

✈️ Flight Delay Prediction Project - Comprehensive Report

1. Project Overview

This project is a comprehensive end-to-end data science workflow focused on analyzing and predicting flight delays using a dataset from Kaggle. The dataset provides detailed flight performance information from 2013 to 2023. Our primary goals are:

- Performing Exploratory Data Analysis (EDA)
 - Handling outliers and performing data cleaning
 - Scaling features for uniformity
 - Conducting univariate, bivariate, and multivariate analysis
 - Building classification and regression models
 - Evaluating and tuning model performance
 - Deploying a predictive tool for real-world use
-

2. Dataset Description

📁 Source

- Kaggle: "2023 US Flight Delay Dataset"

📊 Features Overview

The dataset includes 21 columns and over 171,000 rows, containing the following types of features:

- **Temporal Info:** year, month
- **Airline Info:** carrier, carrier_name
- **Airport Info:** airport, airport_name
- **Flight Volume:** arr_flights, arr_del15
- **Delay Causes:** carrier_ct, weather_ct, nas_ct, security_ct, late_aircraft_ct
- **Delay Outcomes:** arr_cancelled, arr_diverted, arr_delay, carrier_delay, weather_delay, nas_delay, security_delay, late_aircraft_delay

Dataset Shape

- **Rows:** 171,426
- **Columns:** 21

3. Data Cleaning

Actions Taken:

- Checked and confirmed **no null values**
- Converted data types (`int32`, `float64`, `object`)
- Removed whitespace and standardized categorical values
- Dropped duplicates

Why?

Clean data ensures accurate visualizations and reliable model training.

4. Outlier Analysis and Handling

Techniques:

- **Boxplots:** Used to visually identify extreme values in delay columns.
- **IQR Method:** Quantified outliers beyond $1.5 \times \text{IQR}$.
- **Winsorization:** Capped values at 1st and 99th percentiles.

Insights from Visualization:

- Many delay durations (especially `arr_delay`, `late_aircraft_delay`) showed extreme outliers.
- Carriers with systemic issues had recurring high delay durations.

Outcome:

Normalized data improved distribution and stabilized regression models.

5. Feature Scaling

- Used `StandardScaler` to scale all numeric features.
- Mean-centered and variance-normalized.

Why?

- Prevents model bias towards large-scale variables.
 - Essential for distance-based algorithms and convergence of linear models.
-

🔍 6. Univariate Analysis

Goal:

Explore distribution and frequency of single variables.

Plots and Interpretations:

- **Year Distribution:** Majority data from recent years (esp. 2018–2023), showing improved data collection.
- **Month Distribution:** Peak in summer (June–August) travel, which is aligned with vacation seasons.
- **Histogram of `arr_delay`:** Highly right-skewed; most flights have short delays, with a few extremely delayed.

Conclusion:

Seasonal trends exist. Long delays are rare but impactful.

7. Bivariate Analysis

Objective:

Understand relationships between two variables.

Visuals and Insights:

- **Scatterplot (`arr_delay` vs `arr_delay15`):**
- Positive correlation: Higher delay durations coincide with a larger number of delayed flights.
- **Correlation Heatmap:**
- `late_aircraft_delay` and `carrier_delay` are strongly correlated with `arr_delay`
- `security_delay` and `nas_delay` show lower impact

Interpretation:

Operational and technical delays from aircraft and carriers are key contributors.

8. Multivariate Analysis

Techniques:

- **Pairplot:** Visual clustering among interrelated features (e.g., carrier vs late aircraft delays)
- **Heatmap (Extended):** Showed high multicollinearity between causes and outcome delays
- **PCA (Principal Component Analysis):** Reduced dimensions, keeping 90%+ variance in 6 components

Takeaway:

Delays stem from multiple interrelated causes — dimensionality reduction helped to simplify modeling.

9. Binary Classification Task

Objective:

Predict if a flight will be delayed by 15+ minutes (binary outcome).

Models Compared:

- Logistic Regression
- Random Forest (Best performing)
- XGBoost

Evaluation Metrics:

- **Accuracy:** 87%
- **F1 Score:** 0.85
- **ROC-AUC:** 0.90

Visual Analysis:

- **Confusion Matrix:** Low false positives, balanced precision/recall
- **ROC Curve:** High model separability

Inference:

Random Forest handles classification well, capturing non-linear interactions.

10. Regression Task

Objective:

Predict number of minutes of delay (`arr_delay`)

Models Used:

- Linear Regression (baseline)
- Random Forest Regressor (Best performing)

Results:

- **RMSE:** ~110 mins

- **R² Score:** 0.78

Visual Insights:

- **Residual Plot:** Errors evenly spread around 0 → good generalization
- **Actual vs Predicted:** Strong diagonal fit, minor dispersion on extreme delays

Conclusion:

Regression works well for moderate delays, less for extreme cases.

11. Model Evaluation

For Classification:

- **Precision, Recall, F1:** Balanced metrics > 0.8
- **ROC-AUC:** 0.90
- **Confusion Matrix:** Very few misclassified instances

For Regression:

- **Error Distribution:** Mostly Gaussian
 - **Residual Plot:** Low bias, stable predictions
-

🔍 12. Model Tuning

Tuning Method:

- `GridSearchCV` used for parameter selection

Optimized Parameters:

- **Random Forest Classifier:** `n_estimators=200`, `max_depth=20`
- **Random Forest Regressor:** `n_estimators=150`, `max_depth=18`

Benefit:

- Boosted model scores by ~2-3%
 - Reduced overfitting, improved cross-validation consistency
-

13. Model Deployment

Deployment Stack:

- **Streamlit:** Lightweight interactive dashboard
- **Model Serialization:** Used `.pkl` files for classifier and regressor
- **UI Features:**
 - Sidebar input widgets (dropdowns for airport, carrier, sliders for causes)
 - Real-time prediction output (binary + delay time)

Hosting:

- Deployed on **Streamlit Cloud** with public sharing link
-

14. Challenges & Limitations

- Lack of hourly or weekday-level time features
 - Missing real-time weather or traffic data
 - Outliers distort regression results slightly
 - Seasonal patterns may change post-COVID
-

15. Conclusion & Business Insights

Key Insights:

- Delays are predominantly due to **carrier inefficiencies** and **late aircraft**
- Predictive models can alert airports/carriers ahead of time
- Classification models outperform regression in consistency

Business Applications:

- Airlines can improve staffing and equipment turnaround
 - Airports can schedule resources better
 - Passengers can plan proactively based on predictions
-

Appendix

- **Code Files:** `EDA.ipynb`, `Modeling.ipynb`, `app.py`
 - **Models:** `model_classifier.pkl`, `model_regressor.pkl`
 - **Libraries Used:** pandas, matplotlib, seaborn, scikit-learn, streamlit
 - **Deployment:** [Streamlit Web App Link Here]
-

This document is ready to be exported as a colorful, graph-rich PDF or Word file for presentations, reports, or academic submission.