

Team: RooKeys Anudeep Arora, Oussama Landoulsi, Lalit Yadav



Predicting Flight Departure Delay

The Erdős Institute, Fall 2022 Bootcamp

Problem Statement & Business impact

Problem: Given a 4-hour period horizon and the scheduled flight data together with the weather data, will a flight be delayed by more than 15 minutes?

Business position: According to the Federal Aviation Administration, the total flight delays cost is \$22 billion yearly.

Target audience: Airline companies.

<u>Goals</u>: Predicting flight delays to improve airline operations by:

- Saving on delays costs and providing quality service to their customers.
- Reducing potential crowding situation at the airport.



Data gathering

Airlines Dataset: Jan 2018 - Aug 2022

- Reporting Carrier On-Time Performance Bureau of Transportation Statistics.
- 109 features in total. Extracting the data for JFK as an origin.
- Reducing usable features to 17.

Weather Dataset: Jan 2018 - Aug 2022

- Web Scraping daily/hour weather data from weather-underground website.
- Selenium library with WebDriver.

Daily	00361	valion	Э						
Time	Temperature	Dew Point	Humidity	Wind	Wind Speed	Wind Gust	Pressure	Precip.	Condition
12:43 AM	56 °F	55 °F	97 %	SSE	18 mph	26 mph	30.22 in	0.1 in	Rain
12:51 AM	td.mat-cell.cdk-cell.	.cdk-column-humidi	ty.mat-column	-humidity.ng-:	star-inserted 90.15	× ³⁰ ph	30.21 in	0.2 in	Light Rain
1:12 AM	56 °F	55 °F	97 %	SSE	15 mph	0 mph	30.21 in	0.0 in	Light Rain
1:45 AM	55 °F	55 °F	100 %	S	16 mph	0 mph	30.22 in	0.1 in	Rain
1:51 AM	56 °F	55 °F	97 %	S	15 mph	0 mph	30.22 in	0.1 in	Light Rain
2:51 AM	55 °F	54 °F	96 %	S	16 mph	0 mph	30.21 in	0.0 in	Cloudy
🕞 🗘 Inspe	cteur 🕥 Console	e D Débogueur	r †↓ Résea	u {} Édite	ur de style 🛛 Pe	rformances 🕕	Mémoire 🗄	Stockage 🕇	Accessibilité
Q Rechercher dans le HTML + <pre></pre>						<pre></pre>			
html ≻ body.en	able-sda					color: • #	Le2023; ▶○#fff:		

Daily Observations



- Highly imbalanced dataset.
- Merge flight and weather data with non uniform time stamps.
- Capturing seasonality and drift for the carrier specific delay.



Exploratory Data Analysis



Carrier delay exhibits **seasonality** and **drift**





Exploratory Data Analysis (contd.)



Exploratory Data Analysis (contd.)



Feature Engineering

- Predetermined information:
 - Hour, Day and Month of the flight's scheduled departure.
- Weather information:
 - temperature, dew point, humidity, wind, wind speed, wind gust speed, pressure, precipitation and condition.
- Historical information:
 - Delayed flights at the same hour, a day before the target flight's departure.
 - Probability of delay for a given carrier (past performance).
 - Probability of delay for a given destination (historical information).
 - Delayed flights (%) in 1 hours time window at T 4 hour for a scheduled departure at time T

Feature Engineering (contd.)

Encoding:

- Normalization using min-max scaler
 - Delayed flights (%) at T 4 hour, Temperature, Dew Point, Humidity, Wind Speed, Wind Gust, Pressure, and Precipitation.
- Cyclic embedding
 - Day of Week, Month, Hour, and Wind Direction.
- One-Hot encoding
 - Weather Condition

Total features to pass into the model = 71

$$F_i(t) \Rightarrow [\sin\left(\frac{2\pi t}{T}\right), \cos\left(\frac{2\pi t}{T}\right)]$$

[►] {"N","NNE","NE" ...}→degrees→cyclic coordinates

Modeling Framework

KPI for stakeholders:

- F2 score (weighing recall over precision) Labels:
- Not classifying a delay correctly is costlier





Model Performance

1.0

Feature Importance



If we remove "Delaym4h_" feature, the prediction horizon would increase to 24 hours. What would be the performance of model?

- Prediction Horizon: 24 hour
- Model: Random Forest
- Precision: 0.39
- Recall: 0.82
- F2 score: 0.67

Conclusion & Future Work

Conclusion:

- Model achieved 80% ability in predicting a delay.
- Feature engineering is critical for model performance.

Future Work:

- More general model including all airports as an origin.
- Delay linked to the airlines organization:
 - Understanding the spread of the delay for all origin-destination pairs.
 - Predict the evolution of the delay in real time.
 - Probabilistic Graphical Models (PGMs), e.g. Bayesian Network.
 - Incorporate weather information at Destination airport.

Special thanks to Roman, Lindsay, Matt, Shuvra, Alec and everyone in the Erdős Institute. Thank you for listening!