

Data 1030 Midterm Presentation

Repository Link

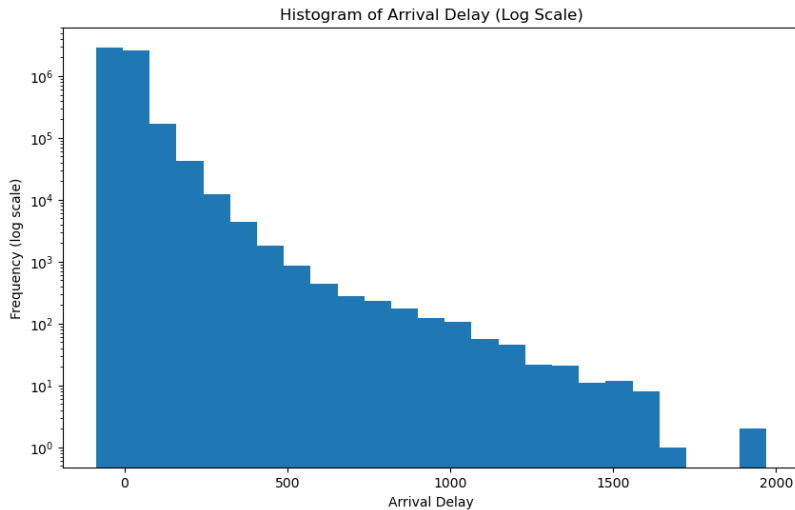
Chai Harsha

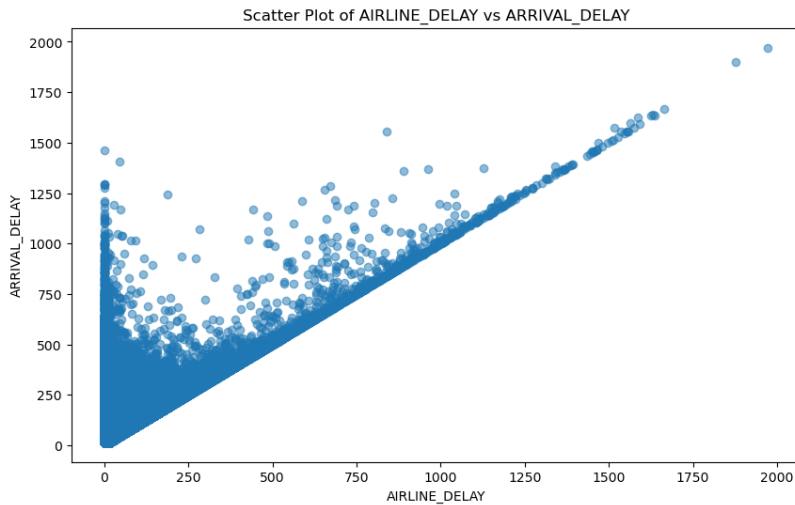
Brown University

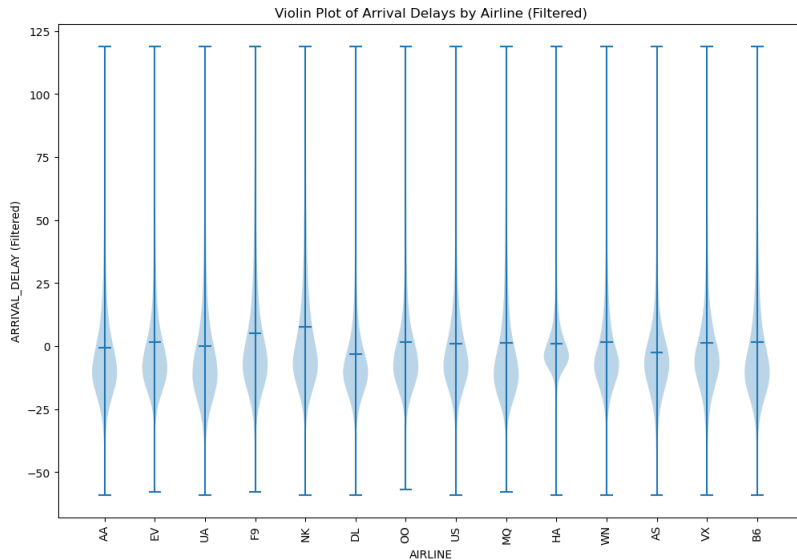
October 20, 2025

Flight Data

- Thousands of flights + thousands of airports = vast potential for delays
- Airport design, airline inefficiencies, etc
- Can we predict delays?
- USDOT Bureau of Transportation Statistics collects on-time performance of major carriers
- [Kaggle dataset - 2015 flight delays](#)
- This is a very large dataset (>5 million rows) and also has many missing values
- Continuous data and some categorical data
- There are other factors not in this database - weather being the obvious one







Preprocessing

- 105071 rows (1.81% of dataset) missing the target feature - dropped
- Too many airports - selected top 20 + “Other”
- Standard scaler on continuous data
- Minmax scaler on time-of-day data
- One-Hot encoded categorical data
- Dropped columns like cancellation reason and air system delay for too many missing data ($> 80\%$)