# Flight Delay and Cancellation Prediction

**Team 4**

Dmitry Veretennikov

Gauri Tiwari

Mostafa Omidi

Poojitha Anuvala Setty

Sai Siva Sri Harsha Vasinenku

**George Mason University**

AIT614-005 Big Data Essentials (Spring 2024)

Dr. Lindi Liao

April 23, 2024

# INTRODUCTION

- Every day in the US, 2.9 million passengers fly from nearly 20,000 airports on 45,000 flights (FAA, 2023).

- In 2023, over 200 million US passengers faced flight delays and cancellations, costing $30-34 billion (Junginger, 2023).

- Our project aims to analyze flight data in the Washington DC area to understand patterns and causes of delays and cancellations using Machine Learning techniques.

# OBJECTIVES

- Determine the timing of high-risk periods for flight delays and cancellations departing from Washington DC.

- Develop a predictive model to estimate the probability of flight delay and cancellation.

# PROJECT TIMELINE

The project is divided into four main tasks, each with designated team members and timelines.

**Data Acquisition and Preprocessing:**

Team Members:
Mostafa, Dmitry
Start Date: 3/13/2024
End Date: 3/19/2024

**Exploratory Analysis:**

Team Members:
Dmitry
Start Date: 3/20/2024
End Date: 3/26/2024

**Building and Comparing ML Models:**

Team Members:
Dmitry, Poojitha, Harsha
Start Date: 3/27/2024
End Date: 4/9/2024

**Report Writing and Presentation Preparation:**

Team Members:
Poojitha, Gauri
Start Date: 4/10/2024
End Date: 4/16/2024

# DATASET

- The proposed dataset (Zelazko, 2023) consists of 3 million flight records. Among these, over 79,000 flights originated from Washington, DC.

- The dataset comprises 32 original features (Appendix 1) categorized into 19 decimals, 8 strings, 4 integers, 1 date

# ARCHITECTURE

Databricks

| Flight Dataset | Data Ingestion | Databricks DBFS / MongoDB | PySpark | Spark MLlib |
|---|---|---|---|---|
| | | Storing Data | Preprocessing | Feature Engineering |
| | | Processing Data | Exploratory Analysis | Model Development |
| | | Data Filtering | | Model Training and Evaluation |

# DATA PROCESSING

- **Data Handling:**
  - Uploading dataset to Databricks DBFS.
  - Filtering and selecting flights from Washington DC.
  - Cleaning data.
  - Converting data to appropriate data types.
  - Introducing new variables.
  - Balancing dataset (oversampling).

# ANALYTIC APPROACHES

## Exploratory Analysis with PySpark:

- Utilize PySpark for data exploration.
- Analyze data distribution, trends, and correlations.
- Visualize insights for understanding flight patterns.

## Machine Learning Model Development using Spark MLlib:

- Employ Spark MLlib for predictive modeling.
- Train models to forecast flight delays and cancellations.
- Experiment with diverse algorithms for optimal performance.

# HARDWARE AND SOFTWARE DEVELOPMENT PLATFORMS

## Databricks:

Cloud-based platform providing both hardware infrastructure and software tools.

Primary platform for data processing, analysis, and model development.

## PySpark:

Apache Spark for distributed data processing.

Used for data manipulation, analysis, and visualization.

## Spark MLlib:

Machine learning library for building predictive models within Spark.

Integrates seamlessly with PySpark for streamlined workflow.

# EXPLORATORY ANALYSIS RESULTS: DELAYS BY AIRLINES

- Allegiant Air has the biggest proportion of delayed flights (50.9%).
- The most punctual airline is Endeavor Air (17.9% of delays).



Proportion of Delayed Flights by Airline (%)

# EXPLORATORY ANALYSIS RESULTS: DELAYS BY DESTINATION AIRPORTS

The airports with maximum number of delays are:

- Hartsfield-Jackson Atlanta International Airport (1,259 delays),
- Boston Logan International Airport (1,171 delays),
- Chicago O'Hare International Airport (898 delays).



Number of Delayed Flights by Destination Airport

# EXPLORATORY ANALYSIS: DELAYS BY HOUR

- Flights from 2 pm to midnight have the biggest proportion of delays (33.3-39.6%).



Proportion of Delayed Flights by Hour

# EXPLORATORY ANALYSIS RESULTS: DELAYS BY THE TIME OF THE DAY



Proportion of Delayed Flights by the Time of the Day (%)

Evenings (37.8%) and nights (37.1%) have more flight delays compared to the other times.

# EXPLORATORY ANALYSIS RESULTS: DELAYS BY THE DAY OF THE WEEK



Proportion of Delayed Flights by Day of the Week (%)

Most of the delayed flights are on Friday(32.6%), followed by Thursday(31.6%).

# EXPLORATORY ANALYSIS RESULTS: DELAYS BY MONTH



Proportion of Delayed Flights by Month (%)

July(35.4%) has the most delayed flights followed by June(35.2%).

# EXPLORATORY ANALYSIS RESULTS: DELAYS BY THE SEASON



Proportion of Delayed Flights by Season(%)

Summer is the season with most delays (34.2%).

# EXPLORATORY ANALYSIS RESULTS: CANCELLATIONS BY THE DESTINATION CITY



Cancelled Flights by the Destination City

The destination city with the maximum number of cancellations is New York (187 flights). However, the biggest proportion of the cancelled flights (40%) have Eagle, CO as a destination city.

# EXPLORATORY ANALYSIS RESULTS: CANCELLATIONS BY MONTH



Cancelled Flights by Month

April (463 flights, or 6.7%) has the most cancelled flights followed by March (428 flights, or 5.6%).

# EXPLORATORY ANALYSIS RESULTS: CANCELLATIONS BY THE TIME OF THE DAY



Cancelled Flights by the Time of the Day

Mornings has the most cancelled flights (781 flights) followed by Evening (751 flights). However, the biggest proportion of flights was cancelled at night (3.8%).

# SOME OTHER EXPLORATORY ANALISYS RESULTS

- Reagan Washington Airport and Dulles International Airport have the same proportion of delayed flights (29.3%) and average delay time (45.2-46.2 min).

- Destination airport with the max proportion of delayed flights – CGI and DAB (100%)

- Destination airport with the minimum proportion of delayed flights – AVP, CID, SBN, SCE (0%)

# MACHINE LEARNING RESULTS

| Model | Cancellation Prediction Accuracy | Delay Prediction Accuracy |
|---|---|---|
| Random Forest | 0.95 | 0.69 |
| Logistic Regression | 0.73 | 0.64 |
| XG Boost | 0.94 | 0.68 |
| Gradient Boosting | 0.92 | 0.66 |
| Artificial Neural Network | 0.62 | 0.54 |

# CANCELLATION PREDICTION: BEST MODELS

Accuracy:

- Random Forest – 0.95
- XG Boost – 0.94



Confusion Matrix for Random Forest (flight cancellation)



Confusion Matrix for XG Boost (flight cancellation)



ROC Curve for Random Forest (flight cancellation)



ROC Curve for XG Boost (flight cancellation)

# DELAY PREDICTION: BEST MODELS

Accuracy:

- Random Forest – 0.69
- XG Boost – 0.68



Confusion Matrix for Random Forest (flight delay)



Confusion Matrix for XG Boost (flight delay)



ROC Curve for Random Forest (flight delay)



ROC Curve for XG Boost (flight delay)

# CONCLUSION

- Through extensive data processing and analysis, we gained valuable insights into flight patterns and disruptions in the Washington DC area.

- Random Forest and XG Boost demonstrated high accuracy in predicting flight cancellations and moderate performance in predicting delays.

- The outcomes of our analysis and modeling have significant implications for travelers, airlines, and airport authorities, enabling informed decision-making and proactive measures to mitigate flight disruptions.

# FUTURE WORK

- Optimize parameters and adjust model architecture to improve performance.

- Identify new predictors to capture additional insights from the data.

- Test the model on other locations to find common and specific factors of flight delays and cancellations

- Include factors like weather conditions, holidays, and airport congestion for improved predictive accuracy.

- Systematically experiment with hyperparameters to find optimal configurations.

# REFERENCES

- Balamurugan, R., Maria, A. V., Baranidaran, G., MaryGladence, L., & Revathy, S. (2022, April 1). *Error Calculation for Prediction of Flight Delays using Machine Learning Classifiers*. IEEE Xplore. https://doi.org/10.1109/ICOEI53556.2022.9776709. Retrieved March 12, 2024.

- Federal Aviation Administration. (2023, April 10). *Air Traffic By The Numbers*. https://www.faa.gov/air_traffic/by_the_numbers. Retrieved March 12, 2024.

- Junginger, J. (2023, October 6). *In numbers: The economic impact of flight disruptions*. AirHelp. https://www.airhelp.com/en-int/blog/in-numbers-the-economic-impact-of-flight-disruptions/. Retrieved March 12, 2024.

- Metropolitan Washington Airports Authority. (2024, February 22). *Washington's airports set new passenger record*. https://www.flydulles.com/news/washingtons-airports-set-new-passenger-record. Retrieved March 12, 2024.

- Publisher Research and Innovative Technology Administration. (2023, December 7). *Airline on-time performance and causes of flight delays*. Catalog. https://catalog.data.gov/dataset/airline-on-time-performance-and-causes-of-flight-delays. Retrieved March 12, 2024.

- Zelazko, P. (2023, December 8). *Flight Delay and Cancellation Dataset (2019-2023)*. Kaggle. https://www.kaggle.com/datasets/patrickzel/flight-delay-and-cancellation-dataset-2019-2023. Retrieved March 11, 2024.

# APPENDIX 1. DATASET SAMPLE