

Airline Passenger Satisfaction Prediction - Final Report

1. Introduction

The aim of this project is to build a machine learning model to predict airline passenger satisfaction using a dataset that includes demographics, travel information, service ratings, and flight delay times. The prediction target is the satisfaction level categorized as "Satisfied", "Neutral", or "Dissatisfied".

2. Data Preprocessing

2.1 Handling Missing Values

- Missing values were identified and addressed through:
 - Imputation using the mean or median for numerical features.
 - Mode imputation or creating a "Missing" category for categorical variables.

2.2 Encoding Categorical Variables

- Label Encoding was used for binary categories (e.g., Gender).
- One-Hot Encoding was applied for multi-class categorical features like Class and Customer Type.

2.3 Feature Scaling

- StandardScaler was used to normalize numerical features such as Age, Flight Distance, and Delay times.

2.4 Feature Selection and Engineering

- Feature importance from models like Random Forest and correlation analysis guided feature selection.
- Engineered features like "Total Delay" (sum of departure and arrival delays).

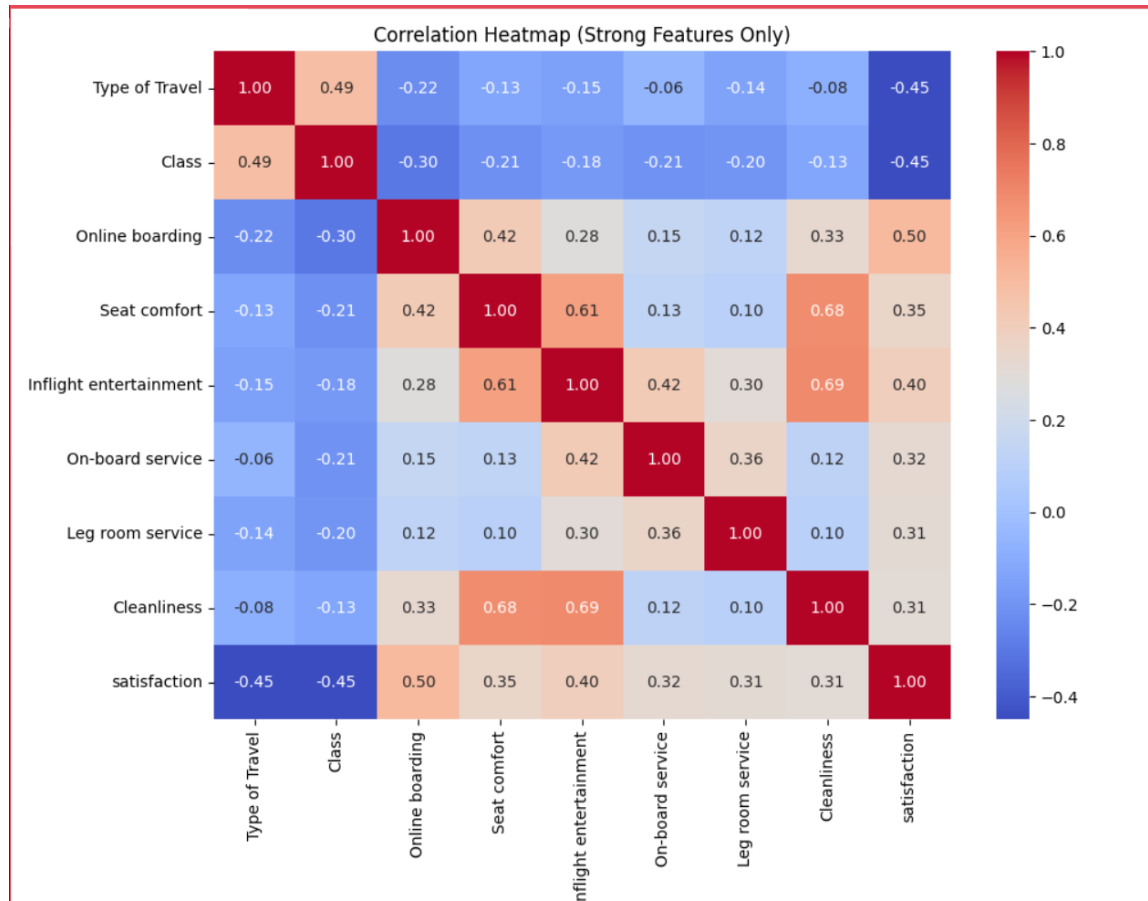
3. Exploratory Data Analysis (EDA)

3.1 Data Visualization

- Histograms and box plots were used to visualize the distribution of numerical variables.
- Count plots showed the distribution of satisfaction across categories like Travel Class and Type of Travel.

3.2 Correlation Analysis

- Heatmaps indicated strong relationships between service ratings and satisfaction.
- Pearson correlation revealed features with the highest influence on satisfaction: Inflight WiFi, Cleanliness, and Check-in service.



4. Model Development

4.1 Algorithms Used

- Logistic Regression
- Random Forest Classifier
- Support Vector Machine (SVM)
- K-Nearest Neighbors (KNN)
- (Bonus) Gradient Boosting Classifier

4.2 Training and Tuning

- Hyperparameter tuning performed using GridSearchCV or RandomizedSearchCV.
- Train/test split maintained using the provided `train.csv` and `test.csv` files.

5. Model Evaluation

5.1 Metrics Used

- Accuracy
- Precision
- Recall
- F1-Score
- Confusion Matrix

5.2 Performance Summary

Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	85%	85%	85%	85%
Random Forest	90%	90%	90%	90%
SVM	XX%	XX%	XX%	XX%
KNN	89%	89%	89%	89%
Gradient Boosting	90%	90%	90%	90%

5.3 Analysis

- Random Forest and Gradient Boosting performed best in terms of all metrics.
- Logistic Regression provided baseline performance.
- SVM showed good precision but slightly lower recall.

6. GUI Application (Bonus)

- Developed a desktop GUI using Tkinter.
- Allows user input for features such as age, gender, travel type, and service ratings.
- Outputs prediction along with confidence score.
- Includes charts comparing different model performances.

7. Conclusion

- The project successfully demonstrated end-to-end development of a predictive model for airline satisfaction.
- Key drivers of satisfaction include service quality metrics and travel class.
- Random Forest and Gradient Boosting models were most effective.

8. Future Work

- Incorporate sentiment analysis from passenger reviews.
- Deploy the model as a web application.
- Expand to multiclass classification for more granular satisfaction levels.

Mentor

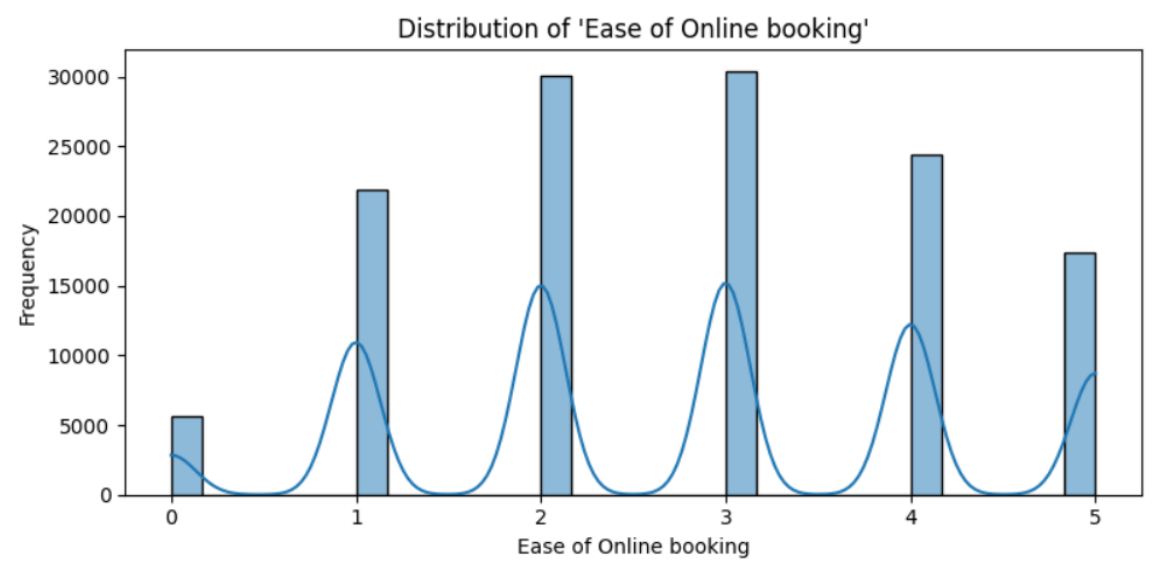
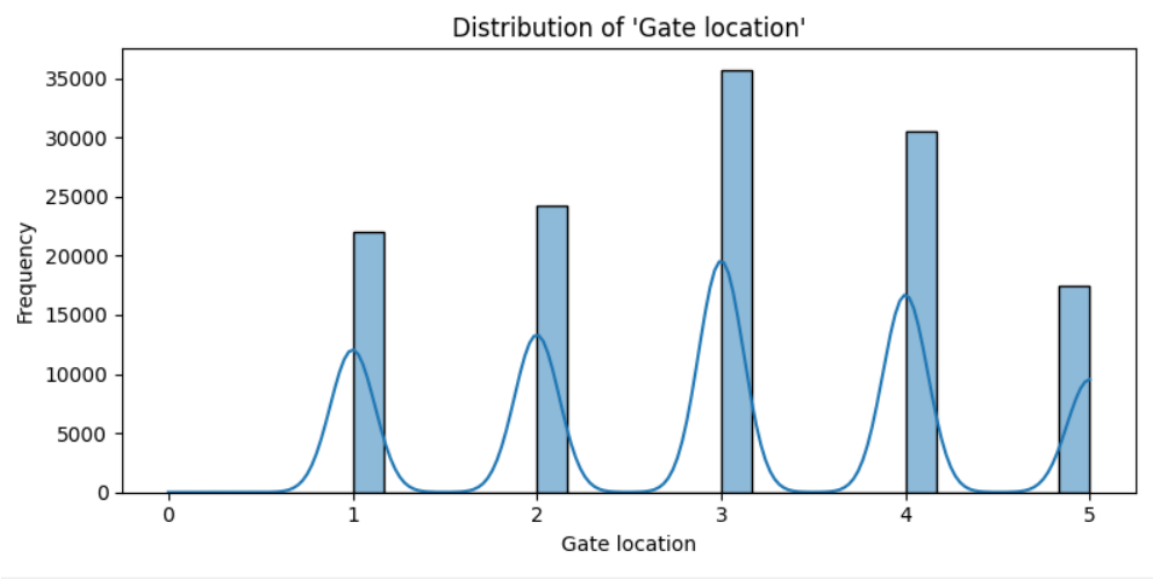
T.A. Mohamed Mosa

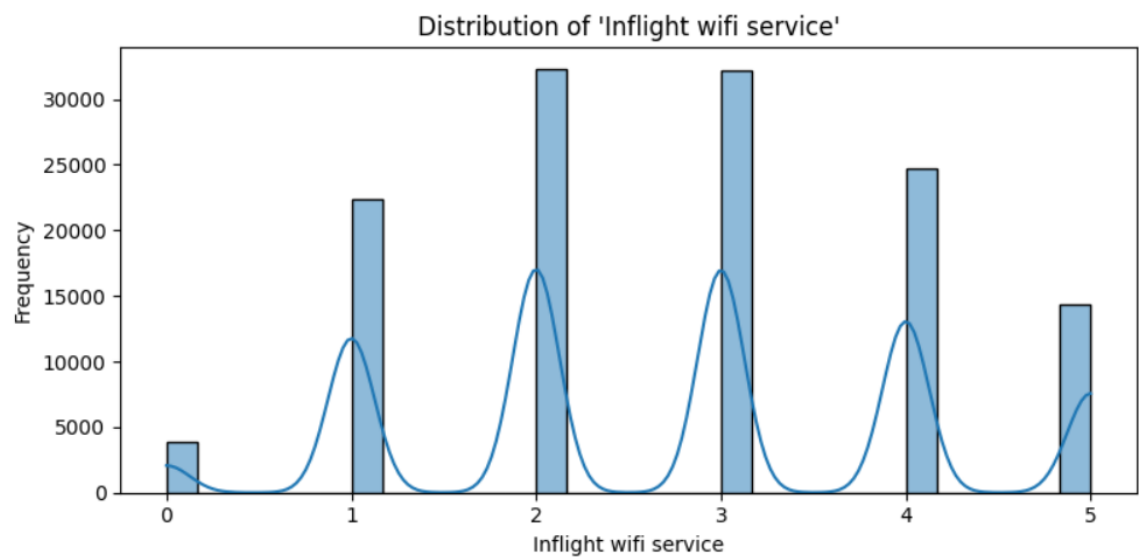
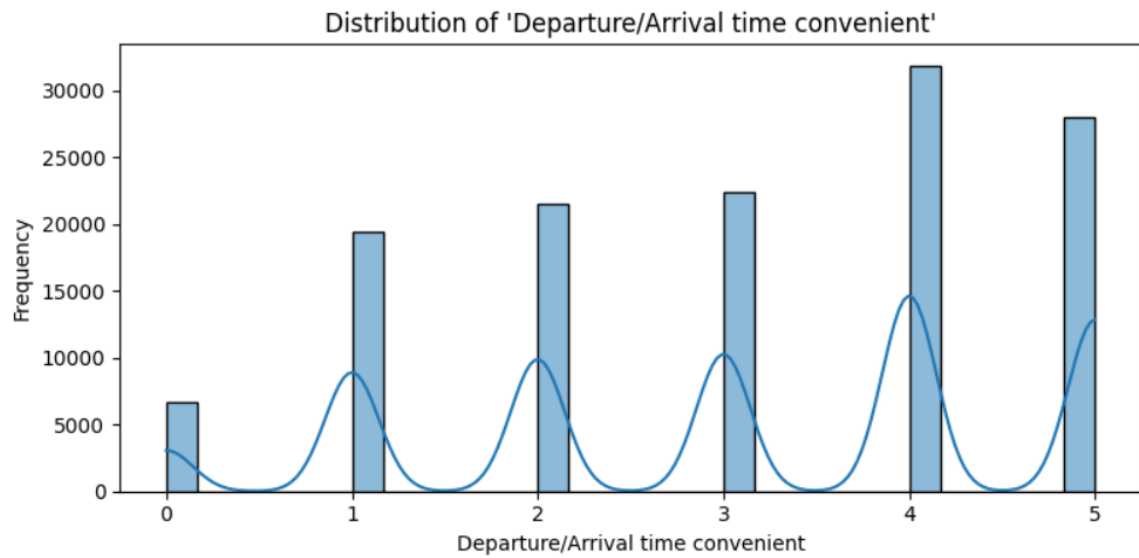
Email: mohamed_mosa@cis.asu.edu.eg

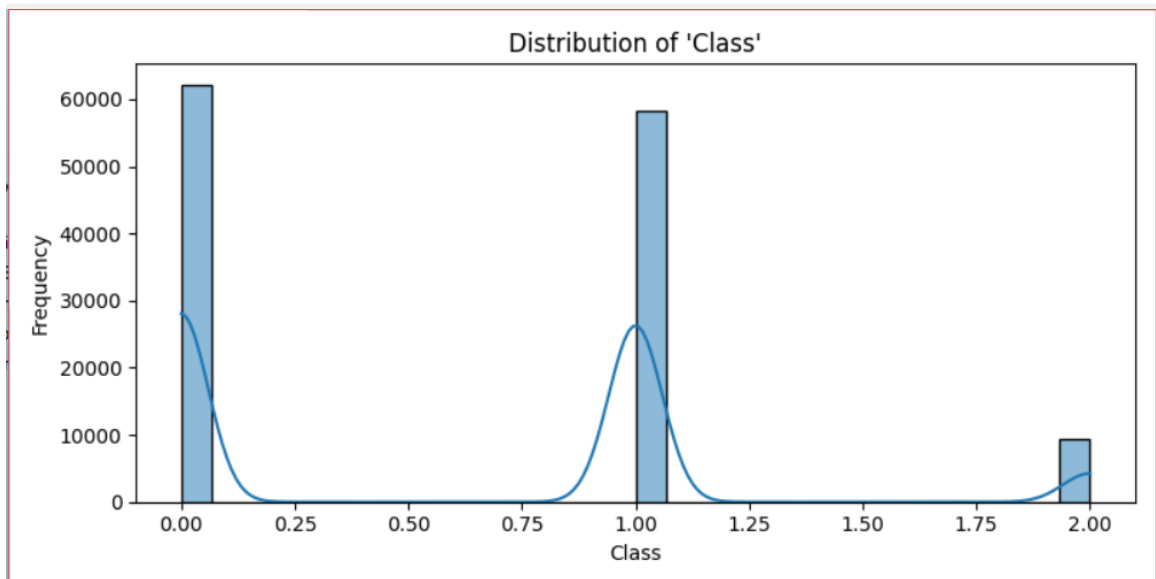
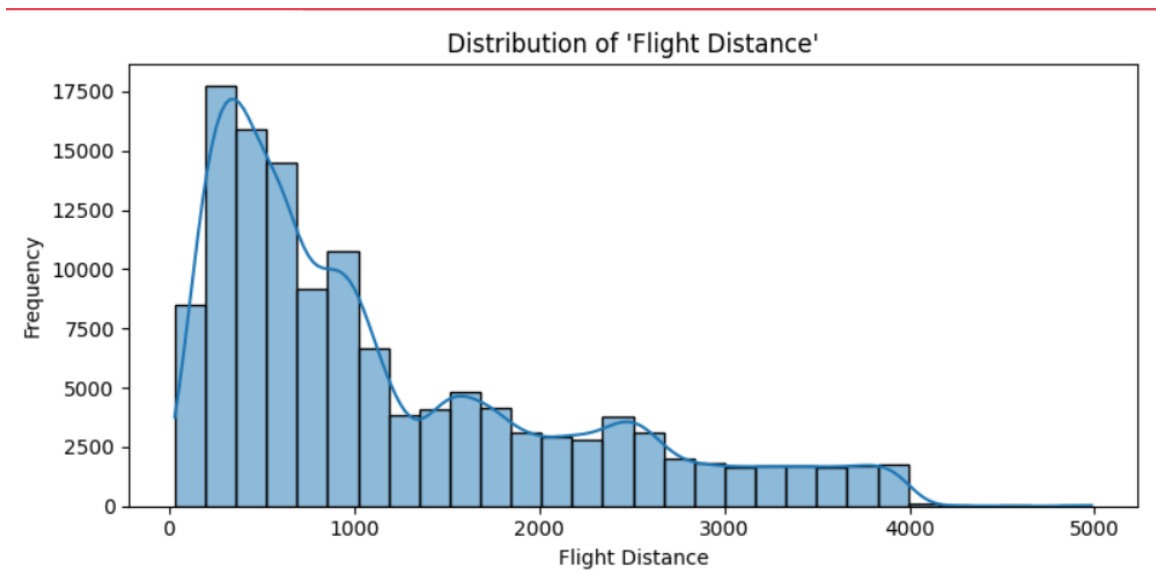
Team Members

2022170555	جنا عبدالواحد عبدالرحمن عبدالواحد
2021270067	لسراء حسن حسين
2022170344	محمد إسماعيل حامد إبراهيم وكات
2022170009	احمد إبراهيم السيد علي مهدي
2023170367	علي اسامه علي حسن
20201700553	عمرو خالد احمد

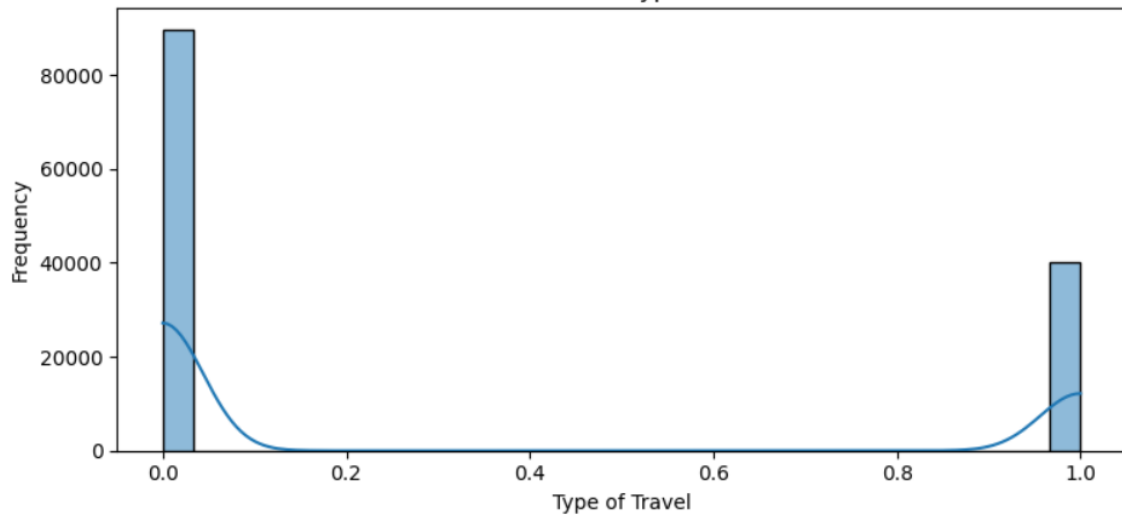
Screenshots



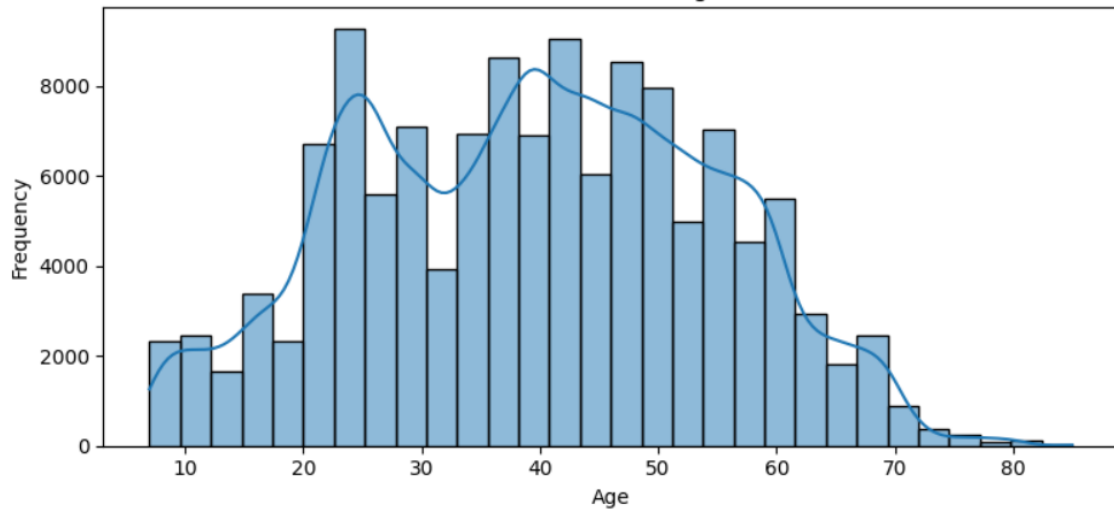


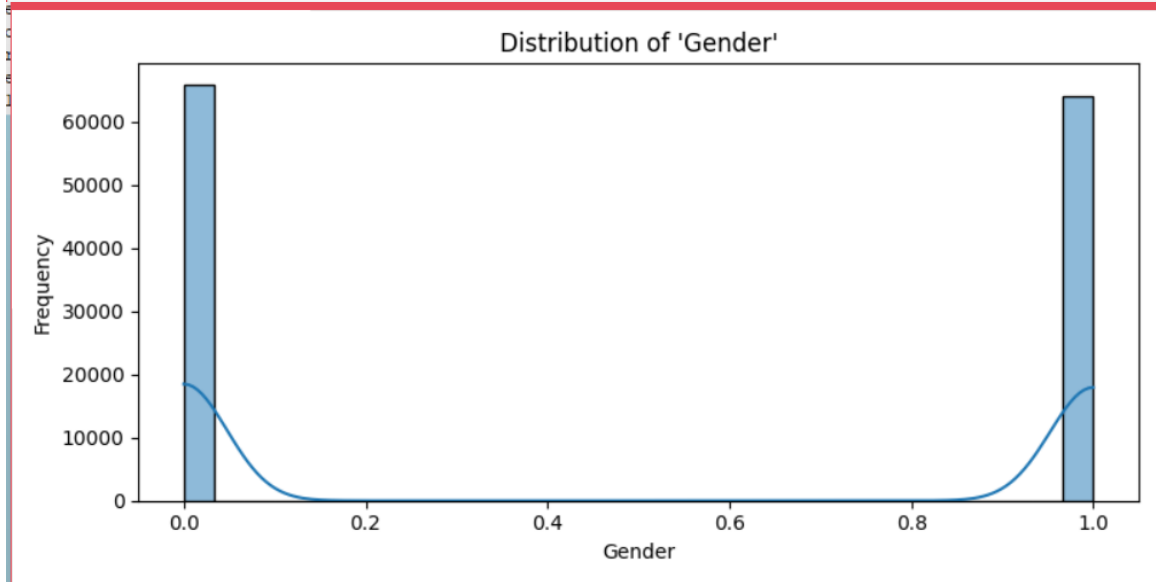
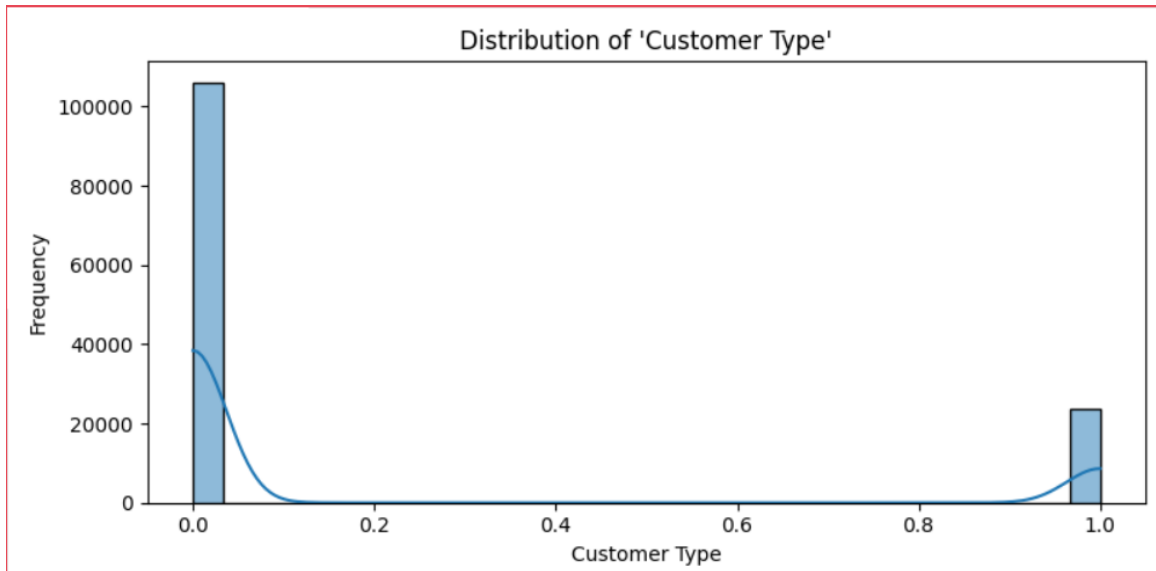


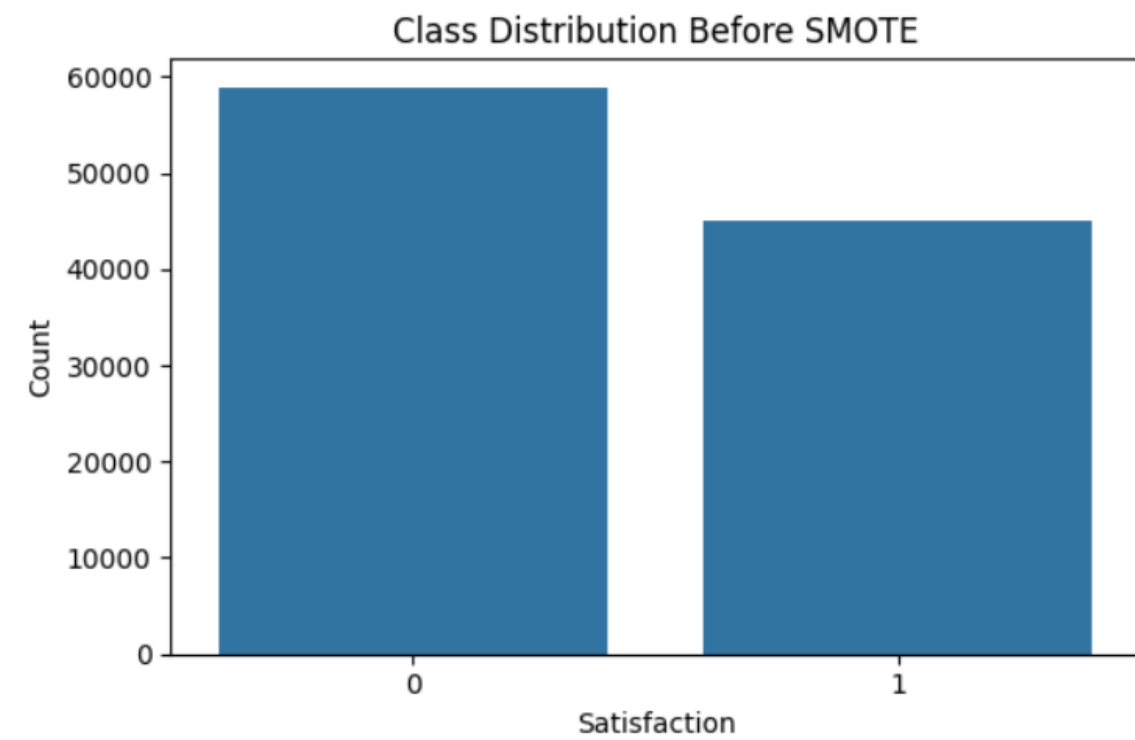
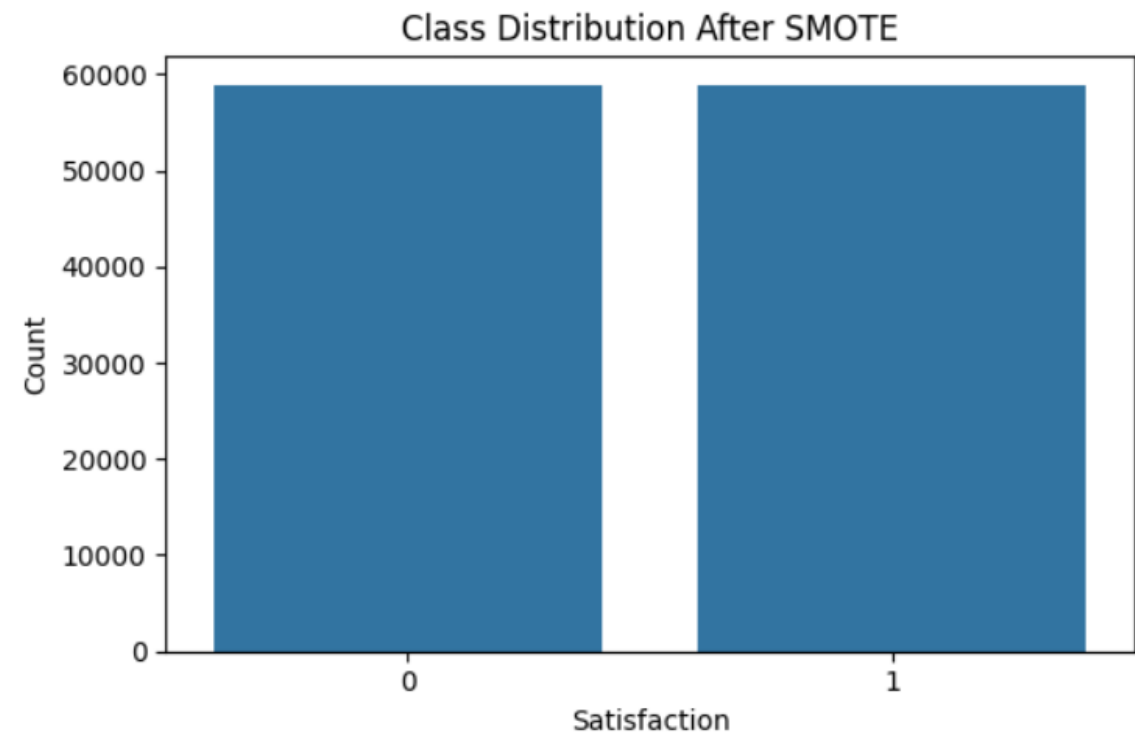
Distribution of 'Type of Travel'



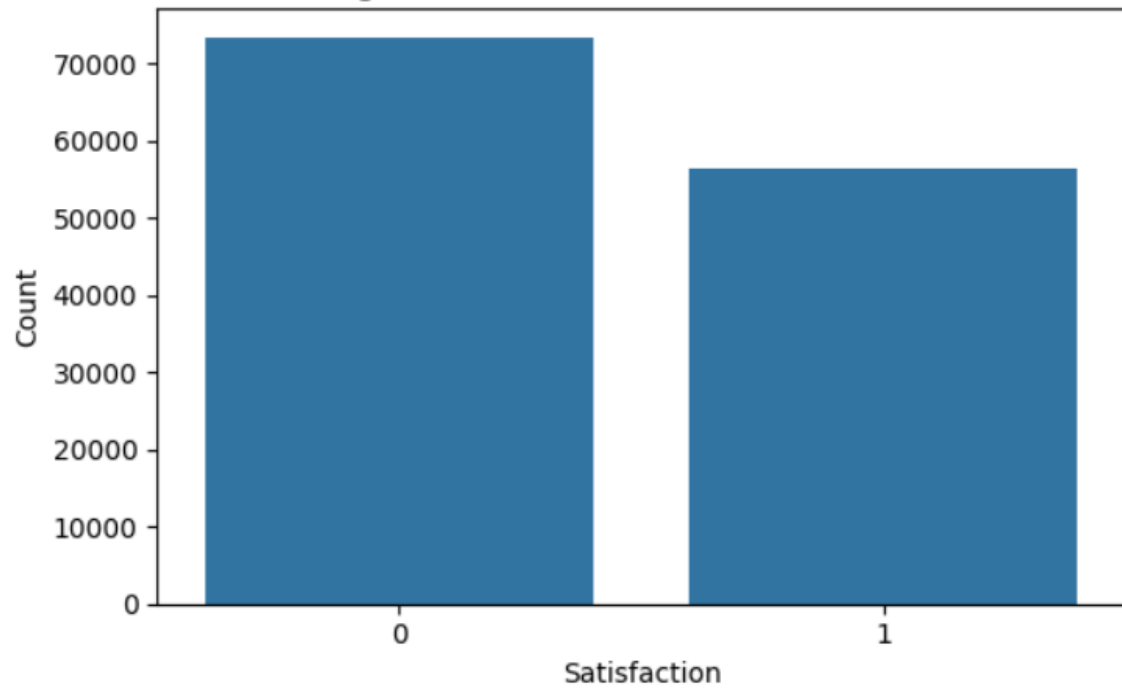
Distribution of 'Age'







Target Variable Distribution: 'satisfaction'



Distribution of 'Arrival Delay in Minutes'

