriori Algorithm:.....	56
<b>Recommendations.....</b>	<b>57</b>
<b>Appendix.....</b>	<b>58</b>
<b>References.....</b>	<b>70</b>

## Table of figures

**Page 10 Figure 1** Missing values in the dataset

**Page 10 Figure 2** Duplicate values in the dataset

**Page 11 Figure 3** Unique values in each feature

**Page 11 Figure 4** Shape of the data after downsampling

**Page 12 Figure 5** Plot of Decreasing order of feature importances from random forest

**Page 14 Figure 6** Cumulative explained variance plot of PCA

**Page 15 Figure 7** Cumulative explained variance ratio plot from SVD

**Page 15 Figure 8** Singular values of the dataset

**Page 16 Figure 9** VIF values of each feature

**Page 17 Figure 10** One hot-encoded dataset

**Page 19 Figure 11** Box plot of Length feature after removal of outliers using IQR method

**Page 20 Figure 12** The covariance matrix of all features

**Page 21 Figure 13** Pearson correlation coefficient matrix

**Page 21 Figure 14** Count of each class of target ‘Delay’

**Page 23 Figure 15** T-statistic and p-values of features

**Page 24 Figure 16** F-test and p-values of each feature

**Page 25 Figure 17** Coefficients and R<sup>2</sup>, Adj. R<sup>2</sup>, AIC, BIC, MSE of the Linear Regression Model

**Page 26 Figure 18** Predicted and Test values of ‘Length’ using Linear Regression

**Page 27 Figure 19** Confidence Interval analysis

**Page 28 Figure 20** OLS summary after eliminating the ‘DayOfWeek\_3’ feature

**Page 29 Figure 21** OLS summary after eliminating the ‘DayOfWeek\_4’ feature

**Page 30 Figure 22** OLS summary after eliminating the ‘DayOfWeek\_5’ feature

**Page 31 Figure 23** OLS summary after eliminating the 'DayOfWeek\_2' feature

**Page 32 Figure 24** OLS summary after eliminating the 'DayOfWeek\_7' feature

**Page 33 Figure 25** OLS summary after eliminating the 'Flight' feature

**Page 34 Figure 26** OLS summary after eliminating 'DayOfWeek\_6' feature

**Page 35 Figure 27** OLS summary after Backward regression

**Page 36 Figure 28** Condition number and Singular values after Backward stepwise regression

**Page 36 Figure 29** Plot of train, test, and predicted values of Backward stepwise regression

**Page 37 Figure 30** Comparison of Linear regression and Backward stepwise regression

**Page 37 Figure 31** Best hyperparameters of Pre-pruned decision tree

**Page 38 Figure 32** Various metrics of Pre-pruned decision tree

**Page 38 Figure 33** Confusion matrix of Pre-pruned decision tree

**Page 38 Figure 34** ROC curve of Pre-pruned decision tree

**Page 39 Figure 35** Optimal alpha value of Post-pruned decision tree

**Page 39 Figure 36** Accuracy vs. Alpha plot

**Page 40 Figure 37** Various metrics of Post-pruned decision tree

**Page 40 Figure 38** ROC curve of Post-pruned decision tree

**Page 40 Figure 39** Confusion matrix of Post-pruned decision tree

**Page 41 Figure 40** Best hyperparameters of Logistic Regression

**Page 41 Figure 41** Various metrics of Logistic Regression

**Page 41 Figure 42** ROC curve of Logistic Regression

**Page 42 Figure 43** Confusion matrix of Logistic Regression

**Page 43 Figure 44** Elbow method for finding optimal k

**Page 43 Figure 45** Best hyperparameters of K-Nearest Neighbors

**Page 43 Figure 46** Various metrics of K-Nearest Neighbors

**Page 44 Figure 47** ROC curve of K-Nearest Neighbors

**Page 44 Figure 48** Confusion matrix of K-Nearest Neighbors

**Page 45 Figure 49** Best hyperparameters of SVM

**Page 45 Figure 50** Various metrics of SVM

**Page 46 Figure 51** Confusion matrix of SVM

**Page 46 Figure 52** ROC curve of SVM

**Page 47 Figure 53** Various metrics of Gaussian Naive Bayes

**Page 47 Figure 54** Confusion matrix of Gaussian Naive Bayes

**Page 48 Figure 55** ROC curve of Gaussian Naive Bayes

**Page 48 Figure 56** Best hyperparameters of Multi-layered perceptron

**Page 49 Figure 57** Various metrics of Multi-layered perceptron

**Page 49 Figure 58** Confusion matrix of Multi-layered perceptron

**Page 49 Figure 59** ROC curve of Multi-layered perceptron

**Page 50 Figure 60** Comparison of all metrics of all classifiers

**Page 50 Figure 61** ROC curve of all classifiers

**Page 52 Figure 62** Silhouette analysis for Optimal k using KMeans clustering

**Page 53 Figure 63** WCSS plot analysis for optimal k using KMeans clustering

**Page 54 Figure 64** K-distance graph for optimal eps for dbscan algorithm

**Page 55 Figure 65** Rules generated from Apriori

## Abstract

Flight delays pose significant challenges in the airline industry, impacting both operational efficiency and passenger satisfaction. Delays can lead to cascading operational disruptions, increased costs for airlines, and frustration among passengers, highlighting the critical need for accurate and timely predictions. This project addresses this pressing issue by focusing on predicting flight delays using a comprehensive airlines dataset and a wide array of machine learning techniques. The ability to anticipate delays allows airlines to proactively manage schedules, allocate resources more effectively, and improve the overall passenger experience by providing timely updates and solutions.

The analysis employs regression, classification, clustering, and association rule mining to uncover patterns and insights that drive delay predictions. Through rigorous feature engineering, preprocessing, and dimensionality reduction techniques such as Random Forest, PCA, and SVD, the dataset is optimized for effective modeling. By identifying the most significant factors contributing to delays, this study provides actionable insights that can assist airlines in mitigating delays and enhancing operational reliability.

This project underscores the broader importance of leveraging predictive analytics within the aviation industry, demonstrating its potential to streamline operations, reduce costs, and boost customer satisfaction. The findings serve as a foundation for data-driven decision-making, offering a roadmap for implementing advanced analytics to address complex challenges in aviation.

## Introduction

The project focuses on leveraging diverse machine learning methodologies to predict flight delays using a reliable airlines dataset based on scheduled flight details. The process begins with comprehensive preprocessing and feature engineering to ensure the dataset is prepared for model development.

Various dimensionality reduction techniques, including PCA, Random Forest, SVD, and VIF, are applied to streamline the feature space. The statistical significance of features is assessed through T-test and F-test analyses. Linear regression and backward stepwise regression are utilized to predict the numerical feature ‘Length’ and refine the model.

For classification, a range of models—such as decision trees, logistic regression, KNN, SVM, Naive Bayes, and neural networks—are evaluated. Grid search and stratified k-fold cross-validation are performed to identify the best hyperparameters, with model selection guided by a detailed analysis of evaluation metrics. Additionally, the project employs KMeans clustering, an unsupervised learning approach, to identify optimal clusters. Techniques like silhouette analysis and within-cluster sum of squares (WCSS) are used to determine the most appropriate number of clusters.

The project concludes with the application of the Apriori algorithm to uncover meaningful association rules within the dataset. This in-depth analysis provides valuable insights into the intricate relationships between features, enhancing the understanding of factors influencing flight delays.

## Description of Dataset

The Airlines dataset comprises 539,383 instances and 9 distinct features, aimed at predicting whether a flight will be delayed based on its scheduled departure details. The dataset includes the following features: id, Airline, Flight, AirportFrom, AirportTo, DayOfWeek, Time, Length, and Delay.

Each record is uniquely identified by the id attribute and contains comprehensive flight information. The Airline attribute specifies the carrier operating the flight, while Flight represents the type of aircraft used. The attributes AirportFrom and AirportTo indicate the departure and arrival airports, respectively.

The DayOfWeek feature expressed numerically from 1 to 7, denotes the day the flight is scheduled. Time records the departure time in minutes past midnight, providing a precise temporal reference. The Length attribute reflects the flight duration in minutes, offering insights into its time span. Finally, the Delay attribute indicates whether the flight experienced a delay, serving as a key factor impacting airline operations and passenger satisfaction.

The dataset includes both categorical features (Airline, Flight, AirportFrom, AirportTo, DayOfWeek, Delay) and numerical features (Time, Length), meeting the project criteria of having more than 50,000 observations, at least two numerical features, and at least two categorical features with multiple categories.

## Phase I: Feature Engineering and EDA

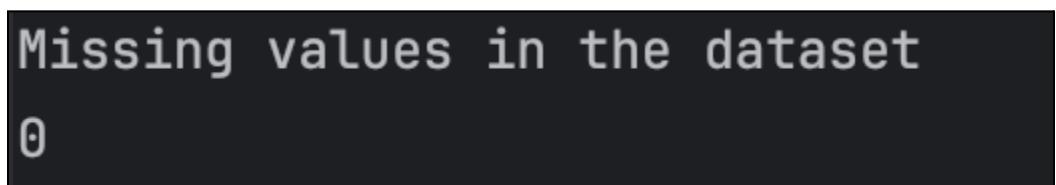
### Data Preprocessing:

#### *Data Cleaning:*

There are no missing values in the dataset. There are also no nan values. The dataset is inherently clean and the information of the scheduled flights is consistent.

**Figure 1**

*Missing values in the dataset*

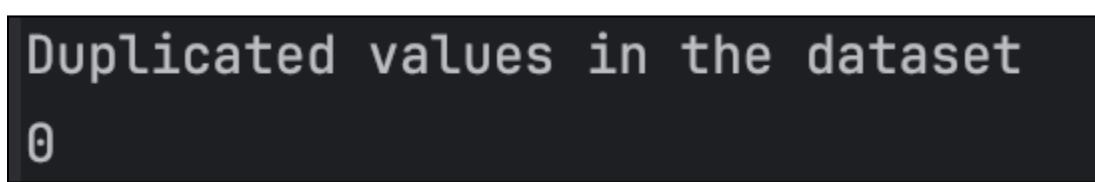


#### *Data Duplicates:*

There are no data duplicates in the dataset.

**Figure 2**

*Duplicate values in the dataset*



#### *Down Sampling:*

There are 18 categories in Airline, 293 categories in AirportFrom and 293 categories in AirportTo, 7 categories in DayOfWeek. Hence, I preferred to predict delays related to top 5 source airports which are from the AirportFrom feature and hence dropped the AirportTo feature. Also, considered the top 5 airlines as the remaining airlines have significantly low observations.

**Figure 3**

*Unique values in each feature*

Number of unique values in each feature of whole dataset	
id	539383
Airline	18
Flight	6585
AirportFrom	293
AirportTo	293
DayOfWeek	7
Time	1131
Length	426

Now, the number of observations is 37,570 with 7 features. Also, dropped id columns as it has no significance in prediction.

**Figure 4**

*Shape of the data after downsampling*

After down sampling, the shape of the dataset  
(37570, 7)

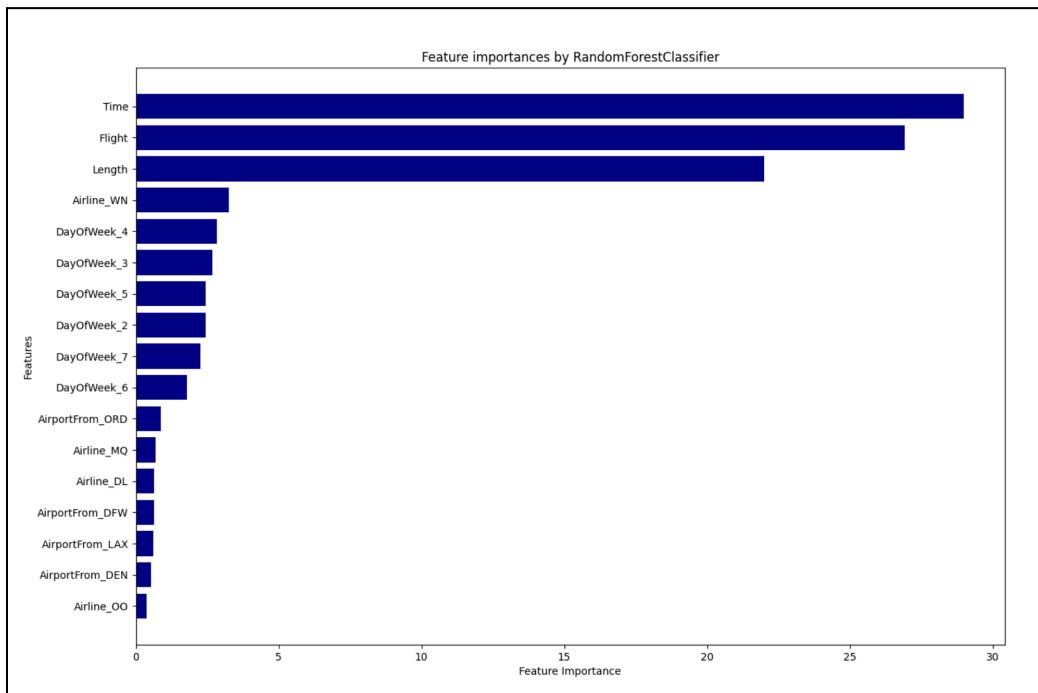
### ***Dimensionality Reduction or Feature Selection:***

#### ***Random Forest analysis:***

Random Forest analysis plays a critical role in dimensionality reduction by identifying the most influential features that impact the model's predictions. This technique evaluates the relevance and contribution of each feature to the target variable using the Random Forest algorithm's built-in “feature\_importances\_” attribute. By quantifying the importance of each feature, this analysis highlights the key variables that significantly affect the model's decision-making process. Prioritizing these essential features helps streamline the feature space, resulting in a more efficient model while maintaining predictive accuracy. Retaining only the most impactful features is crucial for enhancing model performance and reducing computational complexity.

**Figure 5**

*Plot of Decreasing order of feature importances from random forest*



From above figure 5, we can see the importance of features in decreasing order. We can drop features by selecting a cumulative feature importance threshold of 95% and hence dropped features Airline\_MQ, AirportFrom\_DFW, AirportFrom\_LAX, Airline\_DL, AirportFrom\_DEN, Airline\_OO. Hence, the total dimensions are reduced to 11 from 17.

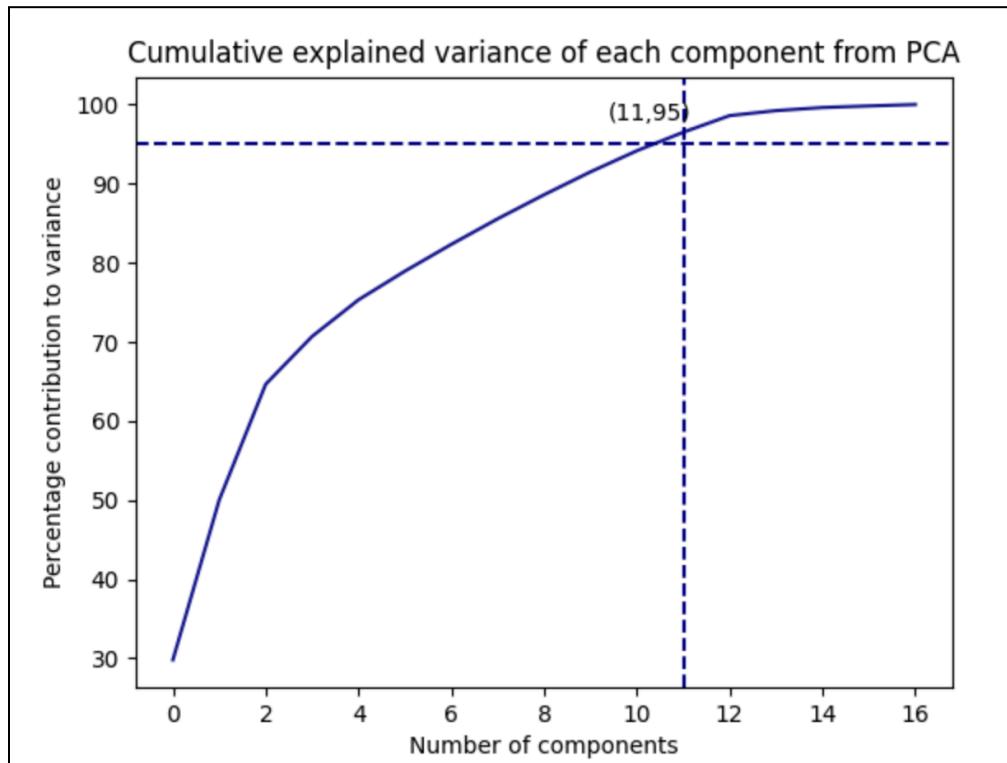
### ***Principal Component Analysis:***

Principal Component Analysis (PCA) is a widely used unsupervised learning technique for reducing the dimensionality of large datasets. It enhances interpretability while minimizing information loss, making it a fundamental tool in dimensionality reduction for data analysis. PCA works by transforming high-dimensional data into a lower-dimensional representation, called principal components. These components retain the most significant information from the original dataset, helping to condense the feature space while preserving key patterns and minimizing redundancy. Unlike methods that identify specific important features, PCA does not indicate which features are important but determines how many are essential to explain the majority of the variance. This dimensionality reduction facilitates easier visualization, and computational efficiency, and often improves the performance of machine learning models.

```
Cumulative explained variance of each component from PCA
[ 29.81  50.05  64.65  70.69  75.32  78.94  82.33  85.56  88.6   91.49
  94.16  96.5   98.6   99.23  99.62  99.82 100. ]
```

**Figure 6**

*Cumulative explained variance plot of PCA*



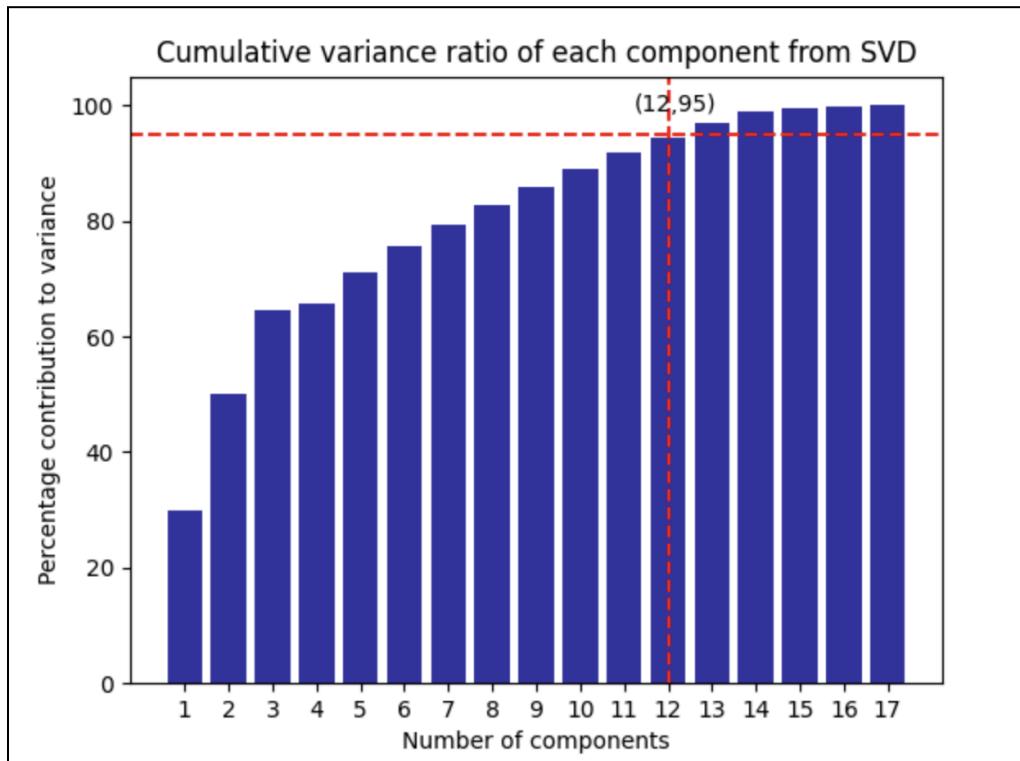
When we choose a threshold of 95% of the total variance, we see that we are left with 11 components from PCA and these features are enough to explain 95% of the total variance of the data.

#### ***Singular Value Decomposition Analysis:***

Singular Value Decomposition (SVD) is a robust dimensionality reduction technique in data analysis. It breaks down a matrix into three component matrices, uncovering patterns and structures within the data. By isolating the most significant components, SVD transforms high-dimensional data into a reduced set of singular vectors. This transformation creates a more compact representation of the original data while retaining critical information. SVD is instrumental in tasks like noise reduction, feature extraction, and improving computational efficiency in machine learning workflows.

**Figure 7**

*Cumulative explained variance ratio plot from SVD*

**Figure 8**

*Singular values of the dataset*

```
Singular values of features [229.87 189.3 160.73 129.39 100.8 90.46 80.09 77.45 75.46 73.35
71.41 68.69 64.29 60.22 33.2 24.01 18.82]
```

Using a 95% threshold for total explained variance, the dimensions can be reduced from 17 to 12 without significant loss of information. This reduction simplifies the dataset, making it more manageable for further analysis..

### ***Variance Inflation Factor:***

Variance Inflation Factor (VIF) is a tool used to identify multicollinearity among predictors in regression analysis. It measures how much the variance of a regression coefficient is increased due to multicollinearity among the predictors.

VIF values that are particularly high (typically above 5 or 10) often indicate strong multicollinearity, suggesting that certain predictors are either redundant or highly correlated. Removing these highly correlated features can help reduce multicollinearity and may serve as a form of dimensionality reduction in regression modeling.

Based on the VIF values observed, the variable ‘Airline\_MQ’ shows significant multicollinearity with other predictors, which could affect the stability and reliability of the coefficient estimates in the regression model.

**Figure 9**

*VIF values of each feature*

VIF analysis
Flight 8.57
Time 1.02
Length 1.28
Delay 3.57
AirportFrom_DEN 3.56
AirportFrom_DFW 2.71
AirportFrom_LAX 2.7
AirportFrom_ORD 2.24
Airline_DL 3.57
Airline_MQ 11.52
Airline_OO 2.4
Airline_WN 1.8
DayOfWeek_2 2.02
DayOfWeek_3 2.02
DayOfWeek_4 1.93
DayOfWeek_5 1.66
DayOfWeek_6 1.8

### ***Collinearity analysis:***

we chose to reduce the features from the random forest as it reduces the most number of features and hence, we reduce collinearity also as these are eliminated in the random forest as well.

### ***Observations:***

The features eliminated from random forest analysis are Airline\_MQ, AirportFrom\_DFW, AirportFrom\_LAX, Airline\_DL, AirportFrom\_DEN, and Airline\_OO and leaves us with 11 dimensions which is the least of all the dimensionality reduction techniques. Hence, I chose to follow random forest feature importances.

### ***Discretization and Binarization:***

One-hot encoding was applied to all categorical columns, with one column from each category dropped to avoid the dummy variable trap. The resulting data frame after one-hot encoding is shown below.

**Figure 10**

*One hot-encoded dataset*

One hot encoded dataset												
	Flight	Time	Length	Delay	AirportFrom_DEN	AirportFrom_DFW	AirportFrom_LAX	AirportFrom_ORD	Airline_DL	Airline_MQ	Airline_OO	Airline_WN
2	2400	20	165	1	0	0	1	0	0	0	0	0
13	2408	55	170	0	0	0	1	0	0	0	0	0
118	1297	345	160	0	0	1	0	0	0	0	0	0
228	2301	360	150	1	0	0	0	1	0	0	0	0
232	548	360	240	0	0	0	1	0	0	0	0	0

DayOfWeek_2	DayOfWeek_3	DayOfWeek_4	DayOfWeek_5	DayOfWeek_6	DayOfWeek_7
0	1	0	0	0	0
0	1	0	0	0	0
0	1	0	0	0	0
0	1	0	0	0	0
0	1	0	0	0	0

***Variable transformation:***

The dataframe was standardized prior to any analysis, model training, or dimensionality reduction techniques. Numerical columns were scaled, and the standardized data is shown below.

**Figure 8**

*Standardized dataset*

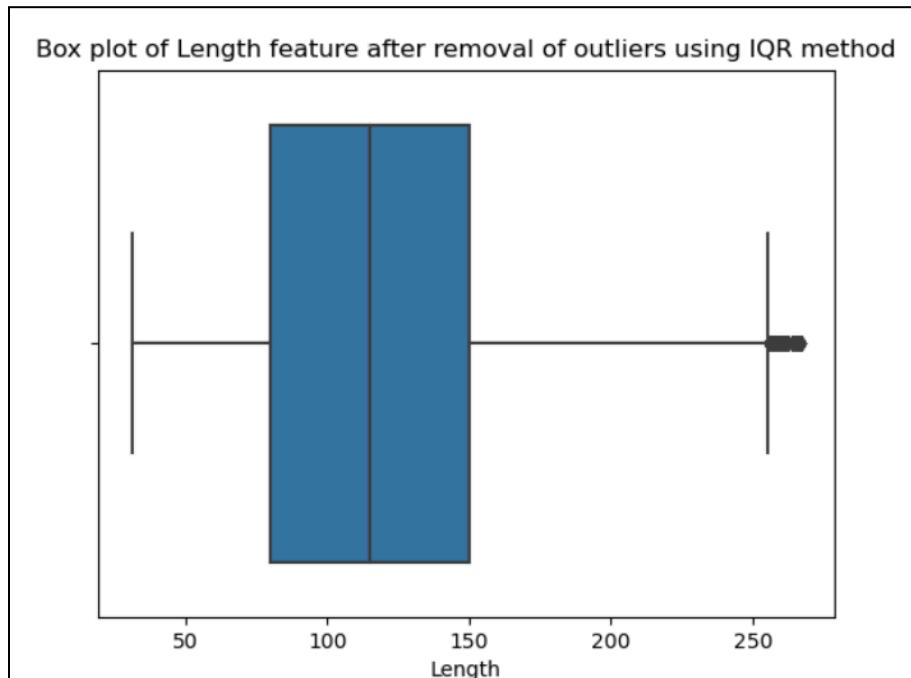
<b>Standardized dataset</b>			
	<b>Flight</b>	<b>Time</b>	<b>Length</b>
2	-0.242986	-3.006184	0.860139
13	-0.239076	-2.879481	0.957252
118	-0.781981	-1.829653	0.763026
228	-0.291363	-1.775351	0.568800
232	-1.147990	-1.775351	2.316836

***Outlier analysis and removal:***

The numerical column ‘Length’ was examined for outliers, which were removed using the IQR method. This approach involves calculating the first quartile (Q1) and third quartile (Q3) of the dataset. The interquartile range (IQR) is then determined, and the lower and upper bounds are defined as Q1 minus 1.5 times the IQR and Q3 plus 1.5 times the IQR, respectively. Any values outside these bounds are classified as outliers and removed.

**Figure 11**

*Box plot of Length feature after removal of outliers using IQR method*



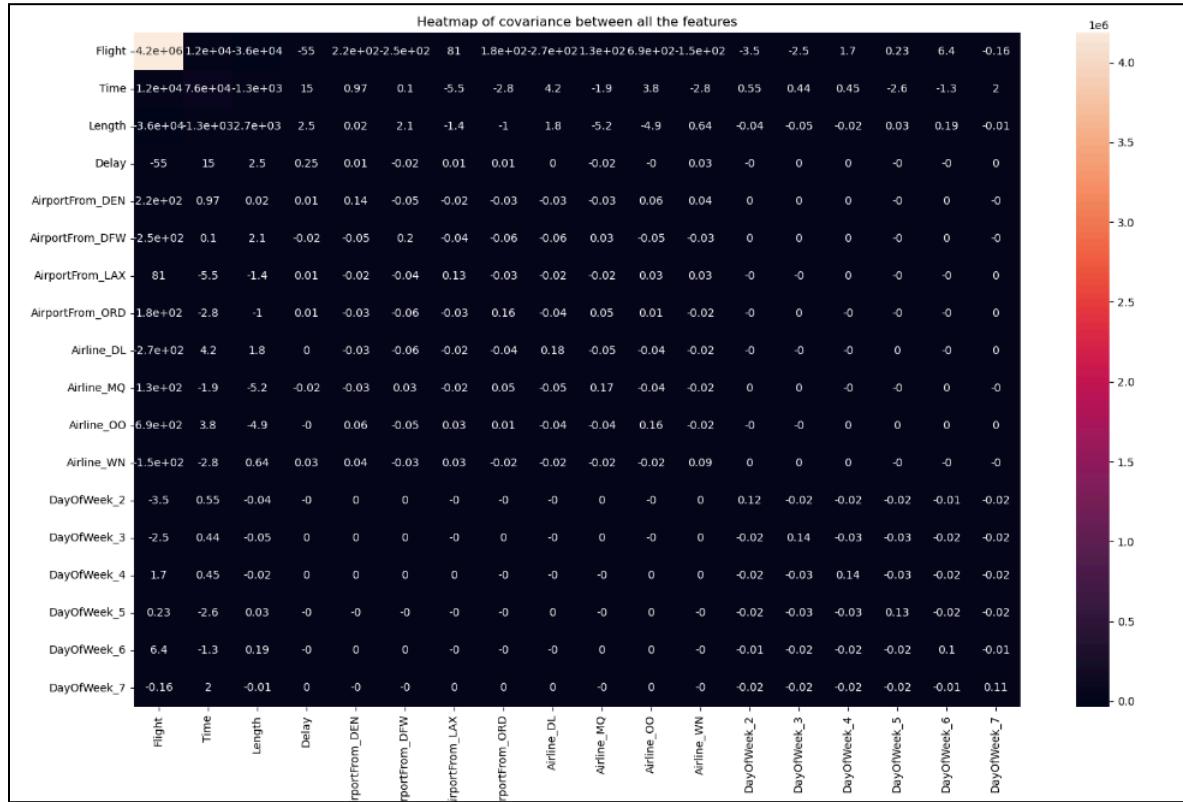
This box plot shows the Length after outlier removal. The outliers currently present are with respect to the updated column where the IQR, Q1 and Q3 changed. Surprisingly, this also removed outliers in the 'Time' column.

### **Covariance Matrix:**

Covariance is a statistical measure that indicates the extent to which two variables change together. It helps assess the direction of the linear relationship between two variables. The heatmap below shows the sample covariance matrix of the dataset across all features. The plot reveals that most values are quite small, likely because the majority of the features are categorical, resulting in minimal relationships between them.

**Figure 12**

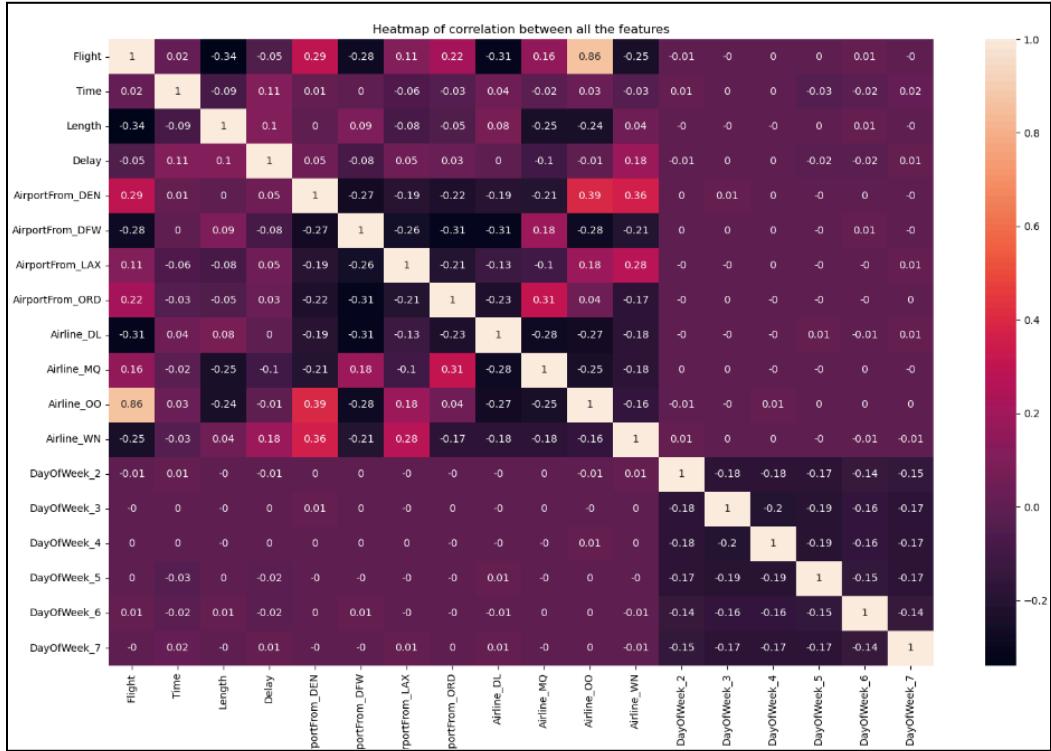
*The covariance matrix of all features*



### Pearson correlation matrix:

**Figure 13**

*Pearson correlation coefficient matrix*



The Pearson correlation matrix indicates that most values are close to 0, with a few features having correlations between 0.2 and 0.4. This suggests minimal correlation between features, primarily because most of them are categorical.

### Balanced or Imbalanced data:

Values in each class of the target variable are almost equal and hence, data is balanced with no bias problem in any further analysis.

**Figure 14**

*Count of each class of target 'Delay'*

The value counts of each class of target feature	
Delay	
0	19806
1	17764

## Phase II: Regression Analysis

### ***T-test analysis:***

The t-test is a statistical technique used to determine whether there is a significant difference between the means of two groups. It evaluates whether the observed differences are due to actual differences in the population or merely random chance. Among the various types of t-tests, the two-sample t-test is the most commonly used and compares the means of two independent groups to determine if they are significantly different. This is achieved by analyzing the group means, variances, and calculating a t-statistic, which is then compared to a critical value from the t-distribution at a chosen significance level (usually 0.05). If the t-value exceeds this critical value, it indicates that the observed difference is unlikely to result from random variation, signifying a statistically significant difference between the groups.

In this analysis, features such as Time, Delay, AirportFrom\_DEN, AirportFrom\_DFW, AirportFrom\_LAX, AirportFrom\_ORD, Airline\_MQ, Airline\_OO, and Airline\_WN exhibit very low p-values (close to 0), highlighting their strong significance in predicting flight delays. Conversely, features like DayOfWeek\_2, DayOfWeek\_3, DayOfWeek\_4, DayOfWeek\_5, DayOfWeek\_6, and DayOfWeek\_7 have higher p-values (above 0.05), indicating they may not significantly influence flight delay predictions based on this analysis.

**Figure 15**

*T-statistic and p-values of features*

t-test analysis			
Feature	t-statistic	p-value	
const	2.22	0.03	
Flight	-1.52	0.13	
Time	-16.24	0.0	
Delay	15.11	0.0	
AirportFrom_DEN	27.45	0.0	
AirportFrom_DFW	12.54	0.0	
AirportFrom_LAX	15.57	0.0	
AirportFrom_ORD	18.46	0.0	
Airline_DL	-3.52	0.0	
Airline_MQ	-45.77	0.0	
Airline_OO	-27.38	0.0	
Airline_WN	-28.71	0.0	
DayOfWeek_2	0.45	0.65	
DayOfWeek_3	-0.17	0.86	
DayOfWeek_4	0.25	0.8	
DayOfWeek_5	0.43	0.67	
DayOfWeek_6	1.48	0.14	
DayOfWeek_7	0.75	0.45	

### **F-test analysis:**

The F-test is a statistical technique used to compare the variability of two or more groups to determine whether their variances differ significantly. It assesses whether the observed variances within these groups are equal or if the differences exceed what can be attributed to random variation. In the context of Analysis of Variance (ANOVA), the F-test is frequently used to analyze the variances of multiple groups simultaneously. The F-statistic is calculated as the ratio of variance between groups to variance within groups, helping to identify if the variability between group means is significantly larger than within-group variability. By comparing the F-value to a critical value from the F-distribution at a chosen significance level (typically 0.05), it can be determined whether the variance differences are statistically significant.

In this analysis, features with high F-test scores, such as Flight, Time, Airline\_DL, Airline\_MQ, and Airline\_OO, demonstrate greater importance in explaining the variability of the target variable. Conversely, features with low F-test scores, such as DayOfWeek\_2, DayOfWeek\_3, DayOfWeek\_4, DayOfWeek\_5, and DayOfWeek\_7, show lesser significance in contributing to the variability. Additionally, p-values close to 0.0 provide strong evidence against the null hypothesis, indicating that the overall set of predictors plays a significant role in explaining the target variable.

**Figure 16**

*F-test and p-values of each feature*

Feature	F-test score	p-value
Flight	89.3	0.0
Time	14.91	0.0
Length	4.03	0.0
AirportFrom_DEN	37.04	0.0
AirportFrom_DFW	48.37	0.0
AirportFrom_LAX	65.57	0.0
AirportFrom_ORD	24.22	0.0
Airline_DL	112.05	0.0
Airline_MQ	65.06	0.0
Airline_OO	122.7	0.0
Airline_WN	23.58	0.0
DayOfWeek_2	0.53	1.0
DayOfWeek_3	0.67	1.0
DayOfWeek_4	0.58	1.0
DayOfWeek_5	0.73	1.0
DayOfWeek_6	1.25	0.01
DayOfWeek_7	0.97	0.62

### ***Linear regression model and prediction of dependent variable:***

Linear regression is a fundamental statistical method used to model the relationship between a dependent variable (y) and one or more independent variables (x). Its main goal is to analyze how variations in the independent variables influence the dependent variable. The model operates under the assumption of a linear relationship between the independent and dependent variables.

**Figure 17**

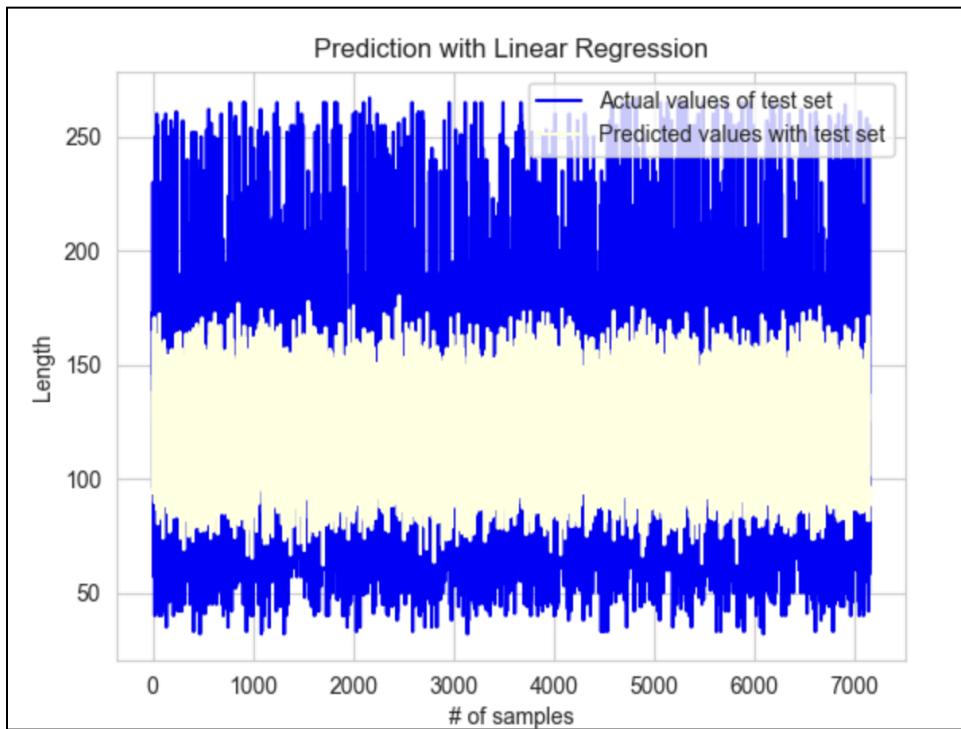
*Coefficients and R2, Adj. R2, AIC, BIC, MSE of the Linear Regression Model*

The coefficients of the Linear Regression model					
Feature	Coefficient				
const	0.0				
Flight	-0.02				
Time	-0.09				
Delay	0.16				
AirportFrom_DEN	0.8				
AirportFrom_DFW	0.35				
AirportFrom_LAX	0.43				
AirportFrom_ORD	0.51				
Airline_DL	-0.09				
Airline_MQ	-1.04				
Airline_OO	-1.22				
Airline_WN	-0.74				
DayOfWeek_2	0.01				
DayOfWeek_3	-0.0				
DayOfWeek_4	0.0				
DayOfWeek_5	0.01				
DayOfWeek_6	0.03				
DayOfWeek_7	0.02				
Mean Squared Error for training set (MSE): 2054.446					
Model	R-squared	Adj. R-squared	AIC	BIC	Mean Squared Error
Linear Regression	0.23	0.22	218035.43	218175.86	2054.45

Linear regression was performed on the original feature set after standardizing the numerical features and the target variable, ‘Length’. The intercept value of 0.0 indicates that the predicted value of the target is zero when all other predictors are zero. The model explains a moderate portion of the variability in the target variable ‘Length’ which is around 23%.

**Figure 18**

*Predicted and Test values of ‘Length’ using Linear Regression*

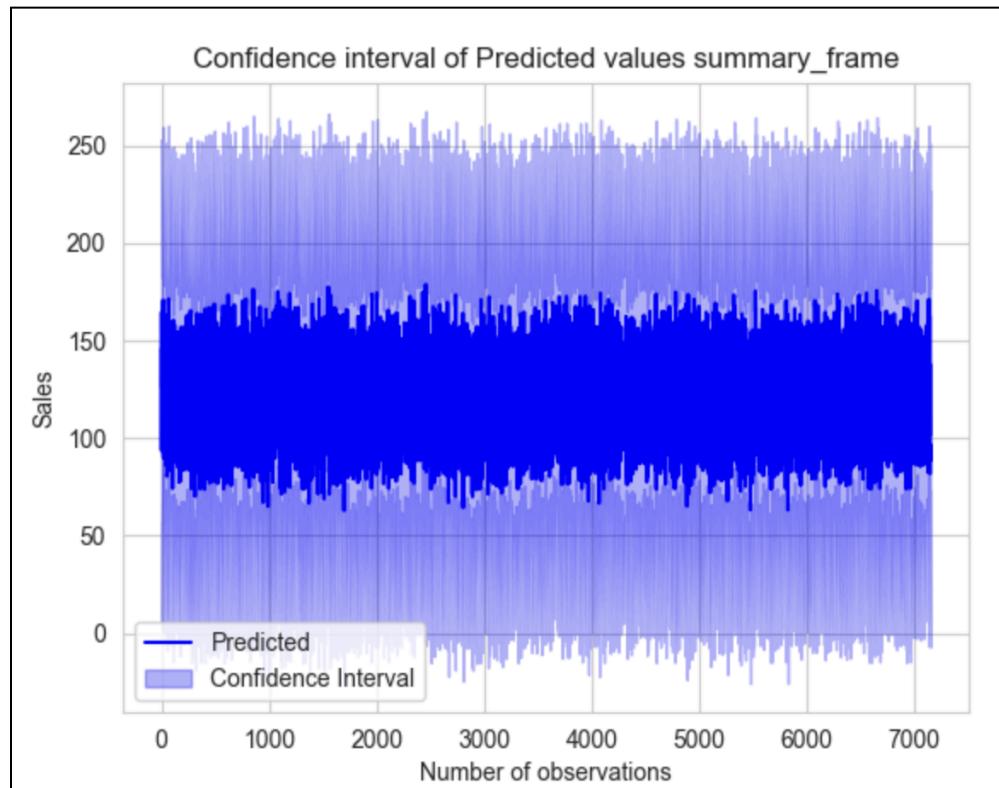


The training and predicted values can be seen in the above graph and the model has bias between predicted and training values.

#### ***Confidence interval analysis:***

The below plot shows the confidence interval of the predicted values which has huge variance from the actual predicted values. This indicates that there is uncertainty in the estimated coefficients and this model is not reliable in predicting the ‘Length’. This uncertainty is expected as our dataset is a classification problem and this is an independent variable and helps in predicting Delay rather than being predicted.

**Figure 19**  
*Confidence Interval analysis*



***Stepwise regression and adjusted R-square analysis:***

This method begins with a model that includes all predictors, gradually removing the least significant variables until no further improvement in the model's performance is achieved. Backward stepwise regression is chosen as some features exhibit higher p-values in the T-test. Features with the highest p-value, exceeding the threshold of 0.01, are eliminated step by step. This process reduces the number of features while retaining those sufficient to explain the variance in the target variable. At each step, the Adjusted R-squared from the OLS summary is monitored, and the process is halted if its value drops, indicating the need to retain the recently removed feature.

Features DayOfWeek\_3, DayOfWeek\_4, DayOfWeek\_5, DayOfWeek\_2, DayOfWeek\_7, Flight, and DayOfWeek\_6 eliminated.

```
eliminated features: ['DayofWeek_3', 'DayofWeek_4', 'DayofWeek_5', 'DayofWeek_2', 'DayofWeek_7', 'Flight', 'DayofWeek_6']
```

**Figure 20**

*OLS summary after eliminating the 'DayOfWeek\_3' feature*

Eliminating the feature with highest p-value which is DayOfWeek_3									
OLS Regression Results									
Dep. Variable:	y	R-squared:	0.225						
Model:	OLS	Adj. R-squared:	0.225						
Method:	Least Squares	F-statistic:	519.3						
Date:	Mon, 02 Dec 2024	Prob (F-statistic):	0.00						
Time:	11:50:43	Log-Likelihood:	-36911.						
No. Observations:	28580	AIC:	7.386e+04						
Df Residuals:	28563	BIC:	7.400e+04						
Df Model:	16								
Covariance Type:	nonrobust								
	coef	std err	t	P> t	[0.025	0.975]			
const	0.0715	0.031	2.286	0.022	0.010	0.133			
Flight	-0.0248	0.016	-1.523	0.128	-0.057	0.007			
Time	-0.0859	0.005	-16.242	0.000	-0.096	-0.076			
Delay	0.1625	0.011	15.114	0.000	0.141	0.184			
AirportFrom_DEN	0.8021	0.029	27.448	0.000	0.745	0.859			
AirportFrom_DFW	0.3504	0.028	12.539	0.000	0.296	0.405			
AirportFrom_LAX	0.4303	0.028	15.572	0.000	0.376	0.484			
AirportFrom_ORD	0.5138	0.028	18.465	0.000	0.459	0.568			
Airline_DL	-0.0948	0.027	-3.522	0.000	-0.147	-0.042			
Airline_MQ	-1.0358	0.023	-45.766	0.000	-1.080	-0.991			
Airline_OO	-1.2170	0.044	-27.379	0.000	-1.304	-1.130			
Airline_WN	-0.7431	0.026	-28.709	0.000	-0.794	-0.692			
DayOfWeek_2	0.0110	0.017	0.641	0.521	-0.023	0.045			
DayOfWeek_4	0.0066	0.016	0.416	0.678	-0.025	0.038			
DayOfWeek_5	0.0102	0.016	0.625	0.532	-0.022	0.042			
DayOfWeek_6	0.0332	0.018	1.823	0.068	-0.002	0.069			
DayOfWeek_7	0.0171	0.017	0.987	0.324	-0.017	0.051			

**Figure 21**

*OLS summary after eliminating 'DayOfWeek\_4' feature*

Eliminating DayOfWeek_4									
OLS Regression Results									
Dep. Variable:	y	R-squared:	0.225						
Model:	OLS	Adj. R-squared:	0.225						
Method:	Least Squares	F-statistic:	553.9						
Date:	Mon, 02 Dec 2024	Prob (F-statistic):	0.00						
Time:	11:50:43	Log-Likelihood:	-36911.						
No. Observations:	28580	AIC:	7.385e+04						
Df Residuals:	28564	BIC:	7.399e+04						
Df Model:	15								
Covariance Type:	nonrobust								
	coef	std err	t	P> t	[0.025	0.975]			
const	0.0740	0.031	2.409	0.016	0.014	0.134			
Flight	-0.0248	0.016	-1.521	0.128	-0.057	0.007			
Time	-0.0859	0.005	-16.241	0.000	-0.096	-0.076			
Delay	0.1625	0.011	15.113	0.000	0.141	0.184			
AirportFrom_DEN	0.8019	0.029	27.445	0.000	0.745	0.859			
AirportFrom_DFW	0.3503	0.028	12.536	0.000	0.296	0.405			
AirportFrom_LAX	0.4302	0.028	15.569	0.000	0.376	0.484			
AirportFrom_ORD	0.5137	0.028	18.462	0.000	0.459	0.568			
Airline_DL	-0.0948	0.027	-3.524	0.000	-0.148	-0.042			
Airline_MQ	-1.0358	0.023	-45.767	0.000	-1.080	-0.991			
Airline_OO	-1.2169	0.044	-27.378	0.000	-1.304	-1.130			
Airline_WN	-0.7430	0.026	-28.707	0.000	-0.794	-0.692			
DayOfWeek_2	0.0086	0.016	0.533	0.594	-0.023	0.040			
DayOfWeek_5	0.0078	0.015	0.512	0.609	-0.022	0.038			
DayOfWeek_6	0.0309	0.017	1.782	0.075	-0.003	0.065			
DayOfWeek_7	0.0147	0.016	0.900	0.368	-0.017	0.047			

**Figure 22**

*OLS summary after eliminating 'DayOfWeek\_5' feature*

Eliminating DayOfWeek_5									
OLS Regression Results									
Dep. Variable:	y	R-squared:	0.225						
Model:	OLS	Adj. R-squared:	0.225						
Method:	Least Squares	F-statistic:	593.5						
Date:	Mon, 02 Dec 2024	Prob (F-statistic):	0.00						
Time:	11:50:43	Log-Likelihood:	-36911.						
No. Observations:	28580	AIC:	7.385e+04						
Df Residuals:	28565	BIC:	7.398e+04						
Df Model:	14								
Covariance Type:	nonrobust								
	coef	std err	t	P> t	[0.025	0.975]			
const	0.0760	0.030	2.497	0.013	0.016	0.136			
Flight	-0.0247	0.016	-1.517	0.129	-0.057	0.007			
Time	-0.0859	0.005	-16.255	0.000	-0.096	-0.076			
Delay	0.1624	0.011	15.106	0.000	0.141	0.183			
AirportFrom_DEN	0.8018	0.029	27.442	0.000	0.745	0.859			
AirportFrom_DFW	0.3503	0.028	12.536	0.000	0.296	0.405			
AirportFrom_LAX	0.4301	0.028	15.567	0.000	0.376	0.484			
AirportFrom_ORD	0.5137	0.028	18.461	0.000	0.459	0.568			
Airline_DL	-0.0948	0.027	-3.523	0.000	-0.148	-0.042			
Airline_MQ	-1.0359	0.023	-45.770	0.000	-1.080	-0.992			
Airline_OO	-1.2170	0.044	-27.381	0.000	-1.304	-1.130			
Airline_WN	-0.7429	0.026	-28.705	0.000	-0.794	-0.692			
DayOfWeek_2	0.0067	0.016	0.425	0.671	-0.024	0.037			
DayOfWeek_6	0.0289	0.017	1.712	0.087	-0.004	0.062			
DayOfWeek_7	0.0128	0.016	0.804	0.422	-0.018	0.044			

**Figure 23**

*OLS summary after eliminating 'DayOfWeek\_2' feature*

Eliminating DayOfWeek_2									
OLS Regression Results									
Dep. Variable:	y	R-squared:	0.225						
Model:	OLS	Adj. R-squared:	0.225						
Method:	Least Squares	F-statistic:	639.2						
Date:	Mon, 02 Dec 2024	Prob (F-statistic):	0.00						
Time:	11:50:43	Log-Likelihood:	-36911.						
No. Observations:	28580	AIC:	7.385e+04						
Df Residuals:	28566	BIC:	7.397e+04						
Df Model:	13								
Covariance Type:	nonrobust								
	coef	std err	t	P> t	[0.025	0.975]			
const	0.0773	0.030	2.550	0.011	0.018	0.137			
Flight	-0.0247	0.016	-1.515	0.130	-0.057	0.007			
Time	-0.0859	0.005	-16.251	0.000	-0.096	-0.076			
Delay	0.1624	0.011	15.102	0.000	0.141	0.183			
AirportFrom_DEN	0.8018	0.029	27.442	0.000	0.745	0.859			
AirportFrom_DFW	0.3503	0.028	12.536	0.000	0.296	0.405			
AirportFrom_LAX	0.4301	0.028	15.566	0.000	0.376	0.484			
AirportFrom_ORD	0.5137	0.028	18.462	0.000	0.459	0.568			
Airline_DL	-0.0948	0.027	-3.525	0.000	-0.148	-0.042			
Airline_MQ	-1.0359	0.023	-45.773	0.000	-1.080	-0.992			
Airline_OO	-1.2171	0.044	-27.385	0.000	-1.304	-1.130			
Airline_WN	-0.7429	0.026	-28.704	0.000	-0.794	-0.692			
DayOfWeek_6	0.0277	0.017	1.665	0.096	-0.005	0.060			
DayOfWeek_7	0.0116	0.016	0.741	0.459	-0.019	0.042			

**Figure 24**

*OLS summary after eliminating 'DayOfWeek\_7' feature*

Eliminating DayOfWeek_7						
OLS Regression Results						
Dep. Variable:	y	R-squared:	0.225			
Model:	OLS	Adj. R-squared:	0.225			
Method:	Least Squares	F-statistic:	692.4			
Date:	Mon, 02 Dec 2024	Prob (F-statistic):	0.00			
Time:	11:50:43	Log-Likelihood:	-36911.			
No. Observations:	28580	AIC:	7.385e+04			
Df Residuals:	28567	BIC:	7.396e+04			
Df Model:	12					
Covariance Type:	nonrobust					
	coef	std err	t	P> t	[0.025	0.975]
const	0.0789	0.030	2.610	0.009	0.020	0.138
Flight	-0.0247	0.016	-1.517	0.129	-0.057	0.007
Time	-0.0858	0.005	-16.243	0.000	-0.096	-0.075
Delay	0.1624	0.011	15.108	0.000	0.141	0.183
AirportFrom_DEN	0.8018	0.029	27.445	0.000	0.745	0.859
AirportFrom_DFW	0.3503	0.028	12.538	0.000	0.296	0.405
AirportFrom_LAX	0.4302	0.028	15.571	0.000	0.376	0.484
AirportFrom_ORD	0.5138	0.028	18.466	0.000	0.459	0.568
Airline_DL	-0.0948	0.027	-3.524	0.000	-0.148	-0.042
Airline_MQ	-1.0360	0.023	-45.776	0.000	-1.080	-0.992
Airline_OO	-1.2171	0.044	-27.385	0.000	-1.304	-1.130
Airline_WN	-0.7431	0.026	-28.712	0.000	-0.794	-0.692
DayOfWeek_6	0.0260	0.017	1.578	0.115	-0.006	0.058

**Figure 25**

*OLS summary after eliminating 'Flight' feature*

Eliminating Flight									
OLS Regression Results									
Dep. Variable:	y	R-squared:	0.225						
Model:	OLS	Adj. R-squared:	0.225						
Method:	Least Squares	F-statistic:	755.1						
Date:	Mon, 02 Dec 2024	Prob (F-statistic):	0.00						
Time:	11:50:43	Log-Likelihood:	-36912.						
No. Observations:	28580	AIC:	7.385e+04						
Df Residuals:	28568	BIC:	7.395e+04						
Df Model:	11								
Covariance Type:	nonrobust								
	coef	std err	t	P> t	[0.025	0.975]			
const	0.0964	0.028	3.449	0.001	0.042	0.151			
Time	-0.0857	0.005	-16.215	0.000	-0.096	-0.075			
Delay	0.1628	0.011	15.144	0.000	0.142	0.184			
AirportFrom_DEN	0.8004	0.029	27.410	0.000	0.743	0.858			
AirportFrom_DFW	0.3529	0.028	12.651	0.000	0.298	0.408			
AirportFrom_LAX	0.4303	0.028	15.572	0.000	0.376	0.484			
AirportFrom_ORD	0.5103	0.028	18.404	0.000	0.456	0.565			
Airline_DL	-0.0984	0.027	-3.670	0.000	-0.151	-0.046			
Airline_MQ	-1.0608	0.016	-67.740	0.000	-1.091	-1.030			
Airline_OO	-1.2768	0.021	-61.630	0.000	-1.317	-1.236			
Airline_WN	-0.7419	0.026	-28.680	0.000	-0.793	-0.691			
DayOfWeek_6	0.0253	0.016	1.537	0.124	-0.007	0.058			

**Figure 26**

*OLS summary after eliminating DayOfWeek\_6 feature*

Eliminating DayOfWeek_6									
OLS Regression Results									
Dep. Variable:	y	R-squared:	0.225						
Model:	OLS	Adj. R-squared:	0.225						
Method:	Least Squares	F-statistic:	830.3						
Date:	Mon, 02 Dec 2024	Prob (F-statistic):	0.00						
Time:	11:50:43	Log-Likelihood:	-36914.						
No. Observations:	28580	AIC:	7.385e+04						
Df Residuals:	28569	BIC:	7.394e+04						
Df Model:	10								
Covariance Type:	nonrobust								
	coef	std err	t	P> t	[0.025	0.975]			
const	0.0995	0.028	3.571	0.000	0.045	0.154			
Time	-0.0858	0.005	-16.243	0.000	-0.096	-0.075			
Delay	0.1626	0.011	15.129	0.000	0.142	0.184			
AirportFrom_DEN	0.8005	0.029	27.412	0.000	0.743	0.858			
AirportFrom_DFW	0.3528	0.028	12.648	0.000	0.298	0.407			
AirportFrom_LAX	0.4303	0.028	15.572	0.000	0.376	0.484			
AirportFrom_ORD	0.5101	0.028	18.397	0.000	0.456	0.564			
Airline_DL	-0.0987	0.027	-3.681	0.000	-0.151	-0.046			
Airline_MQ	-1.0608	0.016	-67.745	0.000	-1.092	-1.030			
Airline_OO	-1.2768	0.021	-61.633	0.000	-1.317	-1.236			
Airline_WN	-0.7423	0.026	-28.697	0.000	-0.793	-0.692			

**Figure 27**

*OLS summary after Backward regression*

OLS Regression Results									
Dep. Variable:	y	R-squared:	0.225						
Model:	OLS	Adj. R-squared:	0.225						
Method:	Least Squares	F-statistic:	830.3						
Date:	Mon, 02 Dec 2024	Prob (F-statistic):	0.00						
Time:	11:50:43	Log-Likelihood:	-36914.						
No. Observations:	28580	AIC:	7.385e+04						
Df Residuals:	28569	BIC:	7.394e+04						
Df Model:	10								
Covariance Type:	nonrobust								
	coef	std err	t	P> t	[0.025	0.975]			
const	0.0995	0.028	3.571	0.000	0.045	0.154			
Time	-0.0858	0.005	-16.243	0.000	-0.096	-0.075			
Delay	0.1626	0.011	15.129	0.000	0.142	0.184			
AirportFrom_DEN	0.8005	0.029	27.412	0.000	0.743	0.858			
AirportFrom_DFW	0.3528	0.028	12.648	0.000	0.298	0.407			
AirportFrom_LAX	0.4303	0.028	15.572	0.000	0.376	0.484			
AirportFrom_ORD	0.5101	0.028	18.397	0.000	0.456	0.564			
Airline_DL	-0.0987	0.027	-3.681	0.000	-0.151	-0.046			
Airline_MQ	-1.0608	0.016	-67.745	0.000	-1.092	-1.030			
Airline_OO	-1.2768	0.021	-61.633	0.000	-1.317	-1.236			
Airline_WN	-0.7423	0.026	-28.697	0.000	-0.793	-0.692			
Omnibus:	1539.012	Durbin-Watson:	2.009						
Prob(Omnibus):	0.000	Jarque-Bera (JB):	1798.650						
Skew:	0.611	Prob(JB):	0.00						
Kurtosis:	3.133	Cond. No.	14.6						

We find that the features const, Time, Delay, AirportFrom\_DEN, AirportFrom\_DFW, AirportFrom\_LAX, AirportFrom\_ORD, Airline\_DL, Airline\_MQ, Airline\_OO, Airline\_WN are enough to perform linear regression and achieve Adj. R-squared as 23%.

Features DayOfWeek\_3, DayOfWeek\_4, DayOfWeek\_5, DayOfWeek\_2, DayOfWeek\_7, Flight, and DayOfWeek\_6 eliminated.

**Figure 28**

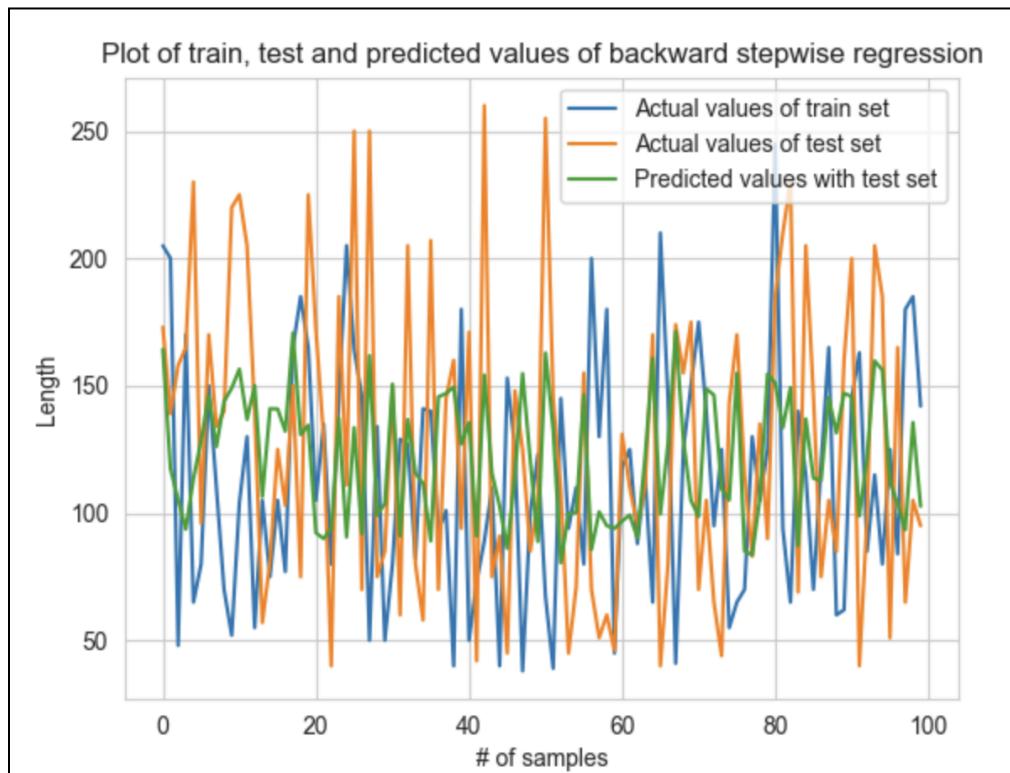
*Condition number and Singular values after Backward stepwise regression*

```
Condition number of the remaining features of the dataset after backward stepwise regression
14.65
Singular values of the remaining features of the dataset after backward stepwise regression
[212.6 169.2 96.36 84.95 81.07 76.89 66.15 61.43 49.75 24.56
 14.52]
```

The condition number and singular values suggest that there's no collinearity and hence, these are our final set of features that are sufficient.

**Figure 29**

*Plot of train, test, and predicted values of Backward stepwise regression*



The above plot shows the train, test, and predicted values after we find the final set of features from backward stepwise regression.

### ***Comparing Linear regression and Backward stepwise regression***

**Figure 30**

*Comparison of Linear regression and Backward stepwise regression*

Model	R-squared	Adj. R-squared	AIC	BIC	Mean Squared Error
Linear Regression	0.23	0.22	218035.43	218175.86	2054.45
Backward Stepwise Regression	0.23	0.22	73849.35	73940.22	2032.57

Both the analyses performed equally well in explaining the target variable's variability but, with backward regression, we significantly reduced AIC and BIC values and also reduced MSE. We are able to achieve the same results with reduced feature space in backward stepwise

## **Phase III: Classification Analysis**

### ***Decision Tree:***

#### ***Pre-pruned decision tree:***

Pruning techniques focus on balancing model complexity and predictive accuracy. To prevent the decision tree from overfitting by growing unchecked and perfectly fitting the training data, pre-pruning introduces constraints on the tree's growth or depth during its construction. This ensures the tree does not expand to its maximum allowable depth. Grid search and stratified k-fold cross-validation were performed to train the model and identify the best hyperparameters.

**Figure 31**

*Best hyperparameters of Pre-pruned decision tree*

```
Best Hyperparameters: {'criterion': 'entropy', 'max_depth': 9, 'max_features': 7, 'min_samples_leaf': 4, 'min_samples_split': 10, 'splitter': 'best'}
```

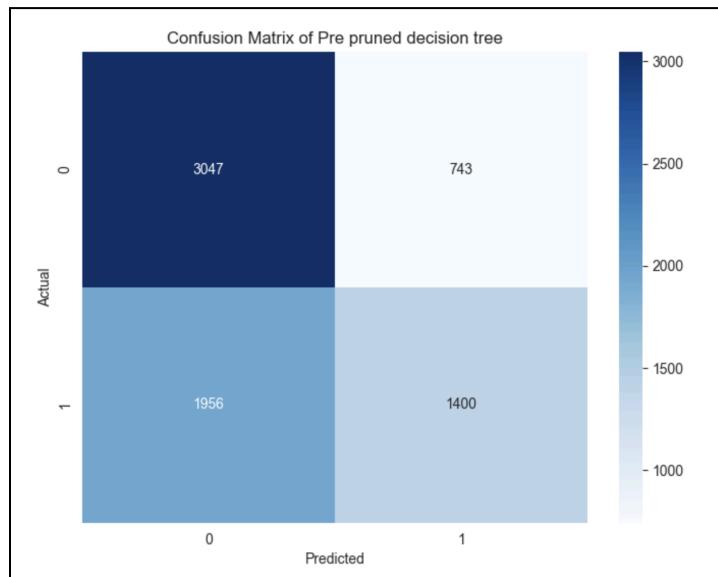
**Figure 32**

*Various metrics of Pre-pruned decision tree*

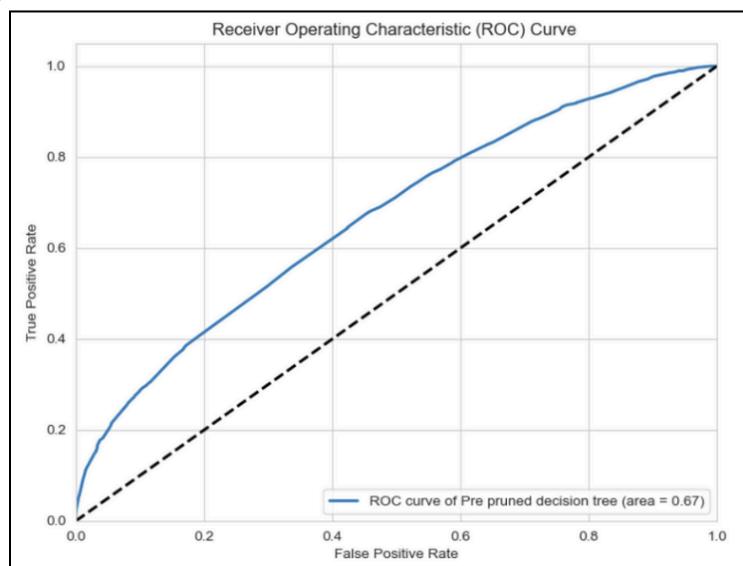
Classifier	Accuracy	Confusion Matrix	Precision	Sensitivity or Recall	Specificity	F-score	AUC
Pre pruned Decision Tree Classifier	0.62	[[3047 743] [1956 1400]]	0.65	0.42	0.8	0.51	0.67

**Figure 33**

*Confusion matrix of Pre-pruned decision tree*

**Figure 34**

*ROC curve of Pre-pruned decision tree*



### ***Post-pruned decision tree:***

Post-pruning, also referred to as tree pruning or cost complexity pruning, is the process of removing or collapsing branches (subtrees) from a fully grown decision tree. Unlike pre-pruning, which applies constraints during the tree's construction, post-pruning is performed after the tree is completely built. The main goal of post-pruning is to reduce the tree's size, minimize overfitting, and enhance its ability to generalize to new data. The parameter `ccp_alpha` governs the tree's complexity, with higher values resulting in more aggressive pruning. The optimal alpha value determined is 0.00014.

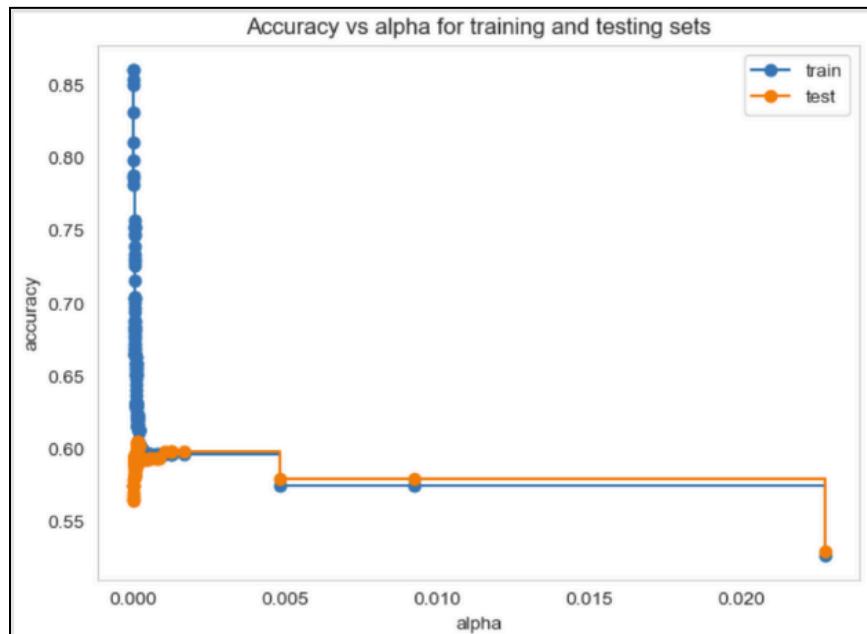
**Figure 35**

*Optimal alpha value of Post-pruned decision tree*

Best alpha value of post pruned decision tree: 0.00014000098283027543

**Figure 36**

*Accuracy vs. Alpha plot*



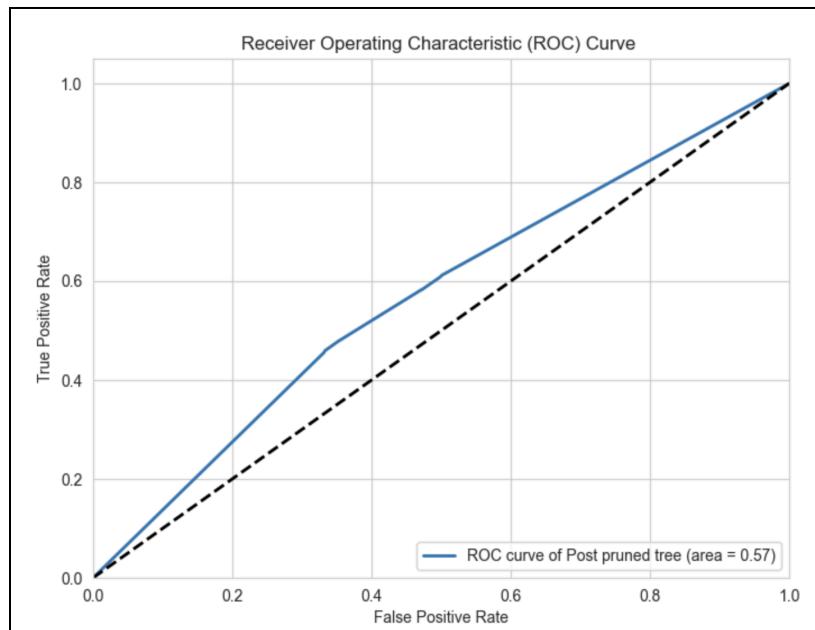
**Figure 37**

*Various metrics of Post-pruned decision tree*

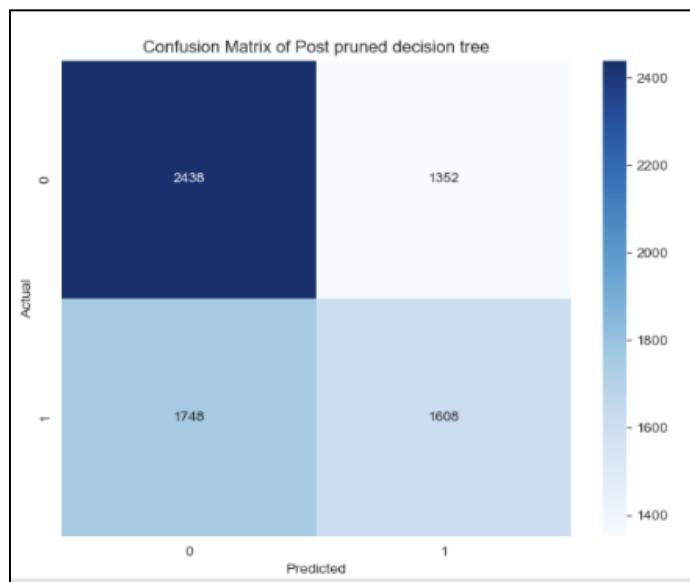
Classifier	Accuracy	Confusion Matrix	Precision	Sensitivity or Recall	Specificity	F-score	AUC
Post pruned Decision Tree Classifier	0.57	[[2456 1334] [1751 1605]]	0.55	0.48	0.8	0.51	0.57

**Figure 38**

*ROC curve of Post-pruned decision tree*

**Figure 39**

*Confusion matrix of Post-pruned decision tree*



### ***Logistic Regression:***

Logistic Regression is a statistical approach commonly applied to binary classification problems, where it predicts the probability of a categorical dependent variable, typically represented as 0 or 1. A grid search was conducted over the parameters ‘C,’ ‘Solver,’ and ‘Penalty,’ combined with stratified k-fold cross-validation, to identify the optimal hyperparameters.

**Figure 40**

*Best hyperparameters of Logistic Regression*

**Best Hyperparameters: {'C': 0.1, 'penalty': 'l2', 'solver': 'saga'}**

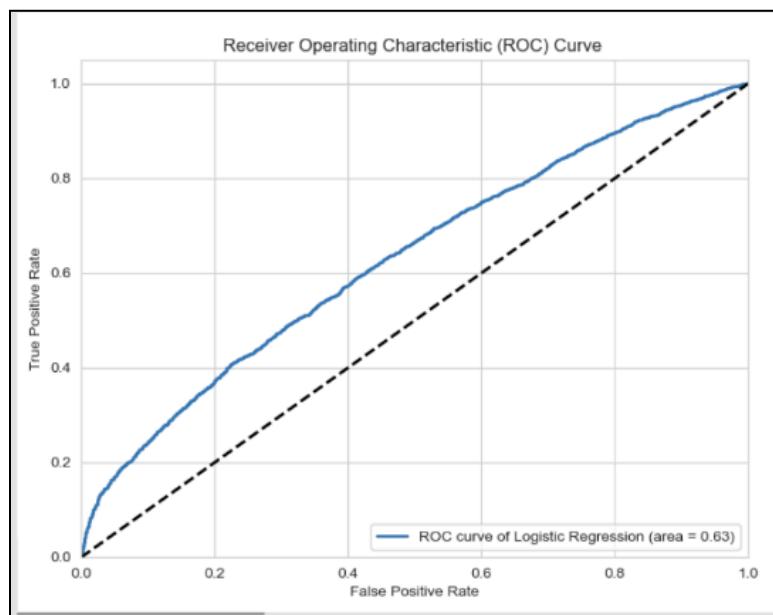
**Figure 41**

*Various metrics of Logistic Regression*

Classifier	Accuracy	Confusion Matrix	Precision	Sensitivity or Recall	Specificity	F-score	AUC
Logistic Regression	0.59	[[2872 918] [1947 1409]]	0.61	0.42	0.76	0.5	0.63

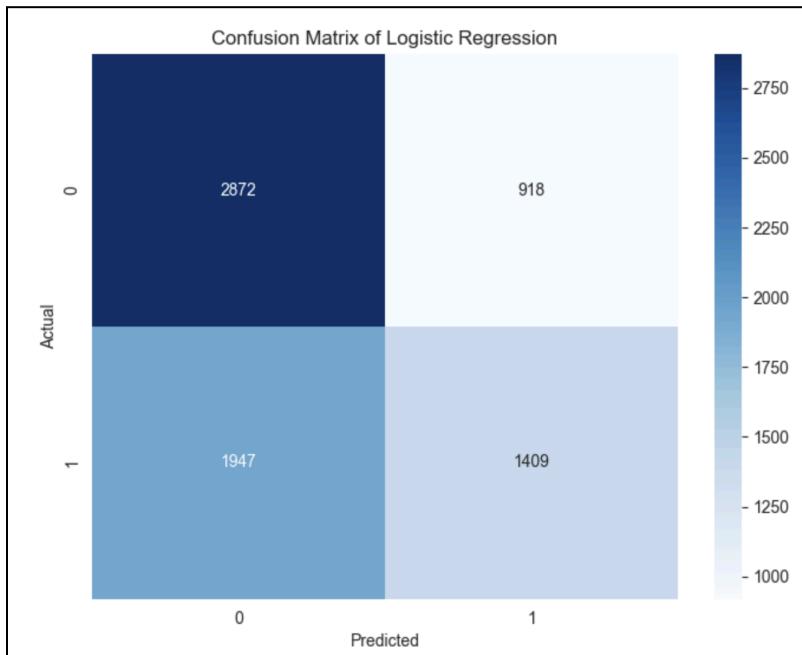
**Figure 42**

*ROC curve of Logistic Regression*



**Figure 43**

*Confusion matrix of Logistic Regression*

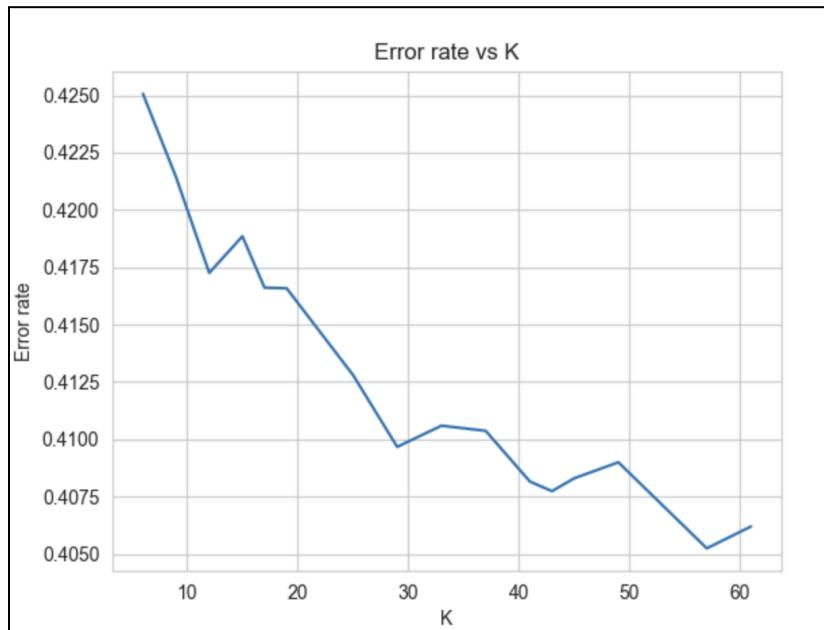


### **K Nearest Neighbors:**

K-Nearest Neighbors (KNN) is a straightforward and efficient supervised learning algorithm applicable to both classification and regression tasks. The 'K' represents the number of nearest neighbors considered for making predictions. To classify or predict a new data point, the algorithm identifies its K closest neighbors using a specified distance metric. Grid search and stratified k-fold cross-validation were employed to determine the optimal hyperparameters. The best value of K, identified from the error rate vs. K graph, closely aligns with the grid search results.

**Figure 44**

*Elbow method for finding optimal k*

**Figure 45**

*Best hyperparameters of K-Nearest Neighbors*

```
Best Hyperparameters: {'n_neighbors': 57}
```

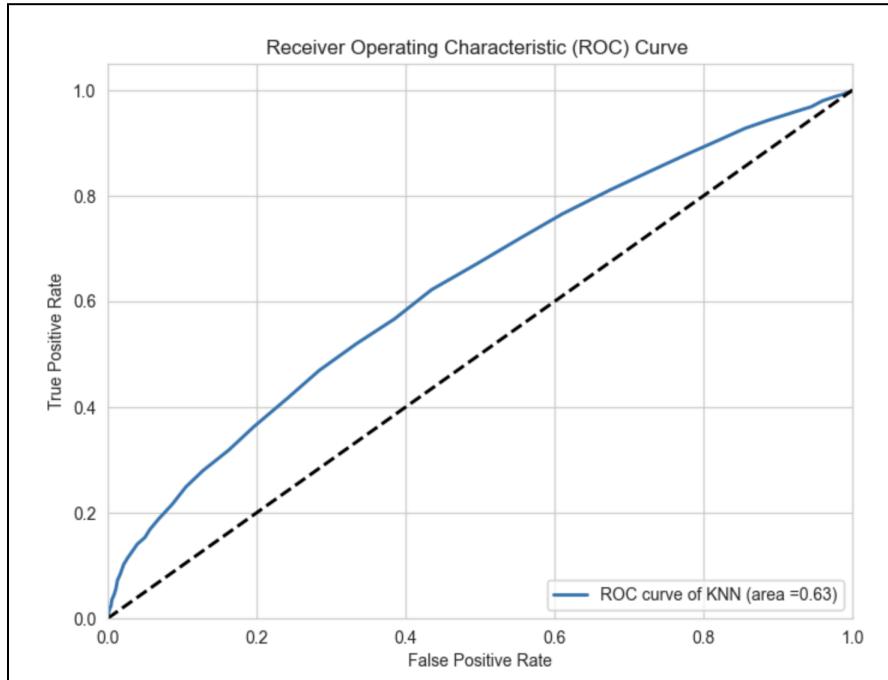
**Figure 46**

*Various metrics of K-Nearest Neighbors*

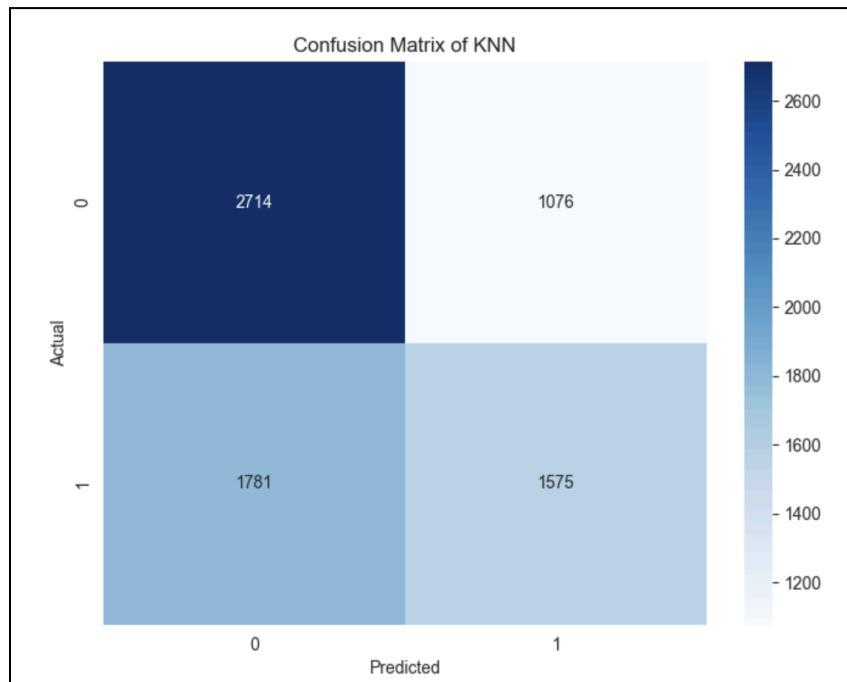
Classifier	Accuracy	Confusion Matrix	Precision	Sensitivity or Recall	Specificity	F-score	AUC
KNN	0.62	[[2714 1076] [1781 1575]]	0.59	0.47	0.72	0.52	0.63

**Figure 47**

*ROC curve of K-Nearest Neighbors*

**Figure 48**

*Confusion matrix of K-Nearest Neighbors*



**SVM:**

Support Vector Machines (SVM) identify the hyperplane that maximizes the margin between the closest data points, known as support vectors, from different classes. A grid search was conducted with linear, polynomial, and RBF kernels, along with stratified k-fold cross-validation, to determine the optimal kernel. The radial basis function (RBF) kernel, commonly used for nonlinear problems, was found to be the best choice.

**Figure 49**

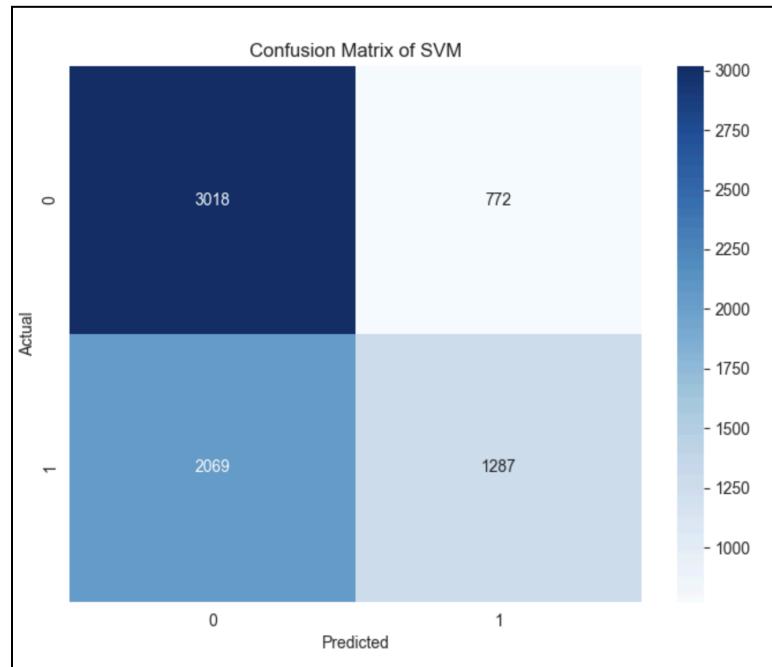
*Best hyperparameters of SVM*

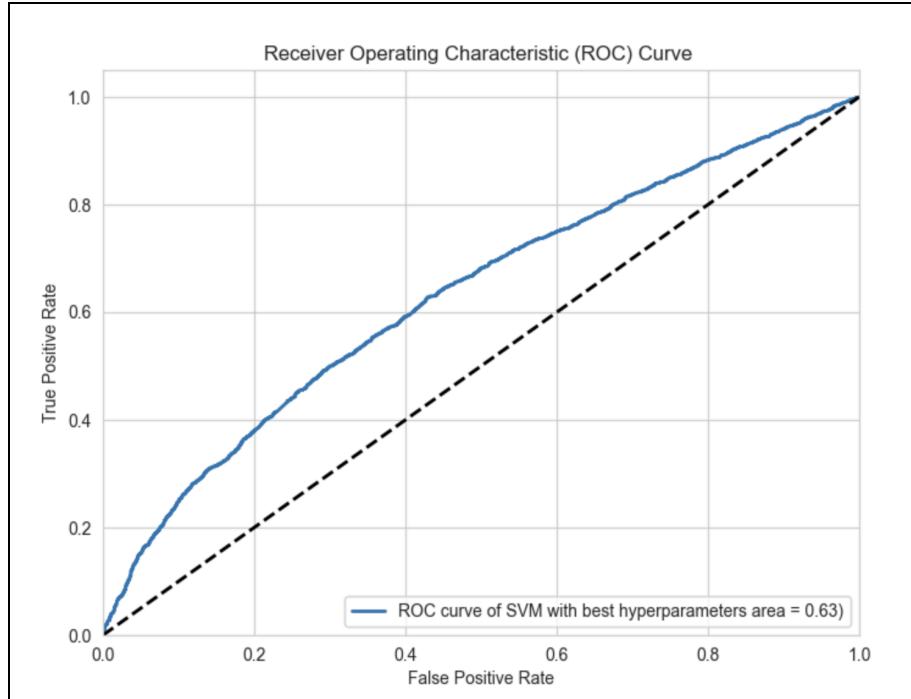
**Best Hyperparameters: {'kernel': 'rbf'}**

**Figure 50**

*Various metrics of SVM*

Classifier	Accuracy	Confusion Matrix	Precision	Sensitivity or Recall	Specificity	F-score	AUC
SVM	0.6	[[3018 772] [2069 1287]]	0.63	0.38	0.8	0.48	0.63

**Figure 51***Confusion matrix of SVM***Figure 52***ROC curve of SVM*



### **Naive Bayes:**

Naive Bayes is a widely used and straightforward probabilistic algorithm primarily designed for classification tasks. It relies on Bayes' theorem and operates under the “naive” assumption that predictors (features) are independent of each other. Since the algorithm has no tunable parameters, grid search is not applicable. Instead, a stratified k-fold cross-validation approach was used to calculate the test accuracy.

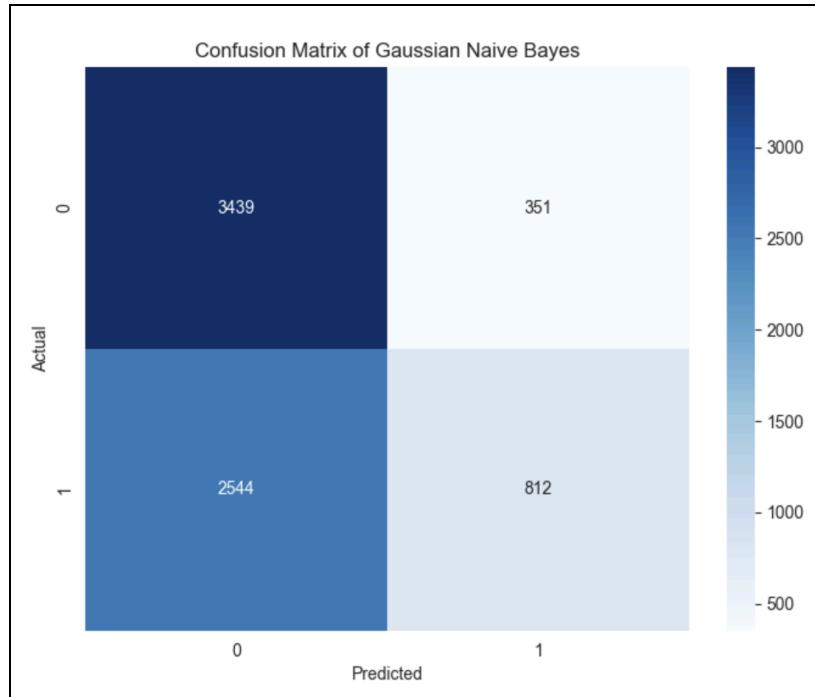
**Figure 53**

*Various metrics of Gaussian Naive Bayes*

Classifier	Accuracy	Confusion Matrix	Precision	Sensitivity or Recall	Specificity	F-score	AUC
Gaussian Naive Bayes	0.58	[[3439 351] [2544 812]]	0.7	0.24	0.91	0.36	0.62

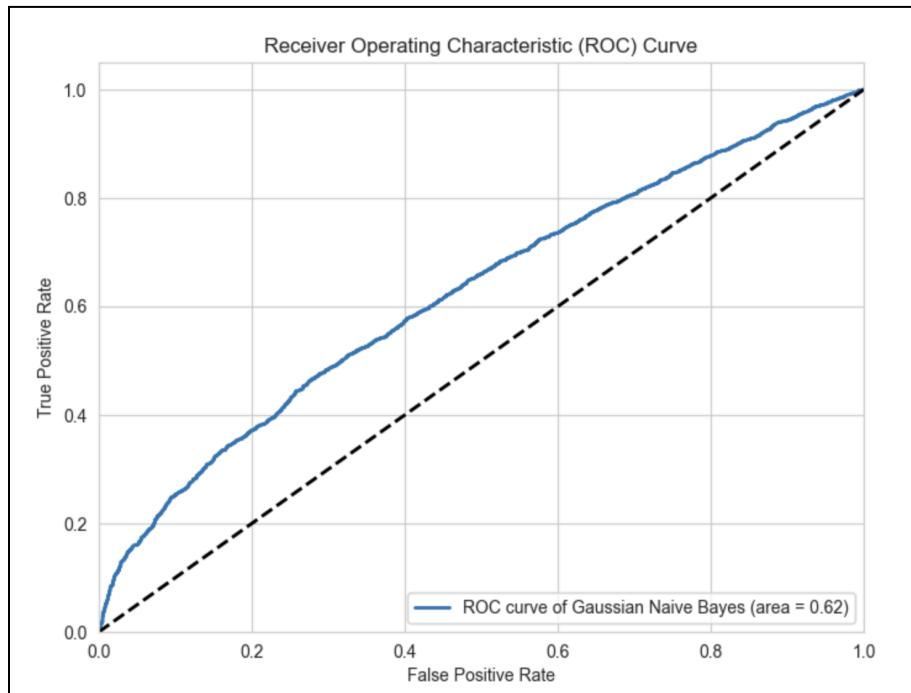
**Figure 54**

*Confusion matrix of Gaussian Naive Bayes*



**Figure 55**

*ROC curve of Gaussian Naive Bayes*



### ***Multi-layered Perceptron:***

A Multi-Layer Perceptron (MLP) is a form of artificial neural network comprising multiple layers of neurons, capable of capturing complex relationships between input and output data. As a feedforward neural network, it maps input data to corresponding outputs. Grid search and stratified k-fold cross-validation were utilized to determine the optimal hyperparameters. The model demonstrates reasonable performance in terms of accuracy, precision, specificity, and AUC, but its recall and F-score are relatively low.

**Figure 56**

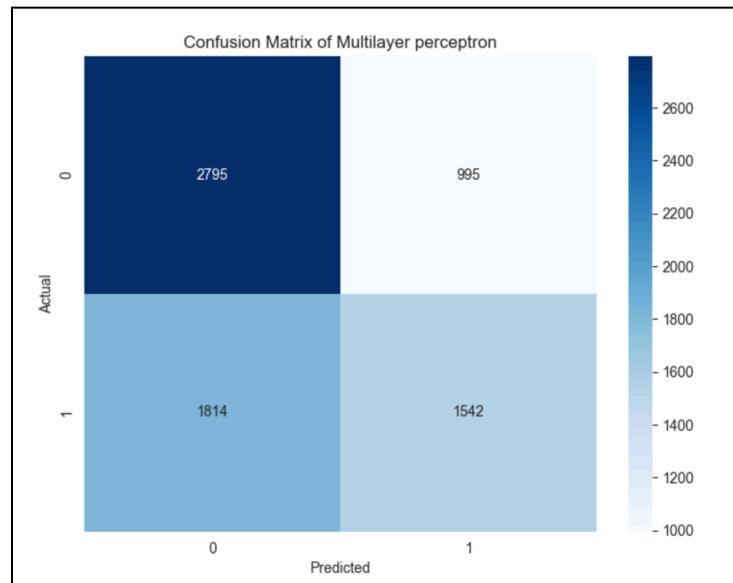
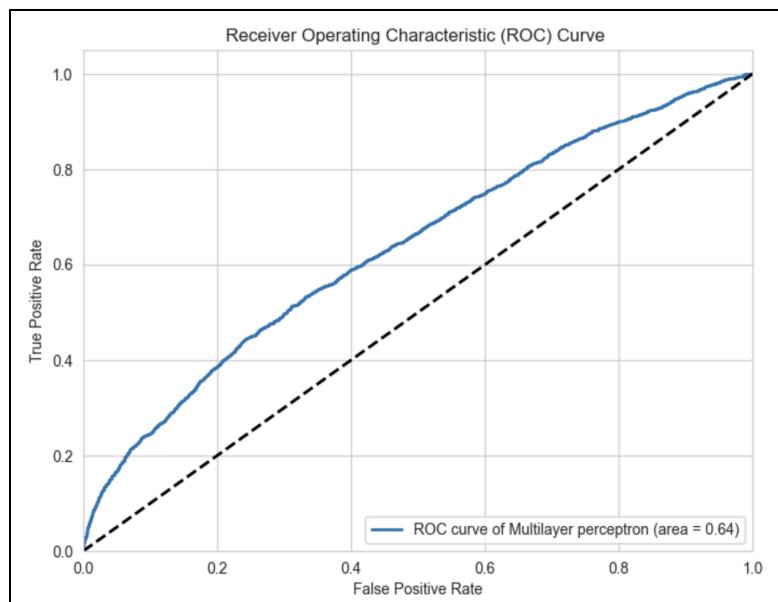
*Best hyperparameters of the Multi-layered perceptron*

**Best Hyperparameters: {'hidden\_layer\_sizes': (15, 15)}**

**Figure 57**

*Various metrics of Multi-layered perceptron*

Classifier	Accuracy	Confusion Matrix	Precision	Sensitivity or Recall	Specificity	F-score	AUC
Multilayer perceptron	0.61	[[2795 995] [1814 1542]]	0.61	0.46	0.74	0.52	0.64

**Figure 58***Confusion matrix of Multi-layered perceptron***Figure 59***ROC curve of Multi-layered perceptron*

### *Analysis of all classifiers of Phase 3:*

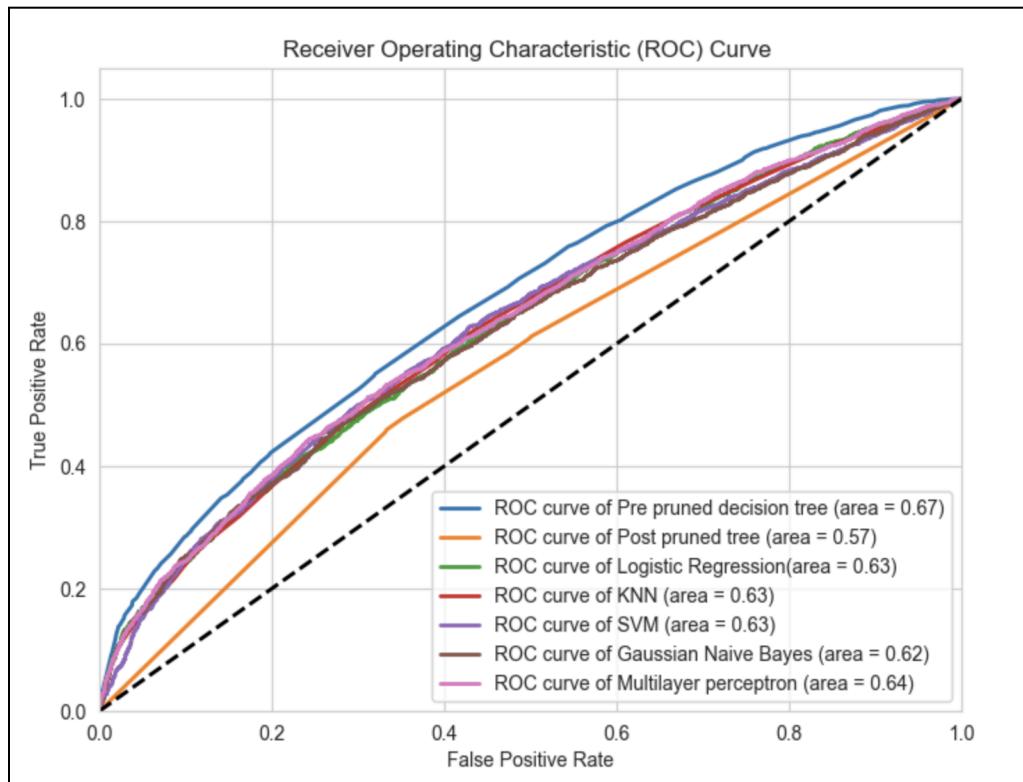
**Figure 60**

*Comparison of all metrics of all classifiers*

Classifier	Accuracy	Confusion Matrix	Precision	Sensitivity or Recall	Specificity	F-Score	AUC
Pre pruned Decision Tree Classifier	0.62	[[3047 743], [1956 1400]]	0.65	0.42	0.8	0.51	0.67
Post pruned Decision Tree Classifier	0.57	[[2456 1334], [1751 1605]]	0.55	0.48	0.8	0.51	0.57
Logistic Regression	0.6	[[2872 918], [1947 1409]]	0.61	0.42	0.76	0.5	0.63
KNN	0.6	[[2714 1076], [1781 1575]]	0.59	0.47	0.72	0.52	0.63
SVM	0.6	[[3018 772], [2069 1287]]	0.63	0.38	0.8	0.48	0.63
Gaussian Naive Bayes	0.58	[[3439 351], [2544 812]]	0.7	0.24	0.91	0.36	0.62
Multilayer Perceptron	0.61	[[2795 995], [1814 1542]]	0.61	0.46	0.74	0.52	0.64

**Figure 61**

*ROC curve of all classifiers*



Pre-pruned Decision Tree vs. Post-pruned Decision Tree: The pre-pruned decision tree has a slightly higher accuracy compared to the post-pruned one (0.62 vs. 0.57). However, the post-pruned tree shows better sensitivity/recall, indicating it's better at correctly identifying positive cases (delayed flights) after pruning. Post-pruning slightly improves the model's performance in correctly identifying delayed flights but results in a decrease in overall accuracy.

Logistic Regression, KNN, SVM, Gaussian Naive Bayes, and MLP: These models demonstrate similar accuracy scores around 0.58 to 0.61. KNN exhibits slightly better sensitivity (0.47) compared to others. Logistic Regression and SVM show decent precision (around 0.61 to 0.63) but moderate sensitivity (around 0.38 to 0.42).

**Interpretation:** Overall, ***The Pre-Pruned Decision Tree emerges as the best model*** due to its balanced performance across key metrics. It achieves the highest accuracy of 0.62, which is crucial for overall prediction reliability, and the highest AUC (0.67), reflecting its ability to distinguish between delayed and non-delayed flights effectively. Additionally, it exhibits strong precision (0.65), meaning it minimizes false positives while maintaining a reasonable sensitivity of 0.42. The model's high specificity (0.80) indicates that it is highly effective in identifying non-delayed flights. Furthermore, its simplicity and interpretability make it an excellent choice for operational settings, such as in the airline industry, where decisions need to be both accurate and explainable.

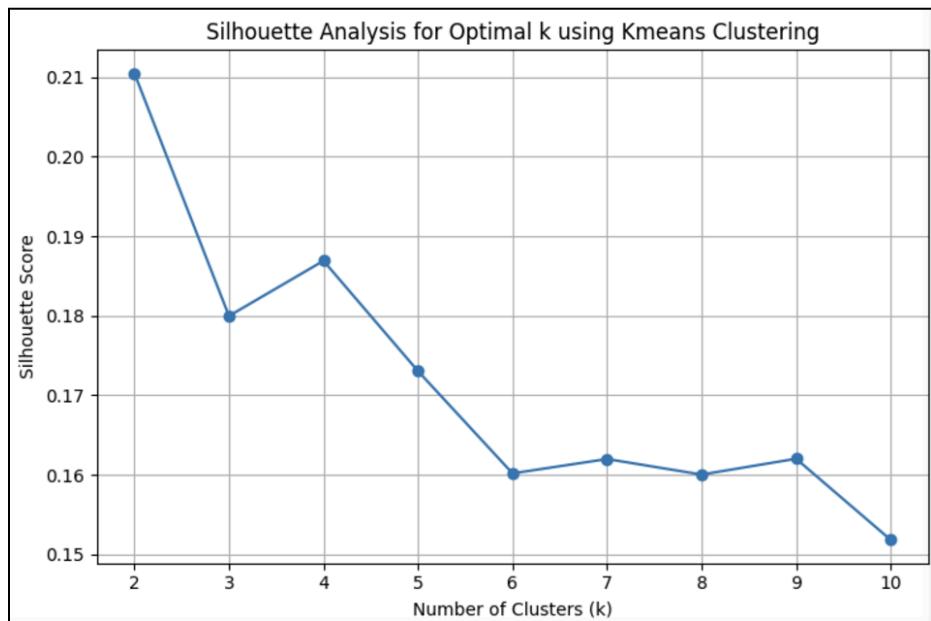
## Phase IV: Clustering and Association

### KMeans Clustering:

K-means clustering is a popular unsupervised machine learning algorithm used for clustering similar data points into groups or clusters. Silhouette Score measures how well-separated the clusters are. Silhouette analysis shows a higher score for  $k = 2$ , but this might happen because the algorithm is combining what should be multiple clusters into fewer larger clusters. A higher silhouette score does not always mean the best clustering if it oversimplifies the data. While the silhouette score for  $k = 2$  is higher, the WCSS elbow at  $k = 4$  suggests that the dataset likely has four natural groupings. The silhouette score at  $k = 4$  is still acceptable, and the drop in silhouette score after  $k = 4$  further validates this choice. Therefore, **optimal value k is 4**.

**Figure 62**

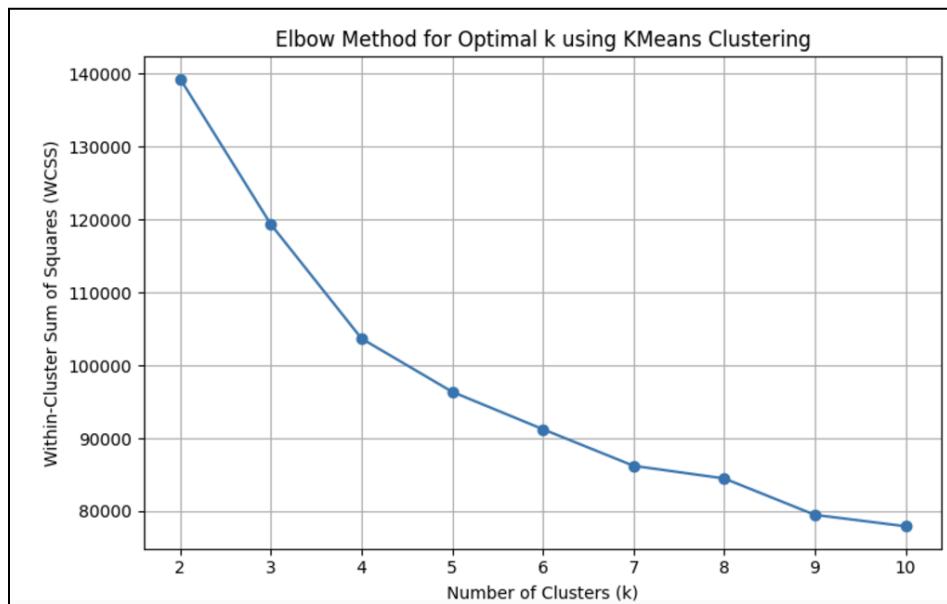
*Silhouette analysis for Optimal k using KMeans clustering*



The silhouette score reaches its highest value at  $k = 2$ , suggesting the data can be grouped into two broad clusters. However, another significant peak is observed at  $k = 4$ , which also represents a well-separated clustering structure. After  $k = 4$ , the silhouette score consistently decreases, indicating that further increasing the number of clusters does not meaningfully improve the clustering quality. This choice of  $K = 4$  is also supported by the elbow plot that is plotted and discussed below.

### Figure 63

*WCSS plot analysis for optimal k using KMeans clustering*



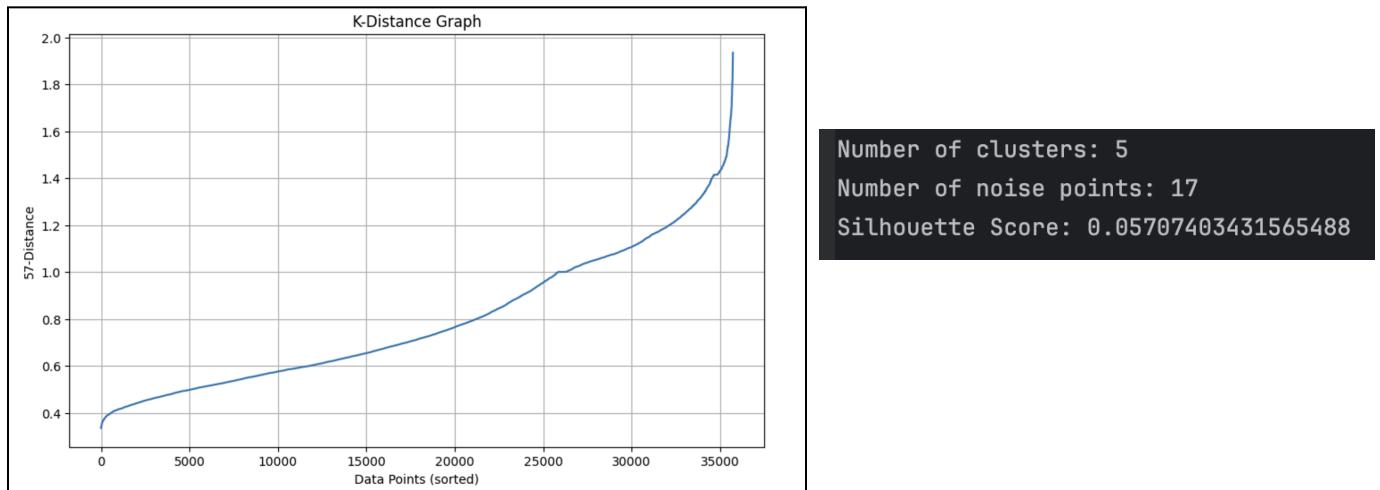
Within-Cluster Sum of Squares (WCSS): Measures the compactness of clusters. Lower WCSS implies more compact clusters. The elbow point is where the rate of decrease in within-cluster variation slows down significantly, forming an “elbow” shape. In this graph, the elbow appears to occur around  $k = 4$ , as the reduction in within-cluster variation beyond this point becomes more gradual. Thus,  $k = 4$  is likely the optimal number of clusters based on this plot.

### **DBSCAN algorithm:**

DBSCAN (Density-Based Spatial Clustering of Applications with Noise) is a density-based clustering algorithm that identifies clusters of high-density points while labeling sparser points as noise. It does not require specifying the number of clusters in advance, making it ideal for datasets with arbitrarily shaped clusters. DBSCAN relies on two key parameters: `eps` (radius of neighborhood) and `min_samples` (minimum points to form a dense cluster). These are determined using the K-Distance graph, where the “elbow point” indicates a suitable `eps`.

**Figure 64**

*K-distance graph for optimal eps for dbscan algorithm*



In this analysis, the K-Distance graph suggested an `eps` of approximately 0.5 with `min_samples` = 5. DBSCAN identified 5 clusters and 17 noise points, achieving a Silhouette Score of 0.057, which indicates weak separation between clusters. This low score might result from overlapping clusters or varying densities within the data. DBSCAN effectively identifies noise and handles non-spherical clusters but can be sensitive to parameter selection and struggles with datasets having highly variable cluster densities.

### ***Apriori algorithm:***

The Apriori algorithm is a classic algorithm used in association rule mining, particularly for discovering frequent item sets in transactional databases or datasets. The support for flights labeled as "Large" is 18.97%, while for "Short" lengths, it's 41.42%. There's an interesting relationship here between flights labeled as "Small" and specific conditions related to both "Medium" length and the "Start of the Week" for the day of the week

**Figure 65**

*Rules generated from Apriori*

Processing 4 combinations   Sampling itemset size 4											
	antecedents	consequents	antecedent support	consequent support	support	confidence	lift	leverage	conviction	zhangs_metric	
36	(Large Aircraft)	(Short Length)	0.18	0.45	0.12	0.65	1.45	0.04	1.59	0.38	
64	(Average Length, StartOfWeek)	(Small Aircraft)	0.18	0.51	0.10	0.56	1.10	0.01	1.12	0.12	
9	(Average Length)	(Small Aircraft)	0.43	0.51	0.24	0.56	1.10	0.02	1.12	0.17	
3	(Average Length)	(Early Morning)	0.43	0.52	0.23	0.55	1.05	0.01	1.06	0.09	
58	(Early Morning, Average Length)	(Small Aircraft)	0.23	0.51	0.13	0.54	1.08	0.01	1.08	0.09	

The presence of "Large Aircraft" (18% of occurrences) is associated with "Short Length" (45% of occurrences) in 12% of cases. The confidence is moderate at 65%, indicating that when there's a large aircraft, there's a 65% chance that the flight will be of short length.

The lift value of 1.45 implies that the occurrence of a large aircraft increases the chance of a short-length flight by 45% compared to what would be expected if the two were independent.

When flights start at the beginning of the week and have an average length, there's a 56% chance they'll involve small aircraft. The lift value of 1.10 suggests a marginal increase in the likelihood of small aircraft in these situations. Flights with an average length and scheduled in the early morning have a 55% chance of occurrence. The lift value is close to 1, indicating no significant association beyond what would be expected by chance. When flights are early in the morning, of average length, there's a 54% chance they'll involve small aircraft. Both the lift and confidence are close to 1, suggesting a mild association without a significant deviation from chance.

## Recommendations

This project involved the practical application of diverse machine learning techniques on an airline dataset to predict flight delays. Throughout the process, I gained valuable experience in implementing various models, fine-tuning hyperparameters, and identifying optimal classifiers. This hands-on exploration not only strengthened my ability to work with real-world datasets but also highlighted the importance of efficient training and hyperparameter optimization in achieving robust model performance.

Looking ahead, I plan to enhance this project by integrating the dataset with additional sources, such as flight schedules, to provide more contextual information. I also aim to experiment with advanced models like XGBoost to further improve prediction accuracy and performance.

Among the evaluated models, ***The Pre-Pruned Decision Tree emerges as the best model*** due to its balanced performance across key metrics. It achieves the highest accuracy of 0.62, which is crucial for overall prediction reliability, and the highest AUC (0.67). In terms of feature importance, attributes such as Time, Length, Flight, and Airline\_WN showed a strong relationship with the target variable, Delay. The prominence of Airline\_WN, which represents the busiest airline in the dataset, likely contributes to higher delays. Additionally, the type of aircraft used also appeared to influence flight delays significantly.

Finally, the KMeans clustering analysis identified four clusters as optimal. While  $k = 2$  achieves the highest silhouette score,  $k = 4$  aligns better with the Elbow Method and provides a more detailed segmentation of the data. This suggests that the dataset contains moderate differentiation into four meaningful subgroups, despite some overlapping characteristics.

## Appendix

Importing necessary libraries:

```

1 %%-----Phase I: Feature Engineering & EDA-----
2 import pandas as pd
3 import seaborn as sns
4 import matplotlib.pyplot as plt
5 import numpy as np
6 from sklearn.preprocessing import StandardScaler
7 from sklearn.model_selection import train_test_split, StratifiedKFold, GridSearchCV
8 from sklearn.ensemble import RandomForestClassifier, StackingClassifier
9 from sklearn.decomposition import PCA, TruncatedSVD
10 from statsmodels.stats.outliers_influence import variance_inflation_factor
11 from sklearn.feature_selection import f_classif
12 from prettytable import PrettyTable
13 import statsmodels.api as sm
14 from sklearn.linear_model import LinearRegression, LogisticRegression
15 from sklearn.metrics import mean_squared_error, accuracy_score, confusion_matrix, roc_auc_score, roc_curve
16 from sklearn.tree import DecisionTreeClassifier
17 from sklearn import tree
18 from sklearn.neighbors import KNeighborsClassifier
19 from sklearn.svm import SVC
20 from sklearn.naive_bayes import GaussianNB
21 from sklearn.neural_network import MLPClassifier
22 from mlxtend.preprocessing import TransactionEncoder
23 from mlxtend.frequent_patterns import apriori, association_rules

```

Data Cleaning/ Check for data duplications and removal:

```

23 #Data cleaning
24 df = pd.read_csv(r'/Users/jyothi/Documents/SEM2/Machine_Learning/Project/Airlines.csv')
25 print(df.columns)
26 print(df.head().to_string())
27 print("\nMissing values in the dataset")
28 print(df.isna().sum().sum() + df.isnull().sum().sum())
29
30 #Check for data duplications and removal
31 print("\nDuplicated values in the dataset")
32 print(df.duplicated().sum())
33 print("\nNumber of unique values in each feature of whole dataset")
34 print(df.nunique())

```

Aggregation and Downsampling:

```

36 #Aggregation
37 df.dayofweek = df[['DayOfWeek', 'Delay']].groupby('DayOfWeek')['Delay'].sum().reset_index().sort_values(by='Delay', ascending=False)
38 print(df.dayofweek.to_string())
39 df.airline = df[['Airline', 'Delay']].groupby('Airline')['Delay'].sum().reset_index().sort_values(by='Delay', ascending=False)
40 print(df.airline.to_string())
41
42 #Down sampling
43 df = df.drop(['id', 'AirportTo'], axis=1)
44 df.airline = df['Airline'].value_counts().sort_values(ascending=False)
45 df = df[df['Airline'].isin(['WN', 'DL', 'OO', 'AA', 'MQ'])]
46 df.airport_to = df['AirportFrom'].value_counts().sort_values(ascending=False)
47 df = df[df['AirportFrom'].isin(['ATL', 'ORD', 'DFW', 'DEN', 'LAX'])]
48 df = df[::2]
49 print("\nThe shape of the dataset after downsampling:")
50 print(df.shape)

```

Dimensionality reduction/feature selection:

Using Random Forest Analysis:

```

104     # Dimensionality reduction
105     # Feature importance using Random forest analysis
106     model_rf = RandomForestClassifier(random_state=5805)
107     model_rf.fit(X, y)
108     importances = model_rf.feature_importances_
109     print("Feature importances from Random forest method")
110     indices = np.argsort(importances)
111     sortedImportance = importances[indices] * 100
112     print(sortedImportance[::-1].round(2))
113     sorted_features = X.columns[indices]
114     plt.figure(figsize=(15, 10))
115     plt.barh(range(len(sortedImportance)), sortedImportance, color='darkblue') # Light blue hex code
116     plt.yticks(range(len(sorted_features)), sorted_features)
117     plt.title('Feature importances by RandomForestClassifier')
118     plt.xlabel('Feature Importance')
119     plt.ylabel('Features')
120     plt.show()

```

Using Principal Component Analysis:

```

122     # Dimensionality reduction/Feature selection
123     # Using PCA
124     pca = PCA()
125     pca.fit(X)
126     cum_var = np.round(np.cumsum(sorted(pca.explained_variance_ratio_, reverse=True)) * 100, decimals=2)
127     print("Cumulative explained variance of each component from PCA")
128     print(cum_var)
129     labels = [i for i in range(1, len(cum_var) + 1)]
130     plt.plot(*args: range(len(cum_var)), cum_var, alpha=0.9, color='darkblue') # Dark blue line
131     x_point = 11
132     y_point = 95
133     plt.axvline(x=x_point, color='darkblue', linestyle='--')
134     plt.axhline(y=y_point, color='darkblue', linestyle='--')
135     plt.annotate(text=f'({x_point},{y_point})', xy=(x_point, y_point), color='black', textcoords="offset points", xytext=(-15, 10), ha='center')
136     plt.xlabel(label: "Number of components", color='black')
137     plt.ylabel(label: "Percentage contribution to variance", color='black')
138     plt.title(label: "Cumulative explained variance of each component from PCA", color='black')
139     plt.tick_params(axis='x', colors='black')
140     plt.tick_params(axis='y', colors='black')
141     plt.show()

```

Using Singular Value Decomposition Analysis:

```

143     # Dimensionality reduction/Feature selection
144     # Using Singular Value Decomposition
145     tsvd = TruncatedSVD(n_components=17)
146     tsvd_result = tsvd.fit(X)
147     print("Explained variance ratio from SVD ", tsvd.explained_variance_ratio_.round(2))
148     print("Singular values of features {tsvd.singular_values_.round(2)}")
149     cum_evr = np.cumsum(100 * tsvd.explained_variance_ratio_)
150     labels = [i for i in range(1, len(cum_evr) + 1)]
151     plt.bar(*labels, height=cum_evr, alpha=0.8, color='darkblue')
152     x_point = 12
153     y_point = 95
154     plt.axvline(x=x_point, linestyle='--', color='red')
155     plt.axhline(y=y_point, linestyle='--', color='red')
156     plt.annotate(text=f'({x_point},{y_point})', xy=(x_point, y_point), textcoords="offset points", xytext=(-15, 10), color='black')
157     plt.xlabel(label: "Number of components", color='black')
158     plt.ylabel(label: "Percentage contribution to variance", color='black')
159     plt.title(label: "Cumulative variance ratio of each component from SVD", color='black')
160     plt.xticks(labels)
161     plt.show()

```

Using VIF and condition number:

```

164     # Dimensionality reduction/Feature selection
165     # Using Variance Inflation Factor
166     print("VIF analysis")
167     vif_fea = [variance_inflation_factor(X.values, i) for i in range(X.shape[1])]
168     for i,j in zip(df.columns, vif_fea):
169         print(i,round(j,2))
170     _, S, _ = np.linalg.svd(X, full_matrices=False)
171     print("Singular values of the dataset")
172     print(S.round(2))
173     condition_number = np.linalg.cond(X_train)
174     print("Condition number : ", round(condition_number,2))
175
176     # Data imbalance analysis
177     sns.set_style('whitegrid')
178     sns.countplot(data=df,x='Delay')
179     plt.title('Countplot of # of obs. for each class of Delay')
180     plt.xlabel("Whether flight delayed")
181     plt.ylabel("# of samples")
182     plt.show()

```

Discretization & Binarization: Label Encoding/one hot encoding :

```

67     #Discretization & Binarization: Label Encoding/one hot encoding
68     one_hot_enc_col = ['AirportFrom', 'Airline', 'DayOfWeek']
69     df = pd.get_dummies(df, columns=one_hot_enc_col, drop_first=True)
70     for i in df.columns:
71         if i not in ['Length', 'Time', 'Delay', 'Flight']:
72             df[i] = df[i].map({False:0, True:1})
73     print("\nOne hot encoded dataset")
74     print(df.head().to_string())
75     df_orig= df[:,:]
76     X_orig = df.drop(labels='Delay', axis=1)
77     X_reg = df[:,:]

```

Variable Transformation: Standardization:

```

93     #Variable Transformation: Standardization
94     sc = StandardScaler()
95     for i in ['Flight', 'Length', 'Time']:
96         df[i] = sc.fit_transform(df[[i]])
97     print("\nStandardized dataset")
98     print(df.head().to_string())
99     X = df.drop( labels='Delay', axis=1)
100    y = df['Delay']
101    print(X.shape)
102    X_train, X_test, y_train, y_test = train_test_split(*arrays: X, y, test_size=0.3,random_state=5805, shuffle=True,stratify=y)

```

## Outlier Analysis and Removal:

```

55 #Anomaly detection/Outlier Analysis and removal using IQR method
56 q1 = df['Length'].quantile(0.25)
57 q3 = df['Length'].quantile(0.75)
58 IQR = q3 - q1
59 lower_bound = q1 - 1.5 * IQR
60 upper_bound = q3 + 1.5 * IQR
61 df = df[(df['Length'] < upper_bound) & (df['Length'] > lower_bound)]
62 df_new = df[::]
63 sns.boxplot(data=df, x='Length')
64 plt.title("Box plot of Length feature after removal of outliers using IQR method")
65 plt.show()

```

## Calculating and plotting of correlation matrix:

```

79 # Correlation matrix using heatmap
80 correlation = df.corr().round(2)
81 plt.figure(figsize = (18,11))
82 sns.heatmap(correlation, annot = True, cbar=True)
83 plt.title('Heatmap of correlation between all the features')
84 plt.show()

```

## Calculating and plotting of covariance matrix:

```

86 # Covariance matrix using heatmap
87 correlation = df.cov().round(2)
88 plt.figure(figsize = (18,11))
89 sns.heatmap(correlation, annot = True, cbar=True)
90 plt.title('Heatmap of covariance between all the features')
91 plt.show()

```

## T-test analysis:

```

247 #T-test analysis
248 pt_ttest = PrettyTable()
249 pt_ttest.field_names = ['Feature', 't-statistic', 'p-value']
250 print("t-test analysis")
251 for i in range(len(X_reg_train.columns)):
252     pt_ttest.add_row([X_reg_train.columns[i], round(model_ols.tvalues[i],2), round(model_ols.pvalues[i],2)])
253 print(pt_ttest)

```

F-test analysis:

```

198     #F-test analysis
199     f_test, p_values = f_classif(X_reg, y_reg)
200     pt = PrettyTable()
201     pt.field_names = ['Feature', 'F-test score', 'p-value']
202     for i in range(len(f_test)):
203         pt.add_row([X.columns[i], round(f_test[i],2), round(p_values[i],2)])
204     print(pt)

```

Linear regression:

```

205     # Prediction with Linear Regression
206     model_linear = LinearRegression()
207     X_reg_train = sm.add_constant(X_reg_train)
208     X_reg_test = sm.add_constant(X_reg_test)
209     model_linear.fit(X_reg_train, y_reg_train)
210     y_pred_train = model_linear.predict(X_reg_train)
211     y_pred_train_orig = y_pred_train*std_length+mean_length
212     y_pred_test = model_linear.predict(X_reg_test)
213     y_pred_test_orig = y_pred_test*std_length+mean_length
214     print("The coefficients of the Linear Regression model")
215     print("Feature", "Coefficient")
216     for i in range(len(X_reg_train.columns)):
217         print(X_reg_train.columns[i], round(model_linear.coef_[0][i],2))
218     mse = mean_squared_error(y_reg_train*std_length+mean_length, y_pred_train_orig)
219     print("Mean Squared Error for training set (MSE): {mse:.3f}")
220     pt=PrettyTable()
221     pt.field_names=['Model', 'R-squared', 'Adj. R-squared', 'AIC', 'BIC', 'Mean Squared Error']
222     r_squared = model_linear.score(X_reg_train, y_reg_train)
223     n = len(y_reg_train)
224     p = X_reg_train.shape[1] - 1
225     adjusted_r_squared = 1 - (1 - r_squared) * (n - 1) / (n - p - 1)
226     aic = n*np.log(mse) + 2*p
227     bic = n*np.log(mse) + p*np.log(n)
228     pt.add_row(['Linear Regression', round(r_squared,2), round(adjusted_r_squared,2),
229     round(aic,2), round(bic,2), round(mse,2)])
230     print(pt)
231     index=np.arange(len(y_reg_test))
232     plt.title("Prediction with Linear Regression")
233     plt.plot(*args: index, y_reg_test * std_length + mean_length, color='blue', label='Actual values of test set') # Blue for actual values
234     plt.plot(*args: index, y_pred_test_orig, color='lightyellow', label='Predicted values with test set') # Yellow for predicted values
235     plt.xlabel("# of samples")
236     plt.ylabel("Length")
237     plt.legend()
238     plt.show()

```

Confidence Interval analysis:

```

367 #Confidence interval analysis
368 index=np.arange(len(y_reg_test))
369 y_prediction = model_ols.get_prediction(X_reg_test)
370 conf_intr = model_ols.conf_int()
371 ci_pred_frame=y_prediction.summary_frame(alpha=0.05)
372 ci_upper = ci_pred_frame.obs_ci_upper
373 ci_lower = ci_pred_frame.obs_ci_lower
374 plt.plot(*args: index,y_pred_test_orig,label="Predicted",color='blue')
375 plt.fill_between(index, ci_upper*std_length+mean_length,
376 ci_lower*std_length+mean_length, color='blue', alpha=0.3, label="Confidence Interval")
377 plt.xlabel("Number of observations")
378 plt.ylabel("Sales")
379 plt.title("Confidence interval of Predicted values summary_frame")
380 plt.legend()
381 plt.show()
382 model_ols_sr = sm.OLS(y_reg_train, X_reg_train).fit()
383 print(model_ols_sr.summary())

```

Backward stepwise regression:

```

255 ## Backward stepwise regression
256 print("Eliminating the feature with highest p-value which is DayOfWeek_3")
257 X_reg_train = X_reg_train.drop('DayOfWeek_3', axis=1)
258 X_reg_test = X_reg_test.drop('DayOfWeek_3', axis=1)
259 model_ols = sm.OLS(y_reg_train, X_reg_train).fit()
260 y_pred_train = model_ols.predict(X_reg_train)
261 y_pred_train_orig = y_pred_train*std_length+mean_length
262 y_pred_test = model_ols.predict(X_reg_test)
263 y_pred_test_orig = y_pred_test*std_length+mean_length
264 print(model_ols.summary())
265
266 print("Eliminating DayOfWeek_4")
267 X_reg_train = X_reg_train.drop('DayOfWeek_4', axis=1)
268 X_reg_test = X_reg_test.drop('DayOfWeek_4', axis=1)
269 model_ols = sm.OLS(y_reg_train, X_reg_train).fit()
270 y_pred_train = model_ols.predict(X_reg_train)
271 y_pred_train_orig = y_pred_train*std_length+mean_length
272 y_pred_test = model_ols.predict(X_reg_test)
273 y_pred_test_orig = y_pred_test*std_length+mean_length
274 print(model_ols.summary())

```

Pre-pruned decision tree:

```

386 ## Pre-prune decision tree
387 X = df.drop( labels: 'Delay', axis=1)
388 print(X.columns)
389 y = df['Delay']
390 X_train, X_test, y_train, y_test = train_test_split(*arrays: X, y, test_size=0.2,random_state=5805, shuffle=True)
391 pt = PrettyTable()
392 pt.field_names = ['Classifier', 'Accuracy', 'Confusion Matrix', 'Precision','Sensitivity or Recall', 'Specificity', 'F-score', 'AUC']
393 clf=DecisionTreeClassifier(random_state=5805)
394 tuned_parameters = [{'max_depth':[5,7,9,11],
395 'min_samples_split': [2,5,10,15],
396 'min_samples_leaf':[1,2,4,8],
397 'max_features':[2,4,7],
398 'splitter':['best', 'random'],
399 'criterion':['gini','entropy','log_loss']}]
400 stratified_kfold = StratifiedKFold(n_splits=5, shuffle=True, random_state=5805)
401 grid_search = GridSearchCV(clf, tuned_parameters, cv=stratified_kfold,scoring='accuracy')
402 grid_search.fit(X, y)
403 print("Best Hyperparameters:", grid_search.best_params_)
404 best_classifier = grid_search.best_estimator_
405 y_train_pred = best_classifier.predict(X_train)
406 y_test_pred = best_classifier.predict(X_test)
407 pre_acc = accuracy_score(y_test, y_test_pred)
408 print(f'Train accuracy of pre pruned decision tree {round(accuracy_score(y_test,y_test_pred),2)}')
409 print(f'Test accuracy of pre pruned decision tree {round(pre_acc, 2)}')
410 conf_matrix_pre = confusion_matrix(y_test, y_test_pred)
411 tn, fp, fn, tp = conf_matrix_pre.ravel()
412 specificity_pre = tn/(tn+fp)
413 y_prob_pre = best_classifier.predict_proba(X_test)[:, 1] # Probability of the positive class
414 roc_auc_pre = roc_auc_score(y_test, y_prob_pre)
415 recall_pre = recall_score(y_test, y_test_pred)
416 precision_pre = precision_score(y_test, y_test_pred)
417 f1_pre = f1_score(y_test, y_test_pred)
418 # heatmap of confusion matrix
419 plt.figure(figsize=(8, 6))

```

Post-pruned decision tree:

```

444 ## Post pruned decision tree
445 fold_accuracies = []
446 fold_train_accuracies = []
447 path = best_classifier.cost_complexity_pruning_path(X_train,y_train)
448 alphas = path['ccp_alphas']
449
450 # Iterate over different values of ccp_alpha
451 for alpha in alphas:
452     clf = DecisionTreeClassifier(random_state=5805, ccp_alpha=alpha)
453     clf.fit(X_train, y_train)
454     accuracy_test = clf.score(X_test, y_test)
455     fold_accuracies.append(accuracy_test)
456     accuracy_train = clf.score(X_train, y_train)
457     fold_train_accuracies.append(accuracy_train)
458
459 best_alpha_index = np.argmax(fold_accuracies)
460 best_alpha = alphas[best_alpha_index]
461 print(f"Best alpha value of post pruned decision tree: {best_alpha}")
462
463 # Plot accuracy vs. alpha for train and test sets
464 fig, ax = plt.subplots()
465 ax.set_xlabel('alpha')
466 ax.set_ylabel('accuracy')
467 ax.set_title("Accuracy vs alpha for training and testing sets")
468 ax.plot(*args: alphas, fold_train_accuracies, marker="o", label="Train Accuracy", drawstyle="steps-post")
469 ax.plot(*args: alphas, fold_accuracies, marker="o", label="Test Accuracy", drawstyle="steps-post")
470 ax.legend()
471 plt.grid()
472 plt.tight_layout()
473 plt.show()
474 best_alpha = alphas[np.argmax(accuracy_test)]
475 print("Best alpha value of post pruned decision tree", best_alpha)
476
477 clf = DecisionTreeClassifier(random_state=5805, ccp_alpha=best_alpha)
478 clf.fit(X_train, y_train)
479 y_train_pred = clf.predict(X_train)

```

Logistic regression:

```

526 ## Logistic Regression
527 param_grid = {
528     'C': [0.1, 1.0, 10.0],
529     'penalty': ['l1', 'l2'],
530     'solver': ['liblinear', 'saga']
531 }
532 model_log = LogisticRegression(random_state=5805, max_iter=400)
533 grid_search = GridSearchCV(model_log, param_grid, cv=stratified_kfold,
534 scoring='accuracy')
535 grid_search.fit(X, y)
536 print("Best Hyperparameters:", grid_search.best_params_)
537 model_log = grid_search.best_estimator_
538 model_log.fit(X_train, y_train)
539 y_train_pred = model_log.predict(X_train)
540 y_test_pred = model_log.predict(X_test)
541 log_acc_train = accuracy_score(y_train, y_train_pred)
542 log_acc_test = accuracy_score(y_test, y_test_pred)
543 print(f'Train accuracy of Logistic Regression {round(accuracy_score(y_train,y_train_pred),2)}')
544 print(f'Test accuracy of Logistic Regression {round(log_acc_test,2)}')
545 conf_matrix_log = confusion_matrix(y_test, y_test_pred)
546 tn, fp, fn, tp = conf_matrix_log.ravel()
547 specificity_log = tn/(tn+fp)
548 y_prob_log = model_log.predict_proba(X_test)[:, 1]
549 roc_auc_log = roc_auc_score(y_test, y_prob_log)
550 recall_log = recall_score(y_test, y_test_pred)
551 precision_log = precision_score(y_test, y_test_pred)
552 f1_log = f1_score(y_test, y_test_pred)
553 pt2 = PrettyTable()
554 pt2.field_names = ['Classifier', 'Accuracy', 'Confusion Matrix', 'Precision',
555 'Sensitivity or Recall', 'Specificity', 'F-score', 'AUC']
556 pt2.add_row(['Logistic Regression', round(np.mean(log_acc_train),2),

```

KNN classifier:

```

583 ## KNN Classifier
584 param_grid = {'n_neighbors': [6, 9, 12, 15, 17, 19, 25, 29, 33, 37, 41, 43, 45, 49, 53, 57,61],
585 }
586 model_knn = KNeighborsClassifier()
587 grid_search = GridSearchCV(model_knn, param_grid, cv=stratified_kfold, scoring='accuracy')
588 grid_search.fit(X, y)
589 results = grid_search.cv_results_
590 k_values = results['param_n_neighbors'].data
591 mean_test_scores = results['mean_test_score']
592 error_rate = 1-mean_test_scores
593 plt.plot(*args: k_values, error_rate)
594 plt.xlabel('K')
595 plt.ylabel('Error rate')
596 plt.title('Error rate vs K')
597 plt.show()
598 print("Best Hyperparameters:", grid_search.best_params_)

```

SVM classifier:

```

646 ## SVM Classifier
647 parameters = {
648 'kernel': ['linear', 'rbf', 'poly'], # Kernel types to try
649 }
650 # Perform grid search with cross-validation
651 svm_classifier = SVC(probability=True)
652 grid_search = GridSearchCV(svm_classifier, parameters, cv=stratified_kfold,n_jobs=-1)
653 grid_search.fit(X, y)
654 print("Best Hyperparameters:", grid_search.best_params_)

655
656 svm_lr = grid_search.best_estimator_
657 svm_lr.fit(X_train, y_train)
658 y_train_pred = svm_lr.predict(X_train)
659 y_test_pred = svm_lr.predict(X_test)
660 svm_lr_acc = accuracy_score(y_test, y_test_pred)
661 print(f'Train accuracy of SVM with best hyperparameters {round(np.mean(svm_lr_acc),2)}')
662 print(f'Test accuracy of SVM with best hyperparameters {round(accuracy_score(y_test, y_test_pred),2)}')

```

Gaussian Naive Bayes:

```

701 ## Naive Bayes Classifier
702 naive_bayes = GaussianNB()
703 train_scores = []
704 test_scores = []
705 for train_index, test_index in stratified_kfold.split(X, y):
706     X_train_cv, X_test_cv = X.iloc[train_index], X.iloc[test_index]
707     y_train_cv, y_test_cv = y.iloc[train_index], y.iloc[test_index]
708     naive_bayes.fit(X_train_cv, y_train_cv)
709     train_score = naive_bayes.score(X_train_cv, y_train_cv)
710     train_scores.append(train_score)
711     test_score = naive_bayes.score(X_test_cv, y_test_cv)
712     test_scores.append(test_score)
713     print(f'Train accuracy of Gaussian Naive Bayes {round(np.mean(train_scores),2)}')
714     naive_bayes.fit(X_train, y_train)
715     y_train_pred = naive_bayes.predict(X_train)
716     y_test_pred = naive_bayes.predict(X_test)
717     print(f'Test accuracy of Gaussian Naive Bayes {round(np.mean(test_scores),2)}')

```

Multilayer perceptron classifier:

```

755 ## MLP Classifier
756 mlp = MLPClassifier(max_iter=400, random_state=5805)
757 param_grid = {'hidden_layer_sizes': [(5, 5), (5, 10), (10, 10), (15, 15), (20, 20)]}
758 grid_search = GridSearchCV(mlp, param_grid, cv=stratified_kfold, scoring='accuracy')
759 grid_search.fit(X, y)
760 print("Best Hyperparameters:", grid_search.best_params_)

761
762 mlp = grid_search.best_estimator_
763 mlp.fit(X_train, y_train)
764 y_train_pred = mlp.predict(X_train)
765 y_test_pred = mlp.predict(X_test)
766 mlp_acc = accuracy_score(y_test, y_test_pred)
767 print(f'Train accuracy of Multilayer Perceptron with best hyperparameters {round(np.mean(mlp_acc),2)}')
768 print(f'Test accuracy of Multilayer Perceptron with best hyperparameters {round(accuracy_score(y_test, y_test_pred),2)}')

```

K-Means using silhouette analysis:

```

869 #K-Means Clustering
870 import matplotlib.pyplot as plt
871 from sklearn.cluster import KMeans
872 from sklearn.metrics import silhouette_score
873 k_values = range(2, 11)
874 wcss = []
875 silhouette_scores = []
876 for k in k_values:
877     kmeans = KMeans(n_clusters=k, random_state=5085)
878     kmeans.fit(X)
879     wcss.append(kmeans.inertia_)
880     labels = kmeans.labels_
881     silhouette_avg = silhouette_score(X, labels)
882     silhouette_scores.append(silhouette_avg)
883 plt.figure(figsize=(8, 5))
884 plt.plot(*args: k_values, wcss, marker='o', linestyle='--')
885 plt.title("Elbow Method for Optimal k using KMeans Clustering")
886 plt.xlabel("Number of Clusters (k)")
887 plt.ylabel("Within-Cluster Sum of Squares (WCSS)")
888 plt.xticks(k_values)
889 plt.grid(True)
890 plt.show()
891 plt.figure(figsize=(8, 5))
892 plt.plot(*args: k_values, silhouette_scores, marker='o', linestyle='--')
893 plt.title("Silhouette Analysis for Optimal k using KMeans Clustering")
894 plt.xlabel("Number of Clusters (k)")
895 plt.ylabel("Silhouette Score")
896 plt.xticks(k_values)
897 plt.grid(True)
898 plt.show()

```

DBSCAN algorithm:

```

900 #Dbscan Algorithm
901 k = 57
902 nbrs = NearestNeighbors(n_neighbors=k).fit(X)
903 distances, indices = nbrs.kneighbors(X)
904 distances = np.sort(distances[:, -1])
905 plt.figure(figsize=(10, 6))
906 plt.plot(distances)
907 plt.title("K-Distance Graph")
908 plt.xlabel("Data Points (sorted)")
909 plt.ylabel(f"{k}-Distance")
910 plt.grid(True)
911 plt.show()
912 eps = 1.15
913 min_samples = 5
914 dbscan = DBSCAN(eps=eps, min_samples=min_samples)
915 clusters = dbscan.fit_predict(X)
916 df['Cluster'] = clusters
917 n_clusters = len(set(clusters)) - (1 if -1 in clusters else 0)
918 n_noise = list(clusters).count(-1)
919 print(f"Number of clusters: {n_clusters}")
920 print(f"Number of noise points: {n_noise}")
921 if n_clusters > 1:
922     silhouette_avg = silhouette_score(X, clusters)
923     print(f"Silhouette Score: {silhouette_avg}")
924 else:
925     print("Silhouette Score is not defined for a single cluster.")
926 if isinstance(X, np.ndarray):
927     X = pd.DataFrame(X, columns=[f"Feature_{i+1}" for i in range(X.shape[1])])
928 plt.figure(figsize=(10, 7))
929 if X.shape[1] >= 2: # Ensure at least two features for plotting

```

Apriori algorithm:

```

893 df_le = df_apr[['DayOfWeek', 'Flight', 'Length', 'Time']]
894 df_le['DayOfWeek'] = pd.cut(df_le['DayOfWeek'], bins=3, labels=['StartOfWeek',
895 'MidOfWeek', 'Weekend'])
896 df_le['Flight'] = pd.cut(df_le['Flight'], bins=3, labels=['Small Aircraft', 'Medium Aircraft', 'Large Aircraft'])
897 df_le['Length'] = pd.cut(df_le['Length'], bins=3, labels=['Short Length', 'Average Length', 'Long Length'])
898 df_le['Time'] = pd.cut(df_le['Time'], bins=3, labels=['Midnight', 'Early Morning', 'Day'])
899 te = TransactionEncoder()
900 te_ary = te.fit(df_le.values).transform(df_le.values)
901 df_le = pd.DataFrame(te_ary, columns=te.columns_)
902 print(df_le.head(5))
903 df_le = df_le.astype('int')
904 print(df_le.head(5))
905 frequent_itemsets = apriori(df_le, min_support=0.1, use_colnames=True, verbose=1)
906 rules = association_rules(frequent_itemsets, metric="lift", min_threshold=0.2, num_itemsets=len(frequent_itemsets))
907 rules = rules.sort_values(by=['confidence'], ascending=False)
908 print(rules.head(5).to_string())

```

## References

- 1) Airlines Dataset to predict a delay. (n.d.). www.kaggle.com.  
<https://www.kaggle.com/datasets/jimschacko/airlines-dataset-to-predict-a-delay>
- 2) scikit-learn. (2019). scikit-learn: machine learning in Python. Scikit-Learn.org.  
<https://scikit-learn.org/stable/>
- 3) Raschka, S. (n.d.). mlxtend. Rasbt.github.io. <https://rasbt.github.io/mlxtend/>
- 4) Pandas. (2018). Python Data Analysis Library — pandas: Python Data Analysis Library. Pydata.org. <https://pandas.pydata.org/>
- 5) <https://www.geeksforgeeks.org/dbSCAN-clustering-in-ml-density-based-clustering/>
- 6) <https://www.analyticsvidhya.com/blog/2021/06/understand-the-dbSCAN-clustering-algorithm/>