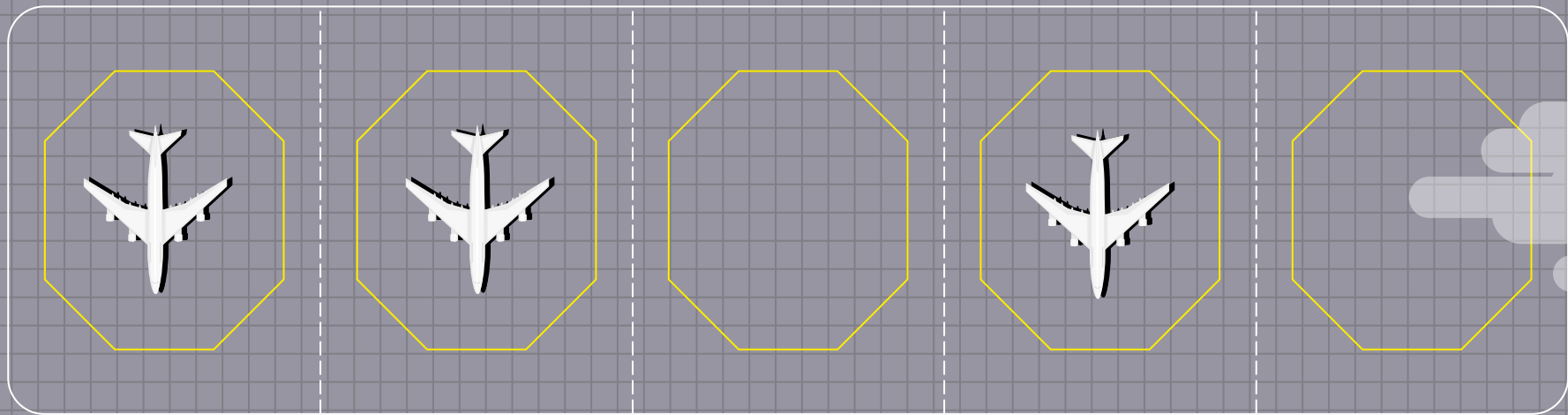


Will My Flight Be Late?

Simon Guichandut, Tim Hallatt, Ketan Sand

Erdos Datascience Bootcamp - Fall 2023



Problem, Stakeholders and KPIs

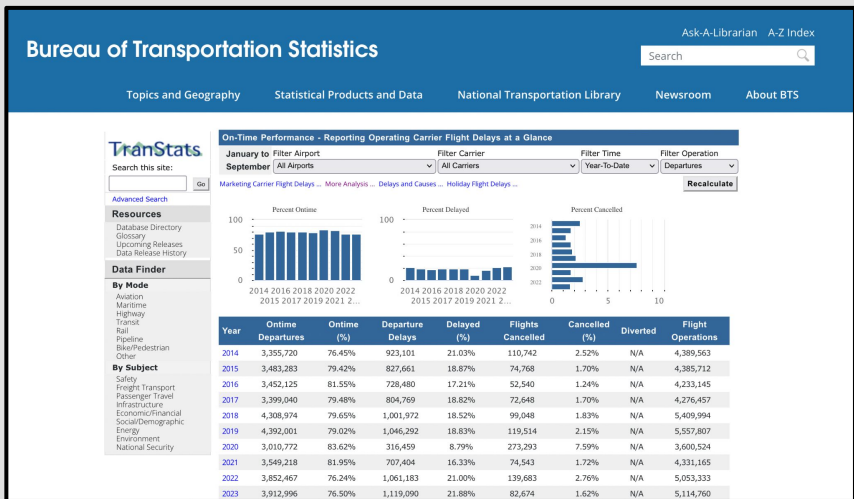
In a single year, there are 200,000 hours of combined flight delays in just the 20 of busiest airports of the United states. Due to all this the US economy suffers a \$32.9 billion annual loss.

Stakeholders - Travelers, Airlines, Airport Authorities, Government (Tourism Sectors), Insurance companies.

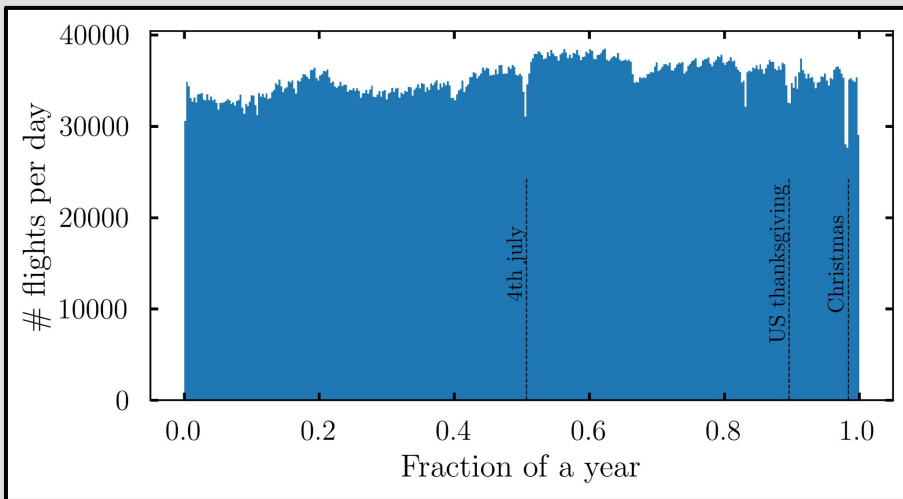
KPIs - Precision/Recall: False Positive vs. False Negative rates, Receiver Operating Characteristic (ROC), Detection Error Tradeoff (DET)



Data Collection



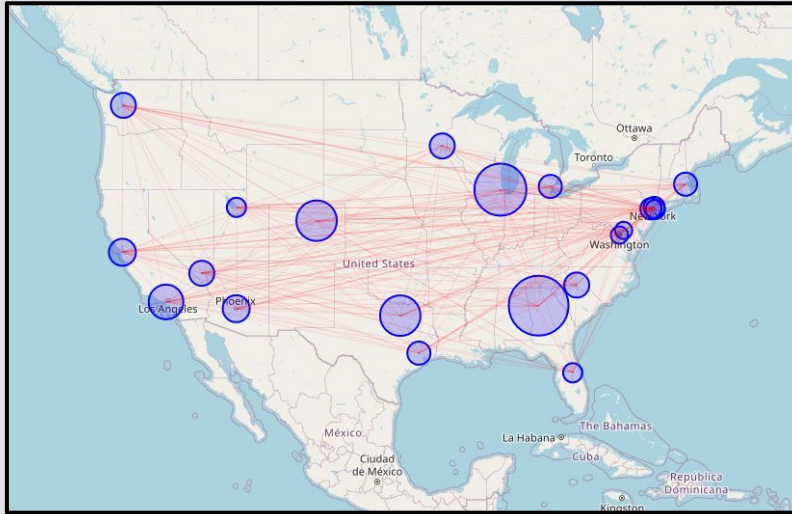
US Bureau of transportation statistics
Data - 2012 - 2023
Automated using Selenium



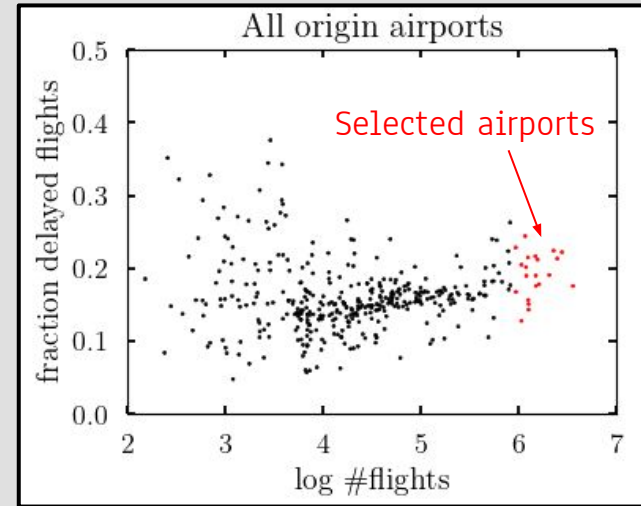
More than 30,000 flights per day!
Dips on the Holidays



Cleaning



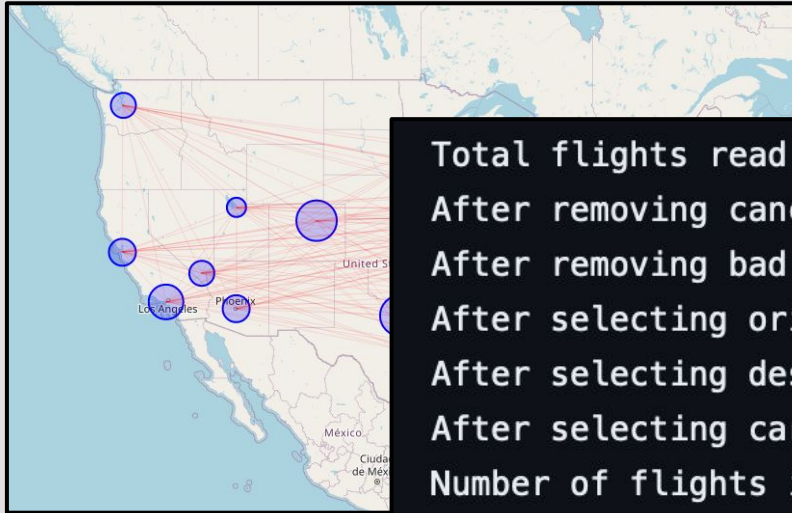
20 Busiest Airports
Removed Cancelled Flight
Top 8 Aircraft Carriers



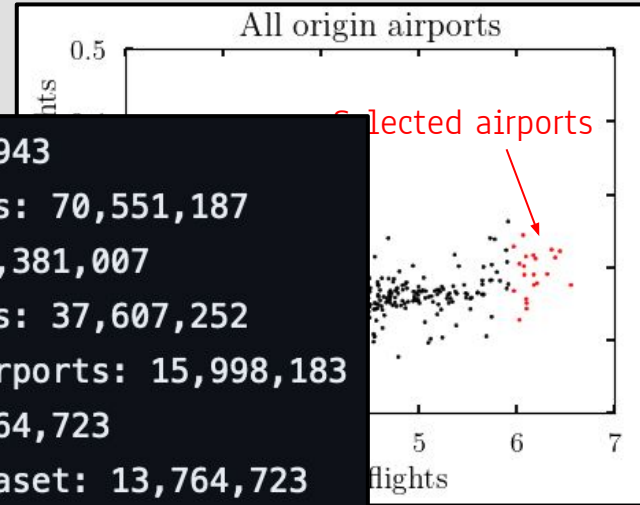
The smallest airports have the most delays: to avoid bias, we take the 20 largest airports to balance number of flights (20% of total) with delays



Cleaning



```
Total flights read in: 71,983,943
After removing canceled flights: 70,551,187
After removing bad columns: 70,381,007
After selecting origin airports: 37,607,252
After selecting destination airports: 15,998,183
After selecting carriers: 13,764,723
Number of flights in final dataset: 13,764,723
```



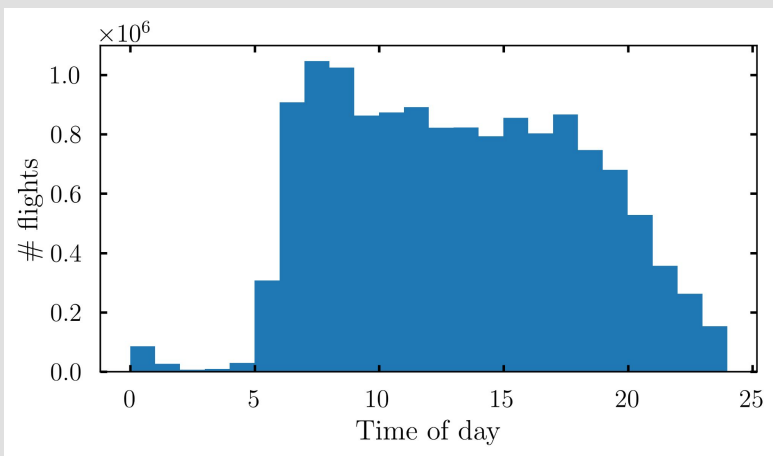
20 Busiest Airports
Removed Cancelled Flight
Top 8 Aircraft Carriers

The smallest airports have the most delays: to avoid bias, we take the 20 largest airports to balance number of flights (20% of total) with delays

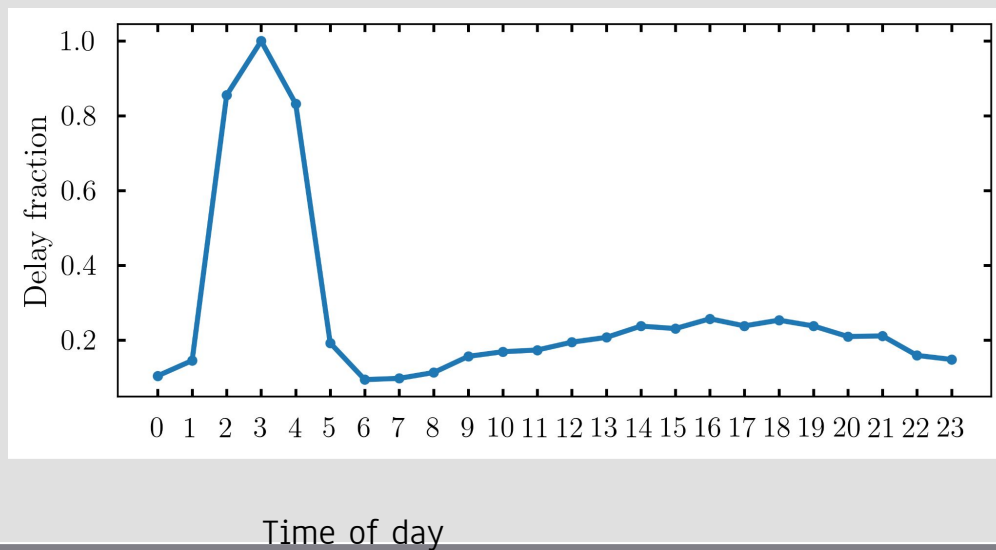


Exploratory Plots

Number of flights

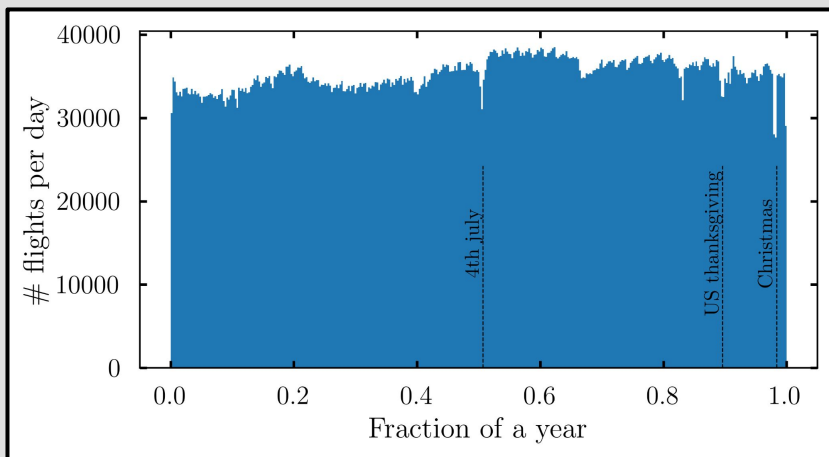


Delay fraction

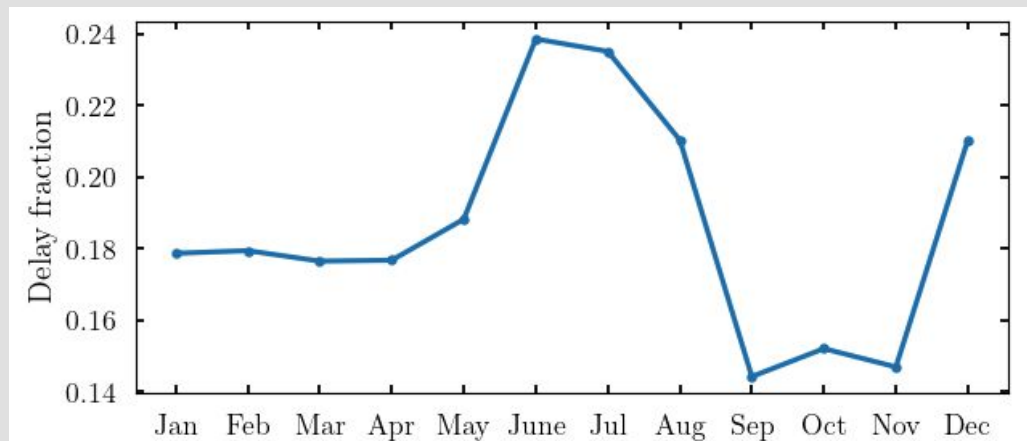


Exploratory Plots

Number of flights

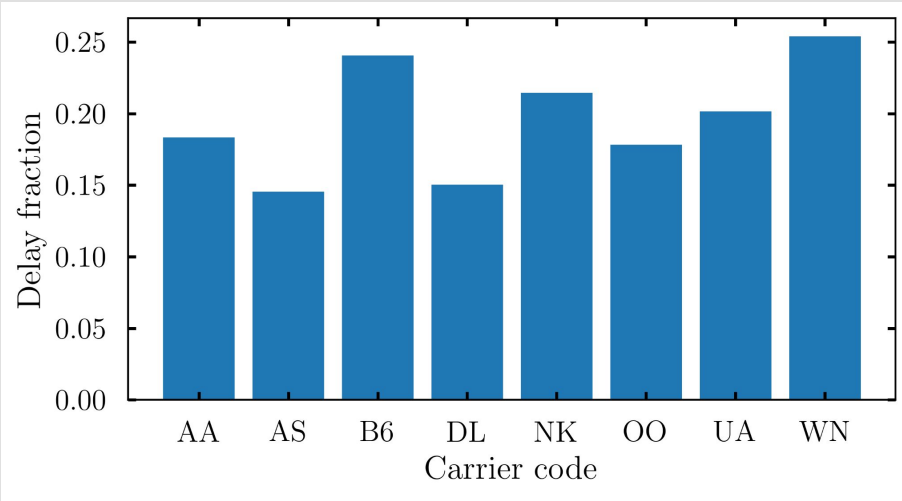


Delay fraction

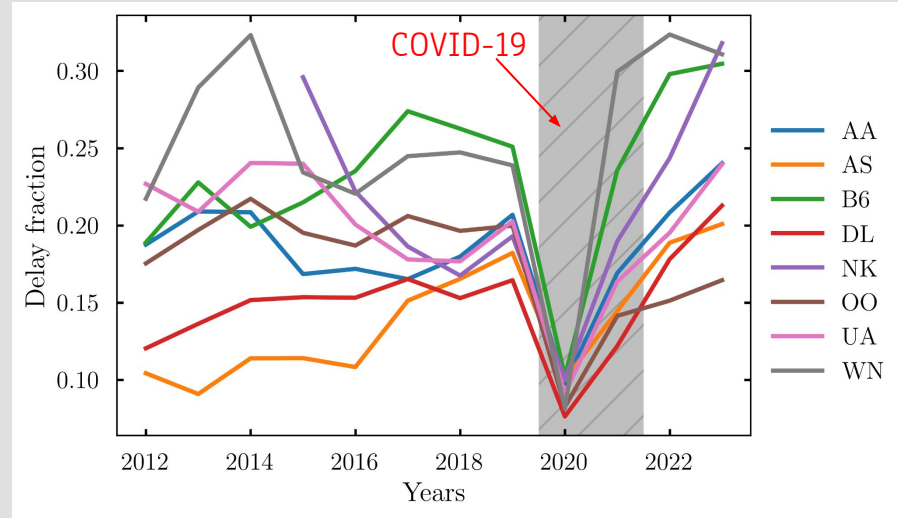


Exploratory Plots

Carrier matters: Southwest (WN)
worst, Alaskan Airlines best (AS)

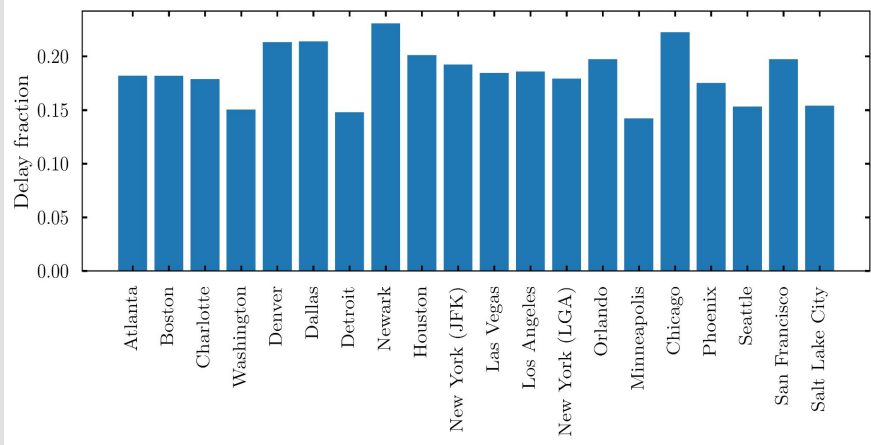


Jetblue (B6) + Southwest (WN) grow worse post-COVID

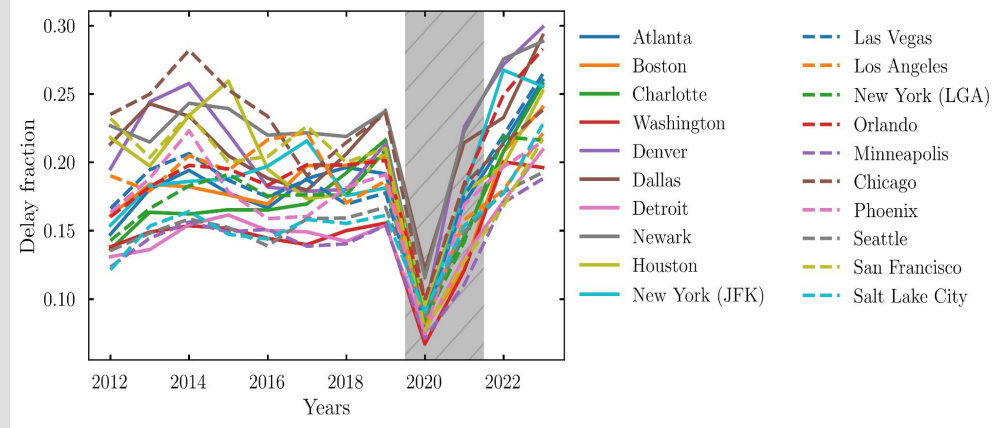


Exploratory Plots

Airport matters, weakly

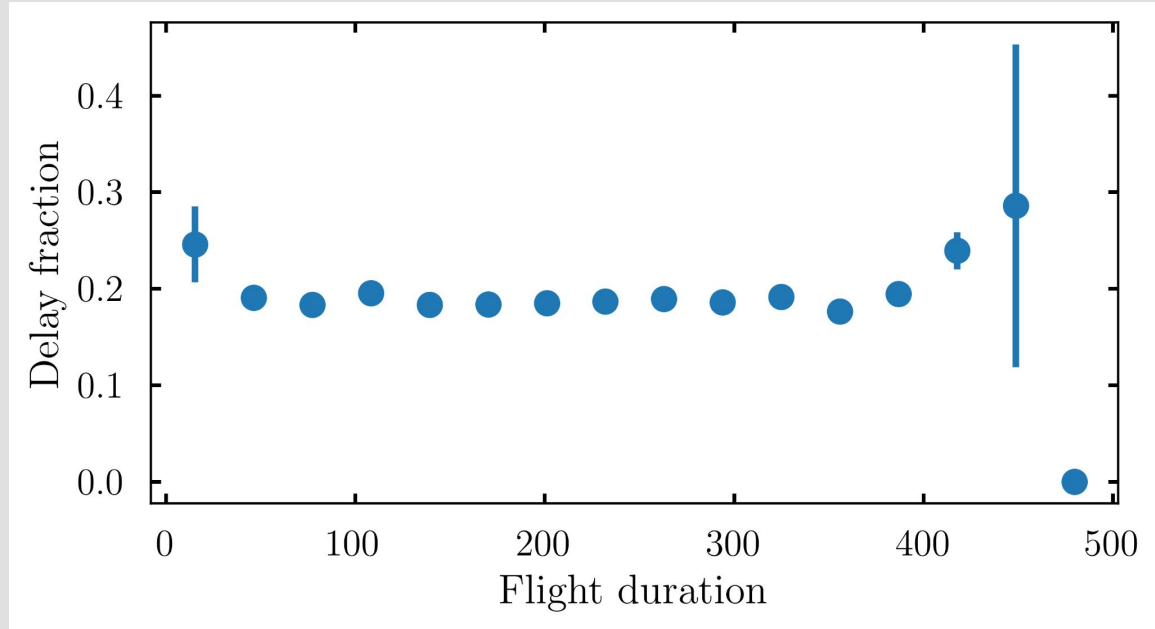


Dallas, Newark, Orlando are bad post-COVID



Exploratory Plots

Flight duration correlates weakly with delay fraction



Workflow

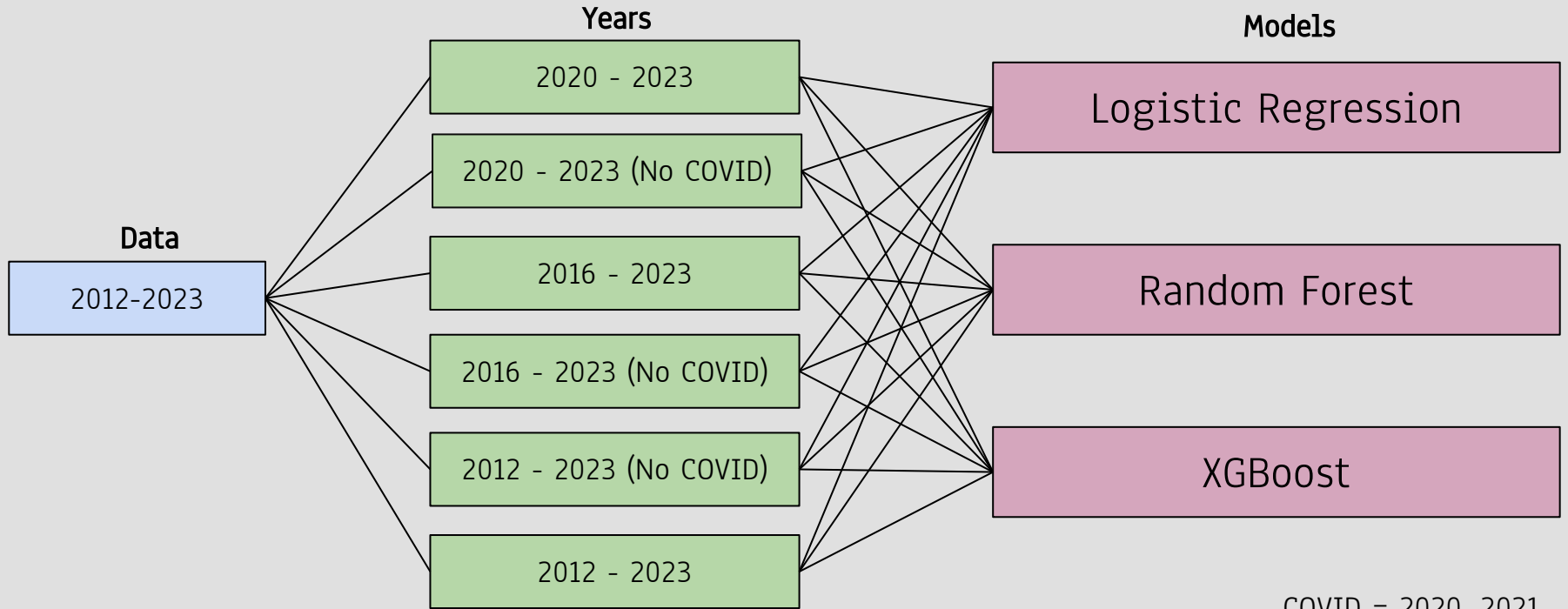
Preparing the data

- **Predicted variable:** Delay times >15min are categorized as TRUE
- **Scaling:** Normalize flight duration to 0-1
- **Categorical variables:** One-hot encode carrier and airport variables
- **Cyclic variables:** Convert day & year fraction to sin & cos components

$$\begin{aligned} \text{day_frac_sin} &= \sin(2\pi \cdot \text{day_frac}) & \text{year_frac_sin} &= \sin(2\pi \cdot \text{year_frac}) \\ \text{day_frac_cos} &= \cos(2\pi \cdot \text{day_frac}) & \text{year_frac_cos} &= \cos(2\pi \cdot \text{year_frac}) \end{aligned}$$



Workflow



COVID = 2020, 2021

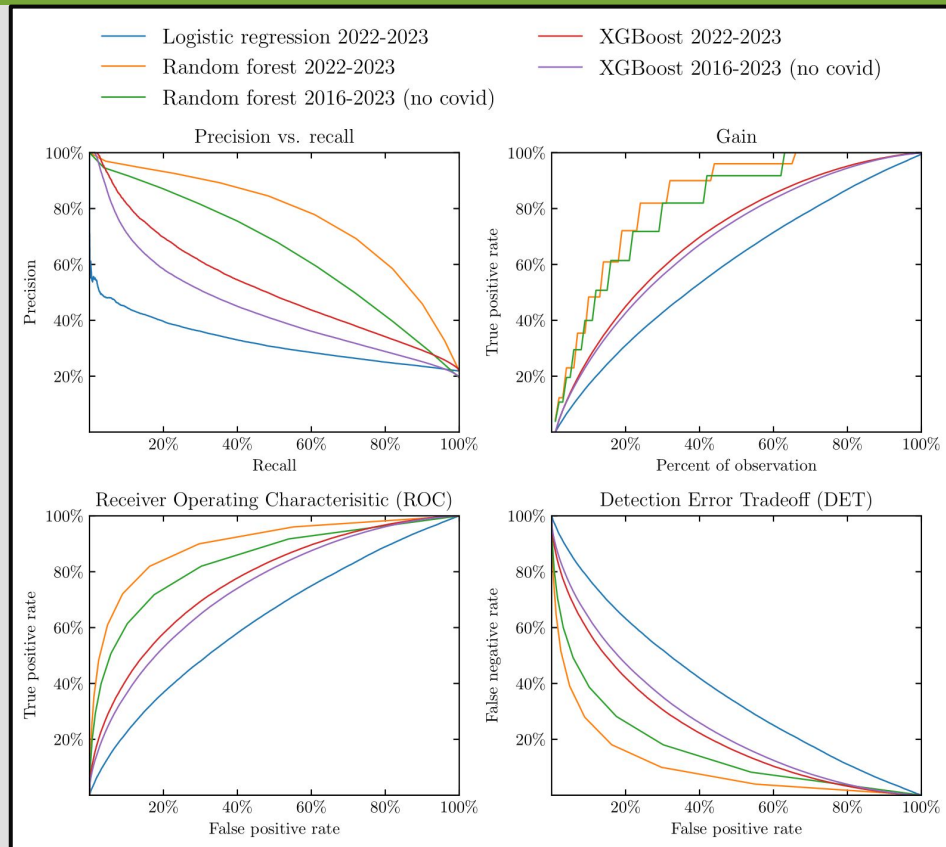


Best Result

Train on 2022 - June 2023
Model - Random Forest

Testing On
July and August 2023

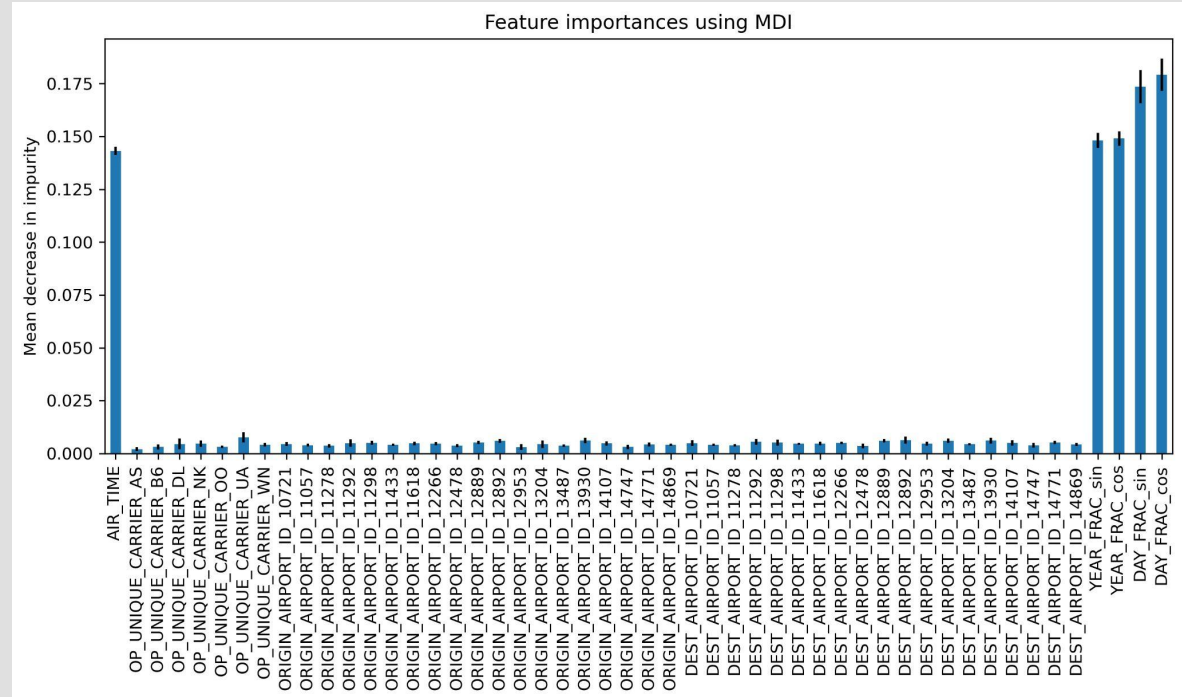
On delayed flight category	Precision	Recall	F1-score	Accuracy
Validation set	0.69	0.72	0.71	0.87
Testing set	0.5	0.7	0.58	0.72



Feature Importance

What determines if your flight will be late?

1. Time of day
2. Time of year
3. Flight duration



Delay Time Range*

Logistic Regression

Accuracy: 0.41

Classification	Report: precision	recall	f1-score	support
1.0	0.40	0.50	0.44	28615
2.0	0.37	0.00	0.01	24765
3.0	0.41	0.65	0.51	29892

Random Forest

Accuracy: 0.64

Classification	Report: precision	recall	f1-score	support
1.0	0.65	0.71	0.68	28615
2.0	0.60	0.51	0.55	24765
3.0	0.67	0.70	0.68	29892

XGBoost

Accuracy: 0.57

Classification	Report: precision	recall	f1-score	support
1.0	0.57	0.65	0.61	28615
2.0	0.53	0.37	0.43	24765
3.0	0.60	0.67	0.63	29892

1 -> 15 mins < Delay time <= 30 mins

2 -> 30 mins < Delay time <= 60 mins

3 -> Delay time >= 60 mins

***OPTIMIZATION NEEDED!**



Web-App

Will My Flight Be Late?

Flight carrier: Delta Air Lines Inc. Flight Date: 2023/06/15

Origin Airport: Newark, NJ: Newark Liberty Internati... Departure Time: 03:00

Destination Airport: New York, NY: John F. Kennedy Intern... Flight duration (hours): 10.50 (range: 0.50 - 12.00)

Evaluate

YES!

Our model thinks your flight will be more than 15 minutes late

Statistics about your flight:

- Over the last 10 years (but excluding pandemic years):
- 14% of flights from your carrier have been delayed (20% in 2023 only)
- 24% of flights from your origin have been delayed (27% in 2023 only)
- 21% of flights to your destination have been delayed (23% in 2023 only)

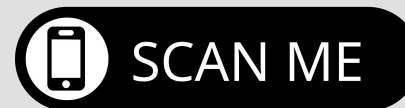
About this app

Current model: Random Forest with 10 estimators

By [Ketan Sand](#), [Simon Guichandut](#) and [Tim Hallatt](#) (links to our linkedin).

We thank the Erdős Institute and our mentor Gleb Zhelezov.

<https://willmyflightbelate.streamlit.app/>



Future Work

- Calibration of model: predict probability, rather than Yes/No
- Delay range prediction: Optimize and Add to the Website
- Hyperparameter tuning: RandomForest takes too long
 - n_estimators, max_features, max_depth, min_samples_split, min_samples_leaf, bootstrap
 - Using RandomForestRegressor and RandomizedSearchCV
 - Notebook in repo
- Train Model that also provides the cause of the delay
- Get more parameters and apply neural networks

Special thanks to our mentor Gleb Zhelezov!!

