



## Flight Delay Analysis

By: Reece Boyce



# Business Problem



## Increased Expense for Crew, Fuel, and Maintenance

Delayed flights cause more expenses for relocating staff  
Fuel expenses increase due to idle time



## Cost of Delays on the Industry

Estimated cost of \$8 billion per year for the industry\*  
Estimated cost of \$17 billion per year for passengers\*



## Decreased Customer Satisfaction

Airlines depend on repeat customers and word of mouth  
Competition is fierce in the industry, airlines that experience delays frequently lose business.



## Customers Do Not Like Uncertainty

Delayed flights cause stress and anxiety for travelling, especially for those customers traveling with connecting flights

\*According to 2010 study commissioned by the Federal Aviation Administration

# Data Gathering

## Full Dataset:

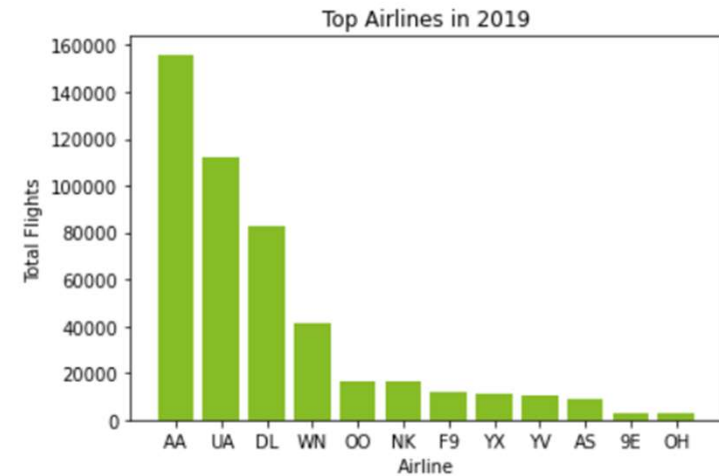
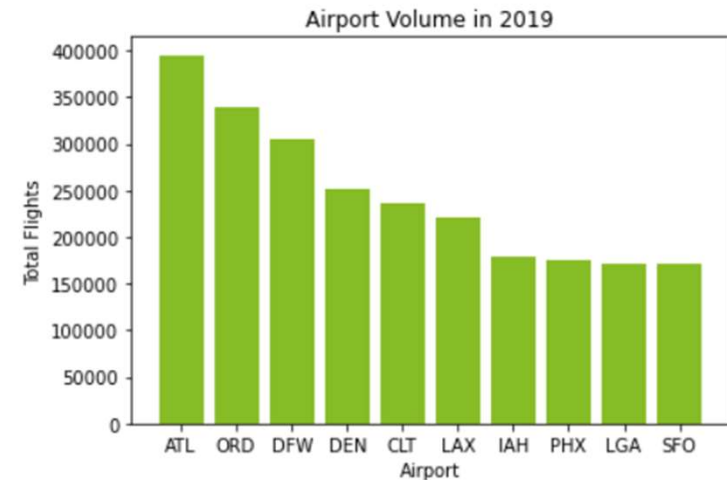
- Source: US DOT Bureau of Transportation Domestic Flight Data
- Includes Data from 2009 – 2019
- +4 million flights per year (40 million in total)

## Selected Data:

- Top domestic carriers that travel nationwide
- Top 10 airports by volume
- 472,000+ data points used in analysis
- 2019 data used in analysis

## Variables:

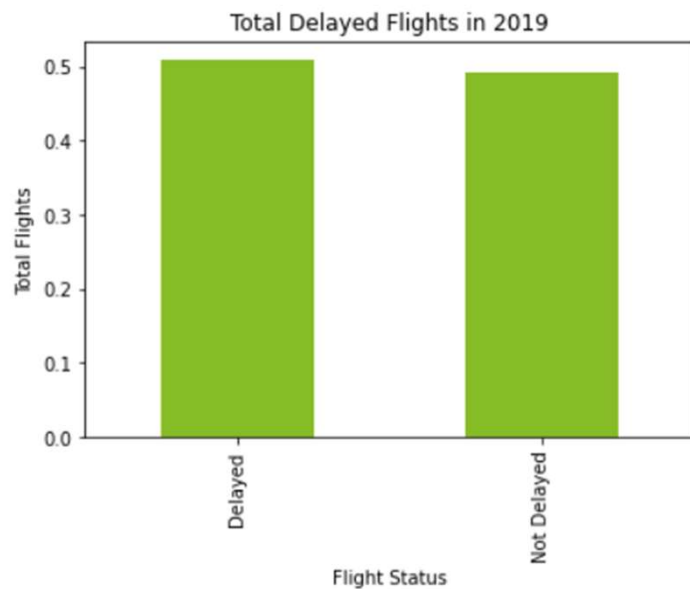
- |                            |                       |
|----------------------------|-----------------------|
| • Year                     | • Origin              |
| • Month                    | • Destination         |
| • Day of the Week          | • Departure time      |
| • Carrier                  | • Time to Taxi out/in |
| • Flight Number            | • Wheels on/off       |
| • Departure/ Arrival times | • Distance            |



# Model Selection



Target Variable: Delay



Problem: Binary Classification

	Computational Needs	Advantages
Logistic Regression	<ul style="list-style-type: none"><li>• Training time: 1.92 seconds</li><li>• Prediction time: 0.01 seconds</li></ul>	<ul style="list-style-type: none"><li>• Less tuning needed, easier to perform on large datasets.</li></ul>
Random Forest	<ul style="list-style-type: none"><li>• Training time: 41.26 Seconds</li><li>• Prediction time: 0.77 Seconds</li></ul>	<ul style="list-style-type: none"><li>• Can handle missing values.</li><li>• Performs better with more explanatory variables.</li></ul>
XGboost	<ul style="list-style-type: none"><li>• Training Time: 17.05 Seconds</li><li>• Prediction Time: 0.32 Seconds</li></ul>	<ul style="list-style-type: none"><li>• Less randomization than RF</li><li>• More parameters than other two methods</li></ul>

# Modeling

## Logistic Regression

True Delays: 47,840  
True Delays Predicted: 30,498

Precision Score: 0.672



Recall Score: 0.637



Accuracy Score: 0.658



F1 Score: 0.654



## XGboost



True Delays: 47,840  
True Delays Predicted: 33,064

Precision Score: 0.757

Recall Score: 0.691

Accuracy Score: 0.731

F1 Score: 0.723

## Random Forest

True Delays: 47,840  
True Delays Predicted: 31,717



Precision Score: 0.663



Recall Score: 0.662



Accuracy Score: 0.659



F1 Score: 0.663

# Recommendations and Way Forward

