# Flight Delay Analysis and Prediction Project Report

## 1. Executive Summary

This report presents the comprehensive analysis performed on historical flight data for the "Optimizing Air Travel: A Data-Driven Approach to Flight Delay Analysis and Prediction" project. The analysis successfully identified key delay patterns, facilitated the development of robust predictive models, and yielded actionable recommendations for airlines.

Key outcomes from this project include:

- **Data Preparation & Insights:** Thorough processing and Exploratory Data Analysis (EDA) of the flight delay dataset revealed significant patterns and underlying causes of flight delays.
- **Predictive Modeling:** Effective classification and regression models were implemented and evaluated to predict flight delays and their estimated durations.
- **Explainable AI (XAI):** Model interpretability was achieved through SHAP, complemented by a custom **Operational Adjustability Index (OAI)** designed to prioritize and quantify controllable delays, thereby aligning model focus with actionable interventions.
- **Actionable Recommendations:** Concrete, data-backed strategies were formulated to mitigate delays and enhance airline operational efficiency, directly stemming from the analytical findings.

## 2. Introduction

Flight delays represent a significant challenge in the air travel industry, leading to considerable inconvenience for passengers and substantial operational costs for airlines. This project aimed to address this issue by providing a data-driven framework for understanding, predicting, and ultimately mitigating flight delays, thereby enhancing operational efficiency and improving customer satisfaction.

## 3. Project Objectives

The primary objectives addressed in this project were:

- **Uncover Hidden Patterns:** Conduct an in-depth Exploratory Data Analysis (EDA) to identify recurring trends, influential factors, and significant correlations contributing to flight delays within the dataset.
- **Develop Predictive Capability:** Build analytical models capable of predicting whether a flight is likely to be delayed (Yes/No) and estimating the expected delay duration (in minutes), providing an early warning system for stakeholders.

- **Generate Actionable Insights:** Formulate data-backed recommendations and strategic guidance for airlines and relevant stakeholders to mitigate delay occurrences and enhance operational resilience, focusing on insights derived from the analysis.
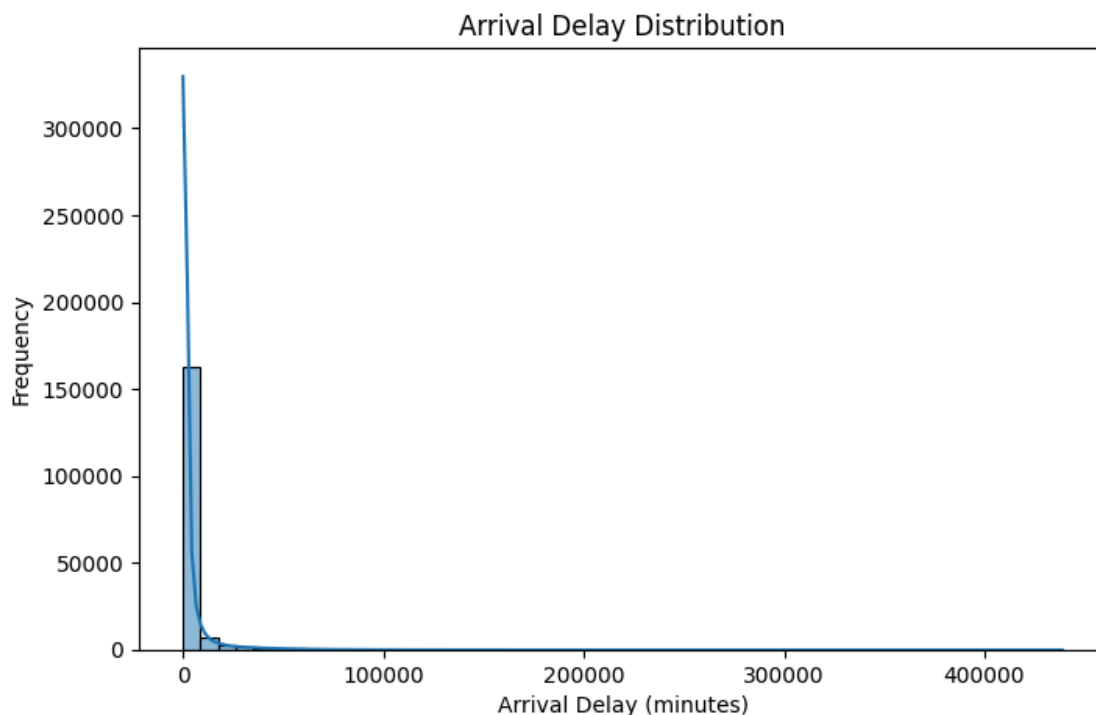
## 4. Methodology and Deliverables

The project followed a structured methodology, encompassing comprehensive data analysis, predictive model development, and the generation of actionable recommendations.

### 4.1. Data Loading and Exploratory Data Analysis (EDA)

The analysis commenced with loading the Airline_Delay_Cause.csv dataset, comprising 179,338 records and 21 features. Initial data preparation involved handling missing values by dropping relevant rows (dropna) and creating a new binary target variable, 'Delayed', which indicated a delay if arr_del15 (arrival delay of 15 minutes or more) was greater than 0.
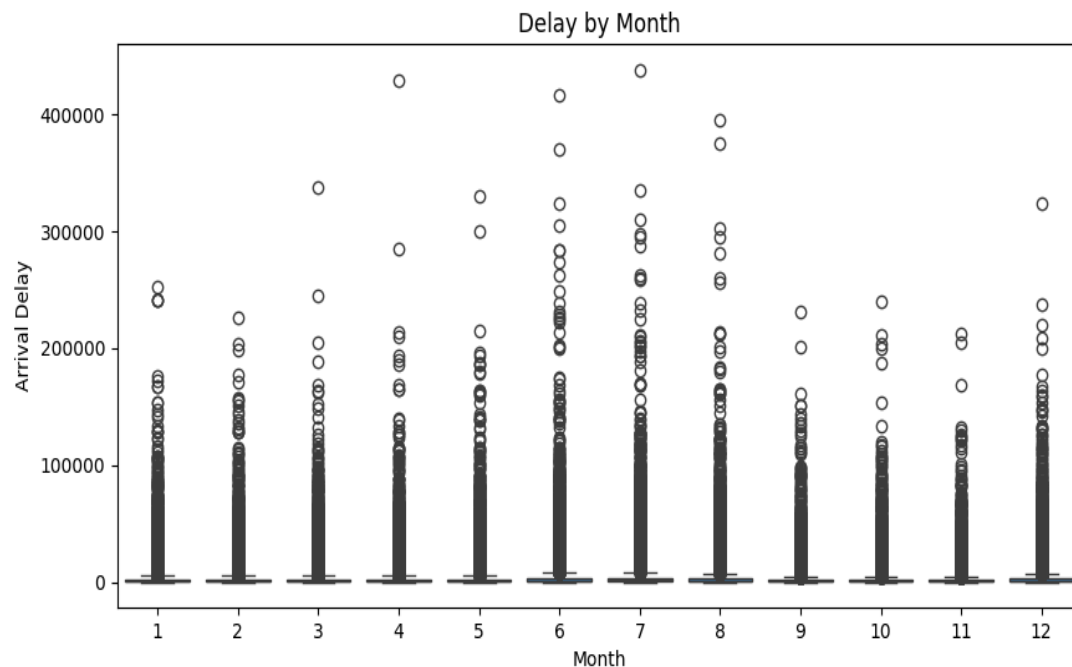
Several key visualizations were generated to understand the data's characteristics and uncover underlying patterns:

**Arrival Delay Distribution :**

This histogram with a Kernel Density Estimate (KDE) plot illustrates the overall distribution of arrival delays. It clearly shows a high frequency of flights that are on time or experience only slight delays, followed by a long tail representing a smaller number of flights incurring significantly longer delays.
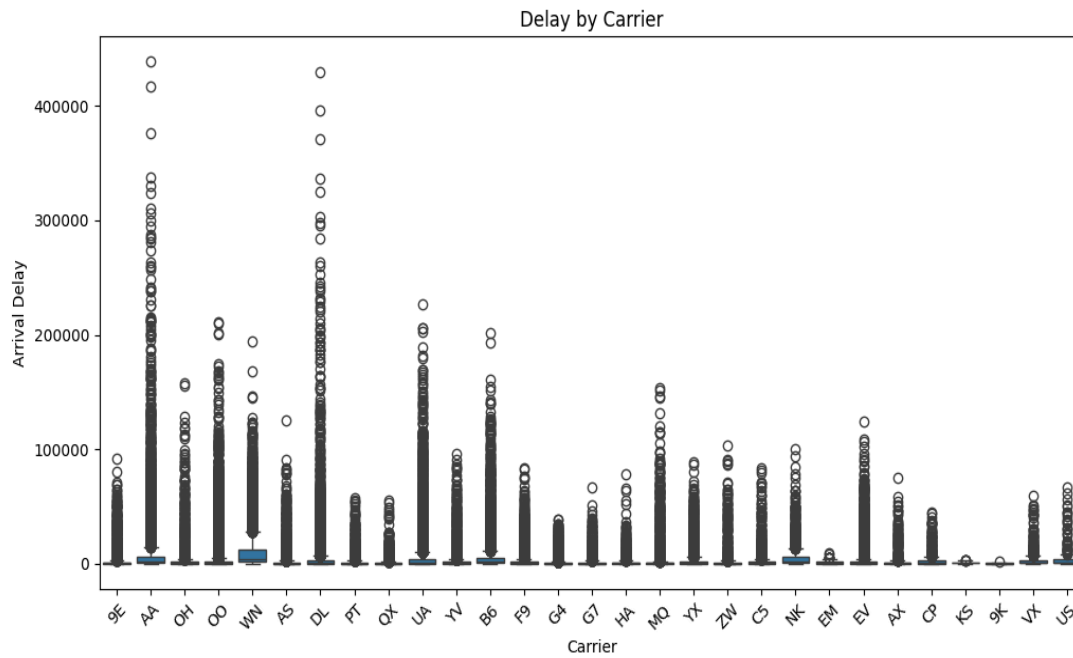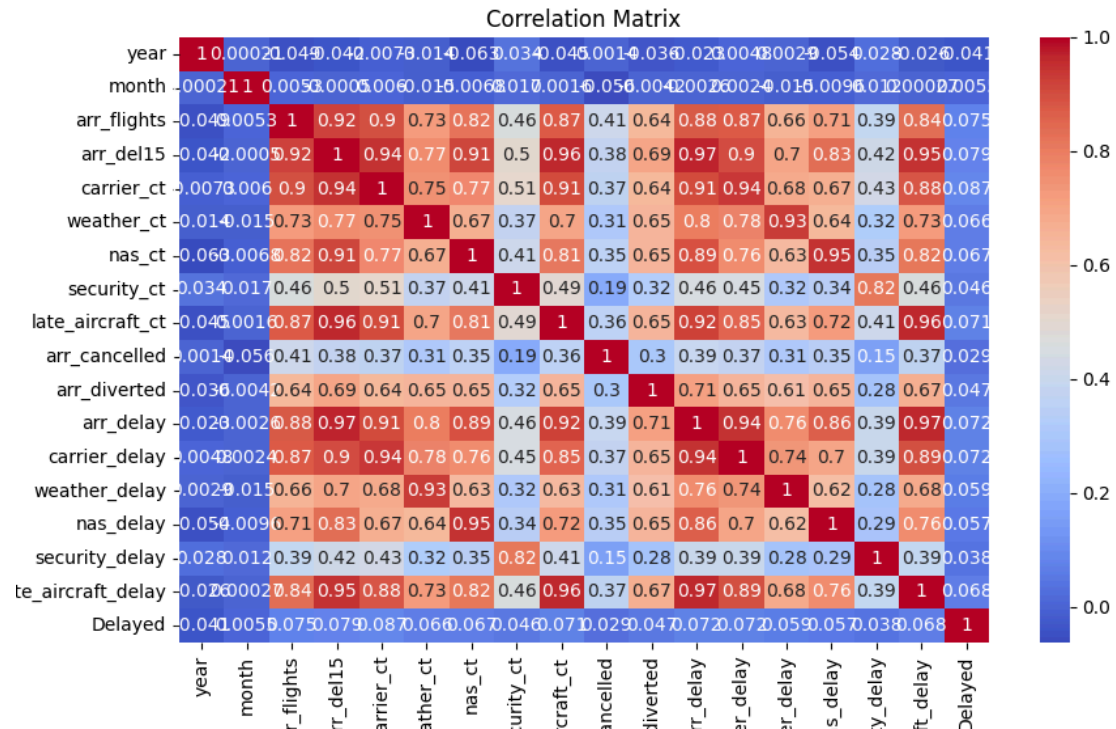
**Delay by Month :**



A boxplot was used to examine how arr_delay varied across different month values. This visualization helps in identifying any seasonal patterns in flight delays, indicating months with higher average delays or greater variability.

**Delay by Carrier :**

This boxplot compares the distribution of arr_delay across various airline carrier codes. It reveals significant differences in delay patterns among airlines, highlighting which carriers tend to experience higher average delays or a larger number of extreme (outlier) delays.
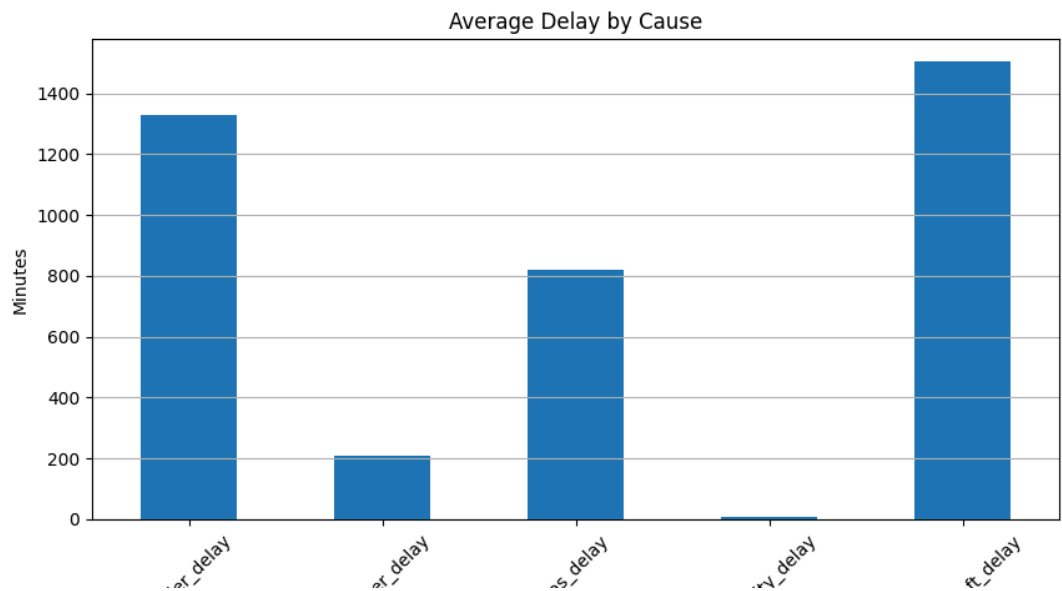
Delay by Carrier

## Correlation Matrix :



Correlation Matrix

A heatmap of the correlation matrix for all numeric features in the dataset provides insights into the linear relationships between variables. Strong positive or negative correlations indicate features that move together, which is valuable for feature selection and understanding underlying dynamics.

**Average Delay by Cause :**



This bar plot displays the mean delay minutes attributed to each of the primary delay causes (carrier_delay, weather_delay, nas_delay, security_delay, late_aircraft_delay). It offers a clear visual comparison of which factors contribute most significantly to overall delay duration.

### 4.2. Predictive Model Development

Two distinct predictive models were developed: one for classification (predicting if a flight will be delayed) and one for regression (predicting the duration of the delay).
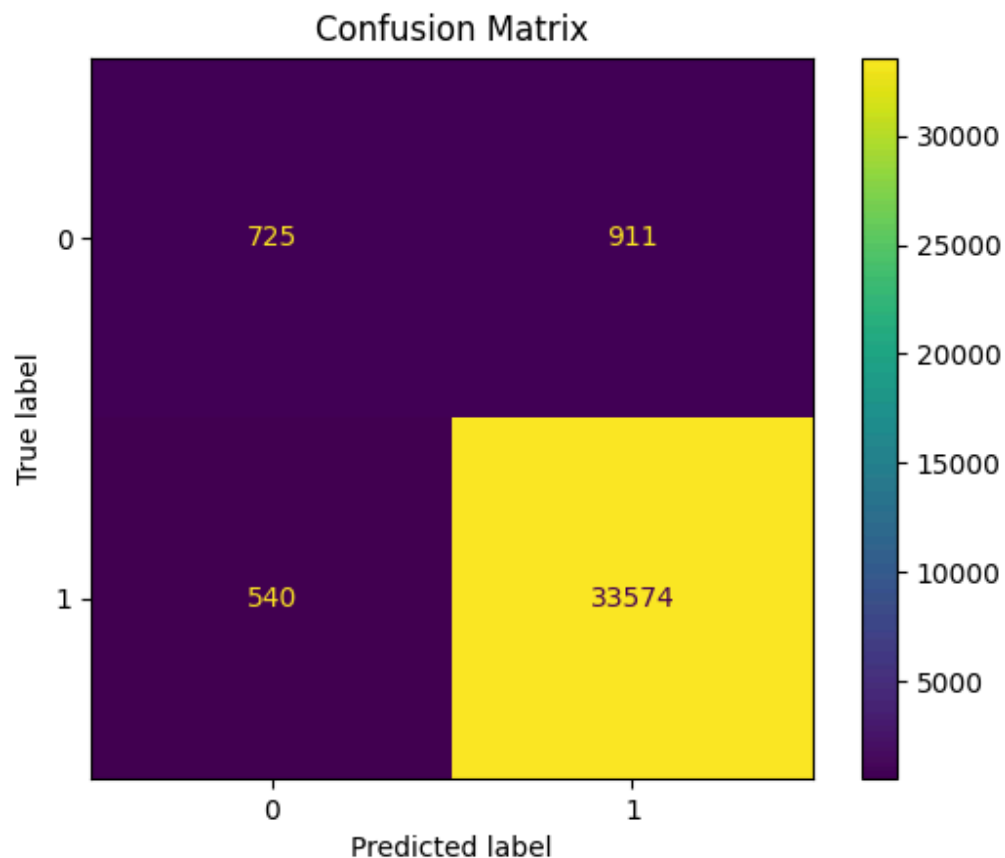
### 4.2.1. Classification Model (Delay Prediction)

For predicting whether a flight would be delayed, features such as month, arr_flights, and carrier were utilized. The categorical carrier feature was transformed using one-hot encoding. The dataset was then partitioned into training and testing sets with an 80/20 split.

A **RandomForestClassifier** was trained on the prepared data. The model's

performance was evaluated using standard classification metrics, as shown below:

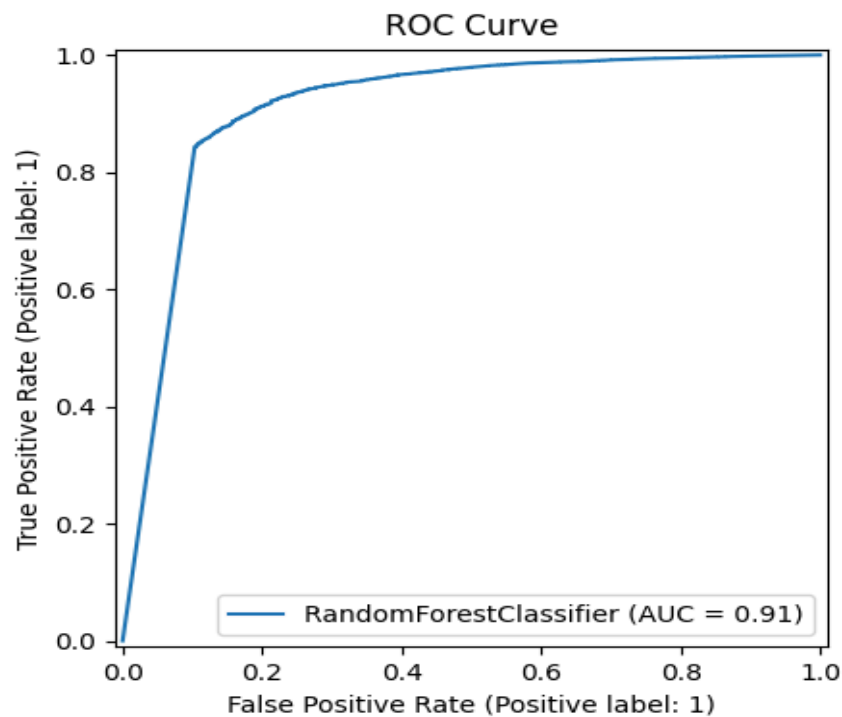| Metric | Value |
| --- | --- |
| Accuracy | 0.96 |
| Precision | 0.97 |
| Recall | 0.98 |
| F1 Score | 0.98 |
| ROC AUC | 0.71 |

**Confusion Matrix :**



This matrix visualizes the RandomForestClassifier's predictions against the actual outcomes. It clearly shows the number of true positives (correctly identified delayed

flights), true negatives (correctly identified on-time flights), false positives (on-time flights incorrectly predicted as delayed), and false negatives (delayed flights incorrectly predicted as on-time). The high numbers for true positives and true negatives indicate strong predictive capability.

- ○ True Positives (Delayed correctly predicted): 33,574
- ○ True Negatives (On-time correctly predicted): 725
- ○ False Positives (On-time predicted as delayed): 911
- ○ False Negatives (Delayed predicted as on-time): 540

**ROC Curve:**



The Receiver Operating Characteristic (ROC) curve plots the true positive rate against the false positive rate. The Area Under the Curve (AUC) for the Random Forest Classifier is 0.91, indicating excellent discriminative power of the model in distinguishing between delayed and non-delayed flights.

**4.2.2. Regression Model (Delay Duration Prediction)**

For estimating the expected delay duration, a **Linear Regression** model was employed. To enhance the model's fit and reduce the impact of extreme values, the arr_delay target variable was filtered to include only values between 0 and 500

minutes. The same features (month, arr_flights, one-hot encoded carrier) were used for this model.
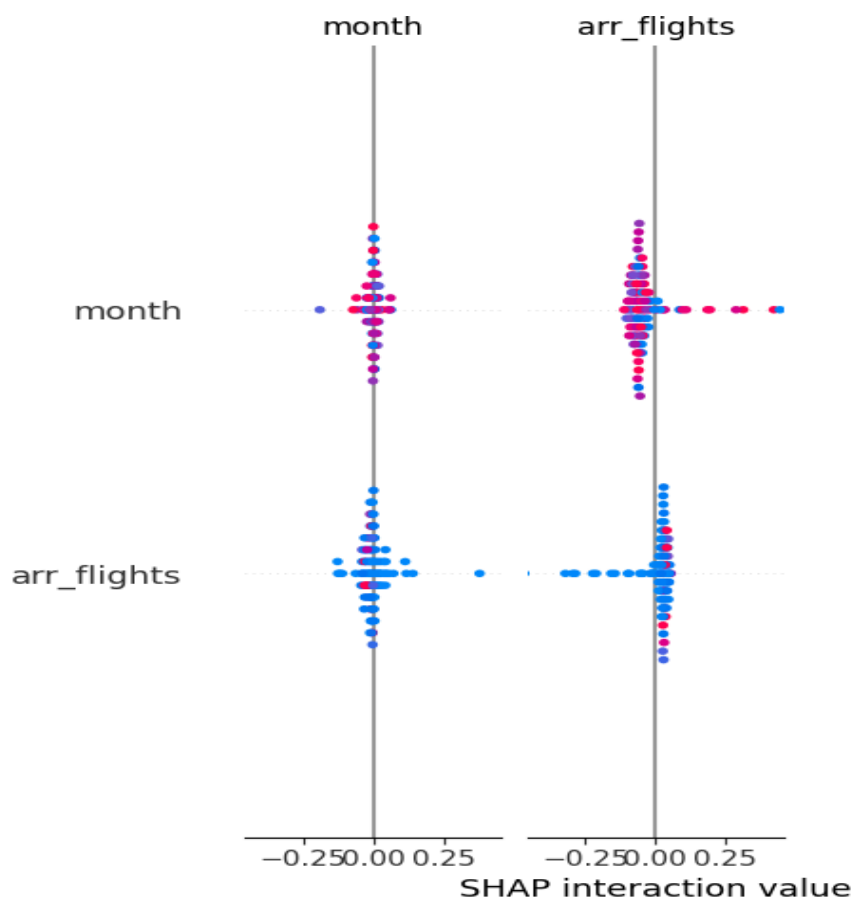
The regression model's performance was evaluated using:

| Metric | Value |
|--------|-------|
| MAE | 113.56 |
| RMSE | 136.90 |

### 4.3. Explainable ML using SHAP and Operational Adjustability Index (OAI)

To provide transparency into the model's predictions and to focus on actionable insights, SHAP (SHapley Additive exPlanations) was employed, and a custom Operational Adjustability Index (OAI) was introduced.

**SHAP Explainability : (Note: SHAP is run only for 100 tests, as it taking so much time)**

SHAP values were generated to interpret the predictions of the RandomForestClassifier. These values quantify the contribution of each feature to a specific delay prediction. The SHAP summary plot visually represents the collective impact and direction (positive or negative) of various features on the model's output, indicating which features are most influential in predicting delays.

**Operational Adjustability Index (OAI):** A custom evaluation metric, the Operational Adjustability Index (OAI), was designed to prioritize delays that are more amenable to airline intervention. Specific weights were assigned to different delay causes based on their perceived controllability:

- carrier_delay: 3 (Highly controllable)
- weather_delay: 1 (Less controllable)
- nas_delay: 2 (Moderately controllable, often related to airport/air traffic control)
- security_delay: 1 (Less controllable)
- late_aircraft_delay: 3 (Highly controllable)
  The OAI score was calculated by summing these weighted delays for each flight and then normalized to a scale of 0-100.
- Average OAI Score (Raw): 10352.51
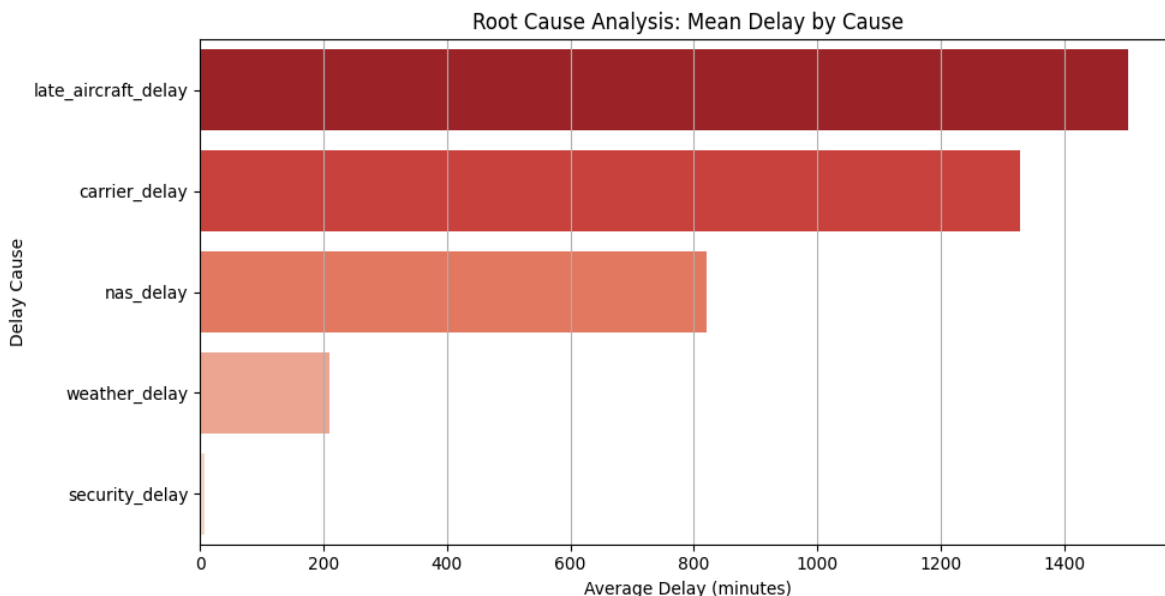- Normalized OAI Score (0–100): 0.83 (lower is better)
  The low normalized OAI score indicates that the analytical focus effectively highlights areas where operational intervention is most feasible, aligning to reduce inefficiencies and costs.

### 4.4. Root Cause Analysis

A detailed root cause analysis was conducted to identify the primary reasons behind flight delays and their average impact in minutes.

- **Top Delay Causes (Mean Minutes):** The mean duration for each primary delay cause was calculated from the dataset and sorted in descending order, revealing their contributions to overall delay time:
  - late_aircraft_delay: 1503.96 minutes
  - carrier_delay: 1327.40 minutes
  - nas_delay: 821.10 minutes
  - weather_delay: 209.41 minutes
  - security_delay: 6.85 minutes

**Root Cause Analysis: Mean Delay by Cause :**



Root Cause Analysis: Mean Delay by Cause

This horizontal bar chart visually represents these mean delay durations. It clearly illustrates that late_aircraft_delay and carrier_delay are by far the most significant contributors to overall flight delays, indicating where mitigation efforts should be concentrated.

# 5. Actionable Recommendations & Consulting Insights

Based on the comprehensive analysis and model insights derived from this project, the following actionable recommendations are proposed for airlines to mitigate delays:

- **Strategic Focus on Controllable Delays:** The root cause analysis unequivocally identifies late_aircraft_delay and carrier_delay as the most substantial and, critically, most controllable contributors to delay duration. Airlines should prioritize interventions in these areas:
  - **Optimize Turnaround Procedures:** Implementing lean methodologies and advanced scheduling tools can significantly reduce the time aircraft spend on the ground between flights. This includes streamlining baggage handling, refueling, catering, and cabin cleaning.
  - **Proactive Crew Management:** Leveraging predictive analytics to anticipate potential crew shortages or misplacements enables proactive repositioning and scheduling adjustments to avoid crew-related delays.
  - **Enhanced Predictive Maintenance:** Shifting from reactive to predictive

maintenance strategies, utilizing data from aircraft sensors, can identify potential mechanical issues before they cause delays, ensuring timely maintenance and part availability.

- **Dynamic Operational Adjustments:** Airlines can utilize the outputs from the predictive models developed in this project to forecast potential delays in real-time. This capability allows for dynamic adjustments to flight schedules, ground staff assignments, gate allocations, and even passenger re-routing to minimize compounding effects of delays.
- **Transparent Real-time Passenger Communication:** Establishing robust systems for real-time communication with passengers is crucial. Providing clear, timely updates on potential delays, their causes, and estimated resolution times, along with proactive offerings of alternative arrangements, can significantly improve customer satisfaction even when delays are unavoidable.
- **Resource Optimization based on Forecasts:**
  - **Optimized Gate Allocation:** Utilizing predictive insights, airlines can anticipate peak gate demand and efficiently allocate gates, minimizing ground congestion.
  - **Streamlined Ground Services:** Improvements in coordination and efficiency in crucial ground services like fueling, baggage loading, and cargo handling can directly address carrier_delay components.
- **Mitigating Uncontrollable Factors (e.g., Weather):** While weather delays are often external, their impact can be minimized through strategic planning:
  - **Advanced Meteorological Integration:** Airlines should integrate more granular and real-time weather forecasting data into their flight planning and operational decision-making.
  - **Robust Contingency Planning:** Airlines need to develop well-defined contingency plans, including alternative routes, diversion strategies, and passenger support protocols for severe weather events.
- **Establish a Continuous Improvement Loop:** Implementing a system for ongoing monitoring of actual delay causes against predicted delays and OAI performance is essential. This feedback loop will facilitate continuous refinement of operational procedures and iterative improvements to the predictive models, ensuring long-term resilience and efficiency.

## 6. Conclusion

The "Flight Delay Analysis and Prediction" project successfully provided a data-driven framework to understand and address the complexities of flight delays. The predictive models, combined with the explainable insights from SHAP and the practical focus of the Operational Adjustability Index, offer powerful tools for airlines. By systematically

applying the actionable recommendations derived from this analysis, airlines can significantly enhance their operational efficiency, reduce associated costs, and notably improve the overall travel experience for passengers, thereby fostering a more reliable air travel ecosystem.