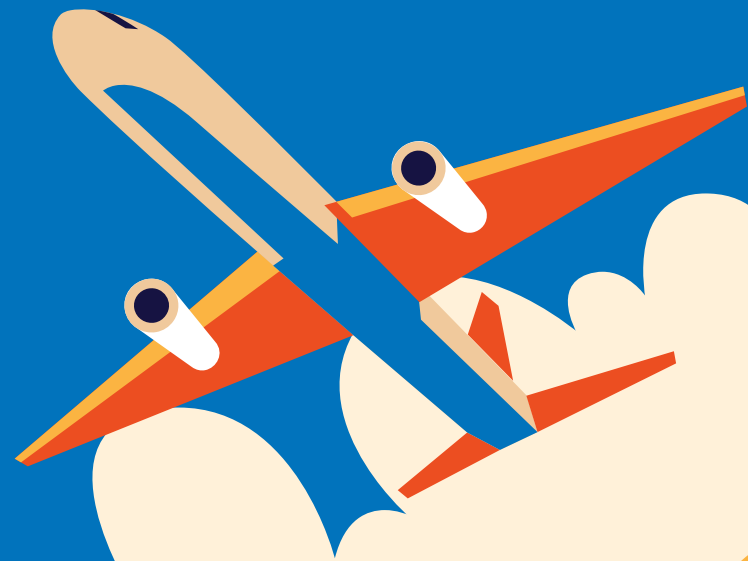



# Estimating Flight Delays

DJ Ammirato, Gene Lam





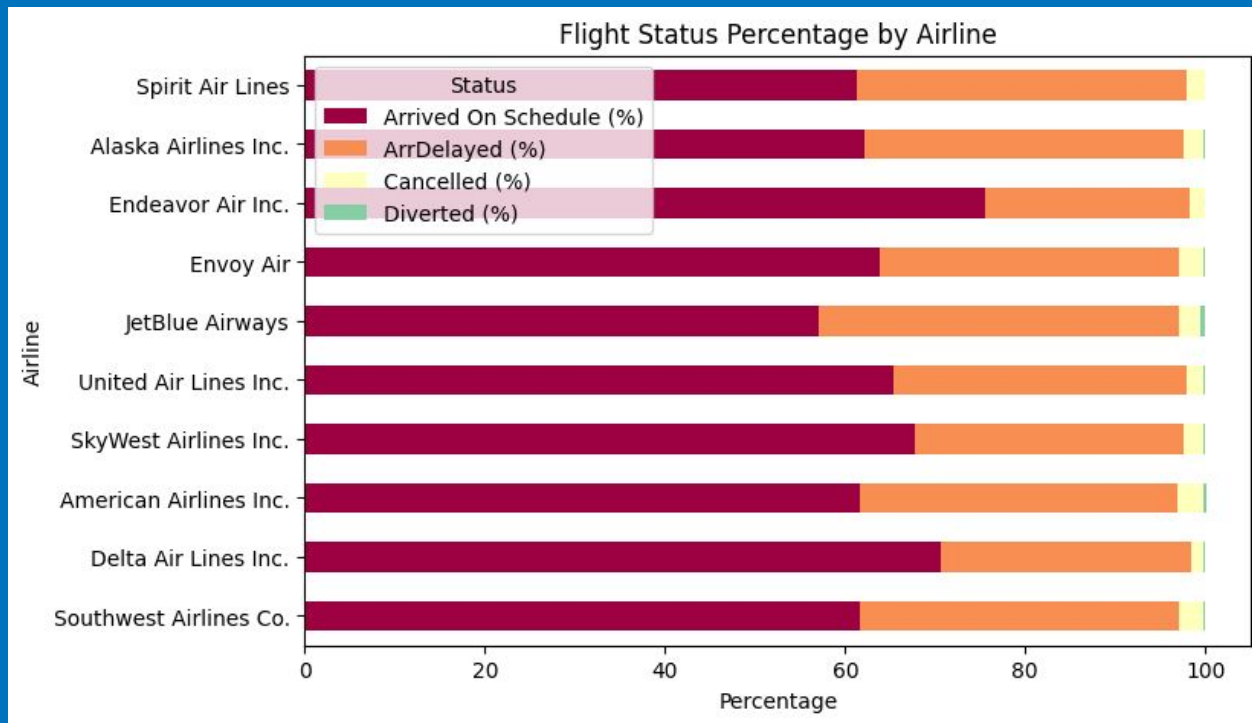
# Objectives

- Can we successfully predict whether or not future flights are going to be delayed? By how long?
  - What features influences whether or not a flight is going to be delayed?
  - Empower consumers to make more informed decisions when buying flight tickets
  - Data comes from Kaggle and Bureau of Transportation Statistics
    - 2018 - 2023 US domestic flights (including US territories)
    - 39M rows, 61 features
    - We sample 100,000 rows and use 28 features
- 



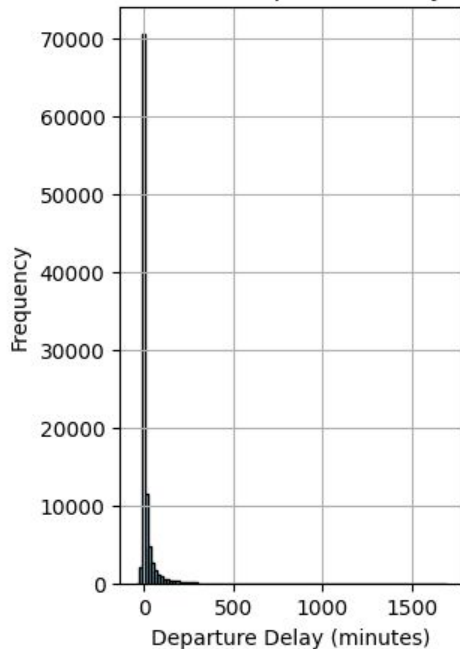
# Exploratory Data Analysis

# Top 10 Airlines in America

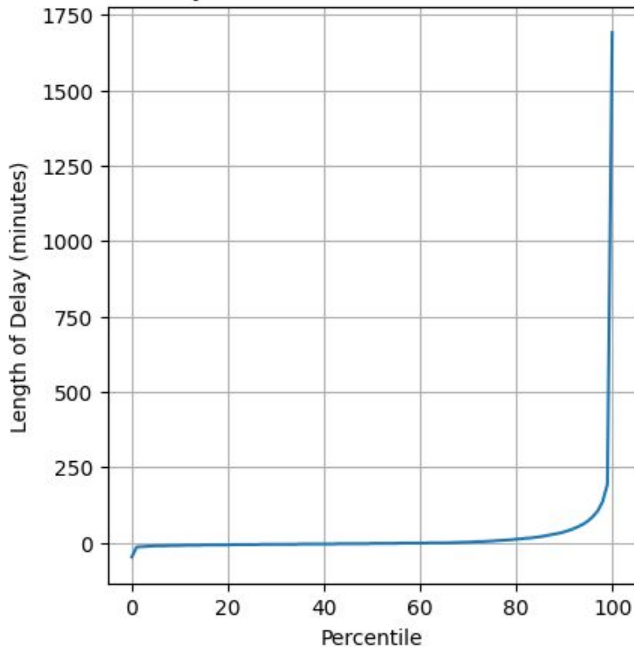


# How long are we waiting?

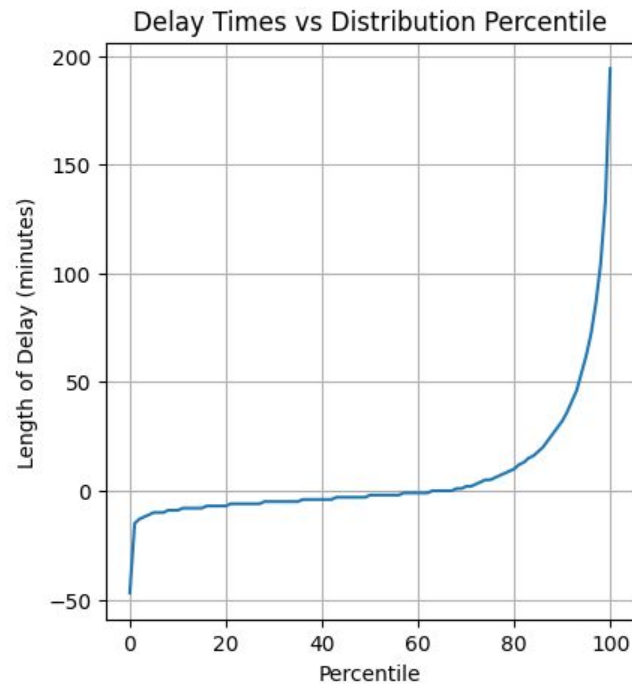
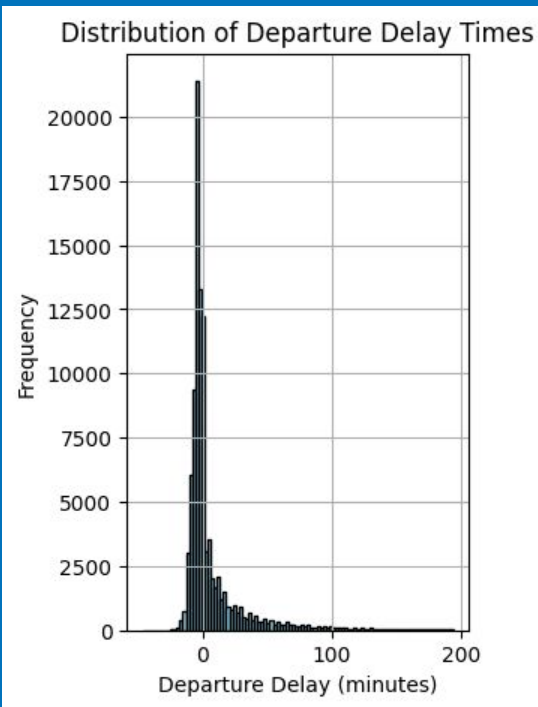
Distribution of Departure Delay Times



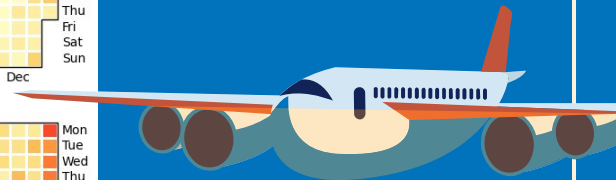
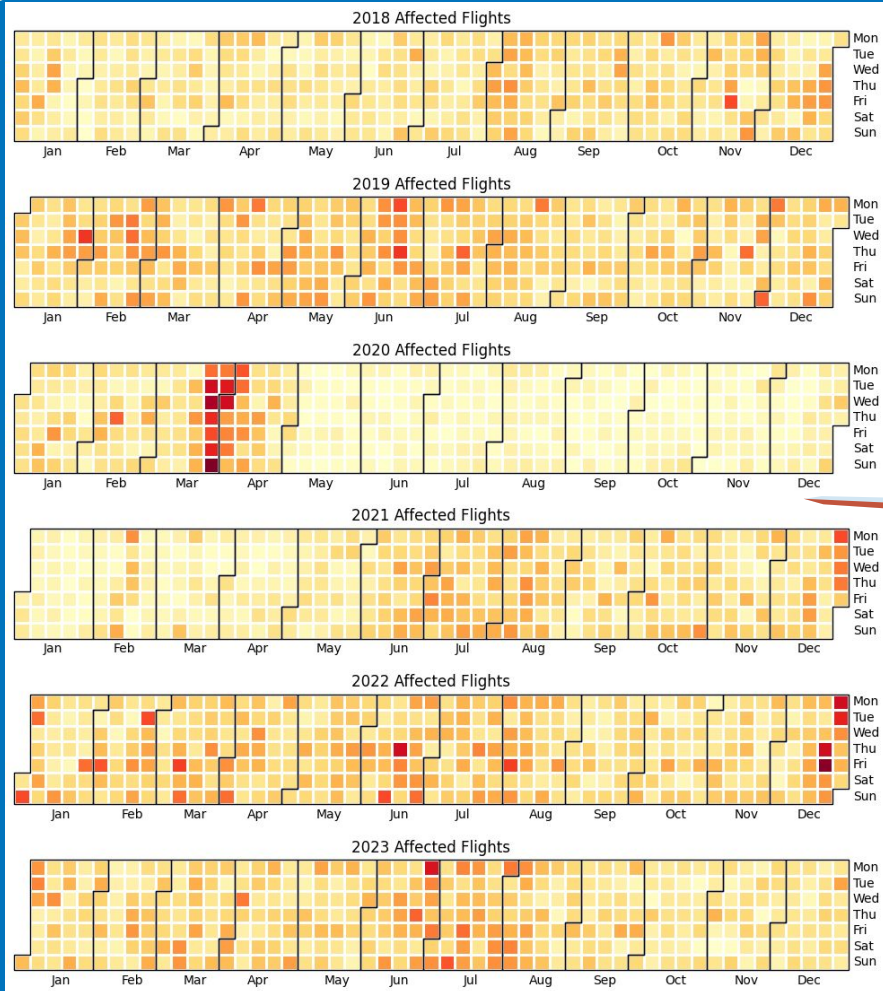
Delay Times vs Distribution Percentile

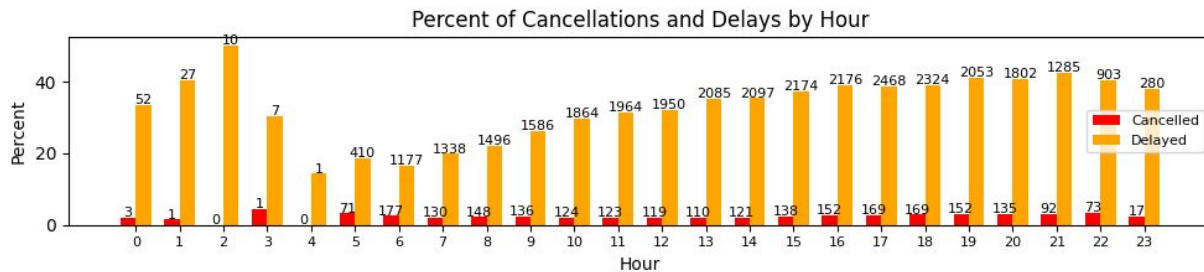
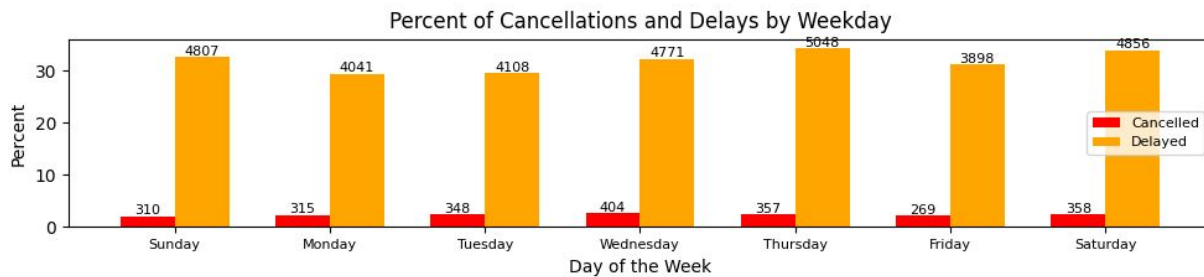
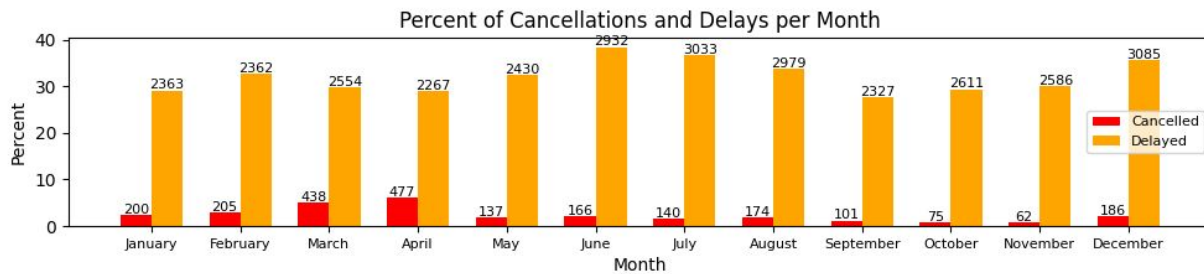


# How long are we actually waiting?



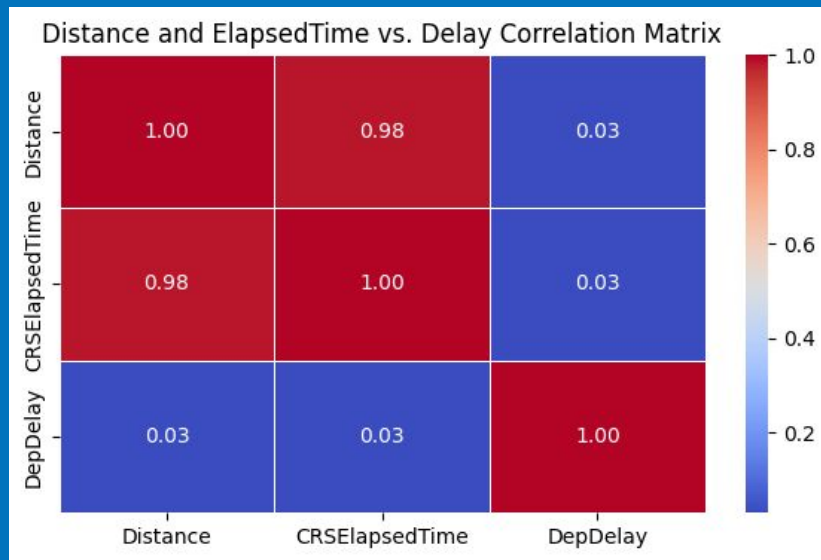
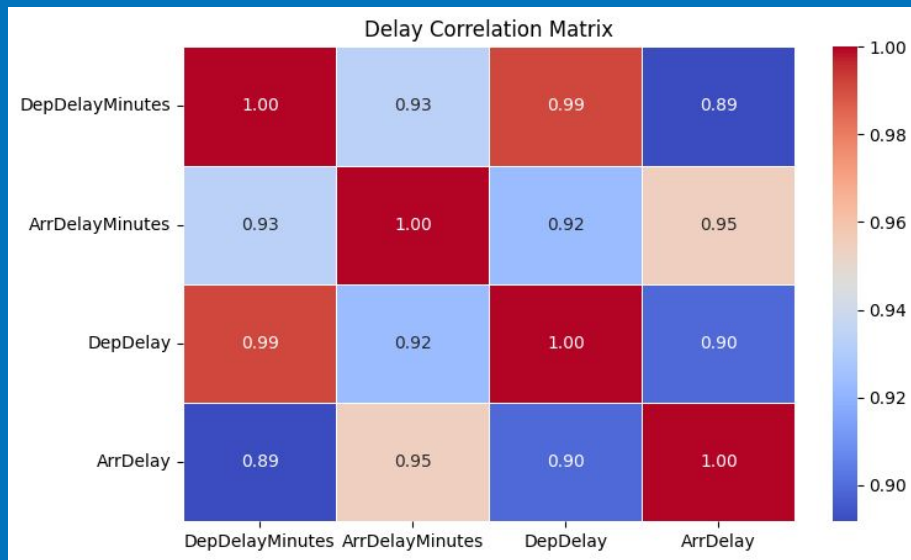
When  
are we  
waiting?







# Correlation of Numeric Features



The background is a solid blue sky. In the top left corner, there is a large, stylized yellow cloud. In the top right corner, there is a large yellow sun partially obscured by a smaller yellow cloud. In the bottom left and bottom right corners, there are more stylized yellow clouds. Three dark blue birds are flying in the sky: one in the upper middle, one on the left side, and one on the right side. The word "Modeling" is written in a large, bold, yellow sans-serif font in the center of the image.

# Modeling

# Performance Metrics

## Classification

- Recall: percentage of positive classifications correctly defined
  - Most important for value proposition
- Switched to F1 due to poor precision

## Regression

- Mean-Squared Error: penalize larger errors to achieve high precision
- R2 score: conveys proportion of variance in delays that is conveyed by our features



# Classification

# Logistic Regression

**59.59%**

Training Accuracy

**64.19%**

Training Recall

**24.91%**

Training Precision

**58.73%**

Testing Accuracy

**63.66%**

Testing Recall

**24.35%**

Testing Precision

# Random Forest

**66.96%**

Training Accuracy

**67.95%**

Training Recall

**30.44%**

Training Precision

**63.80%**

Testing Accuracy

**57.53%**

Testing Recall

**26.12%**

Testing Precision



# Bayesian Optimization

- A method of automatically tuning the hyperparameters of a model
- Fits a series of models and compares their relative performance

# Optimized Logistic Regression

**59.56%**

Training Accuracy

**64.23%**

Training Recall

**24.90%**

Training Precision

**58.66%**

Testing Accuracy

**63.66%**

Testing Recall

**24.31%**

Testing Precision



# Random Forest

**74.03%**

Training Accuracy

**76.52%**

Training Recall

**38.29%**

Training Precision

**65.93%**

Testing Accuracy

**53.48%**

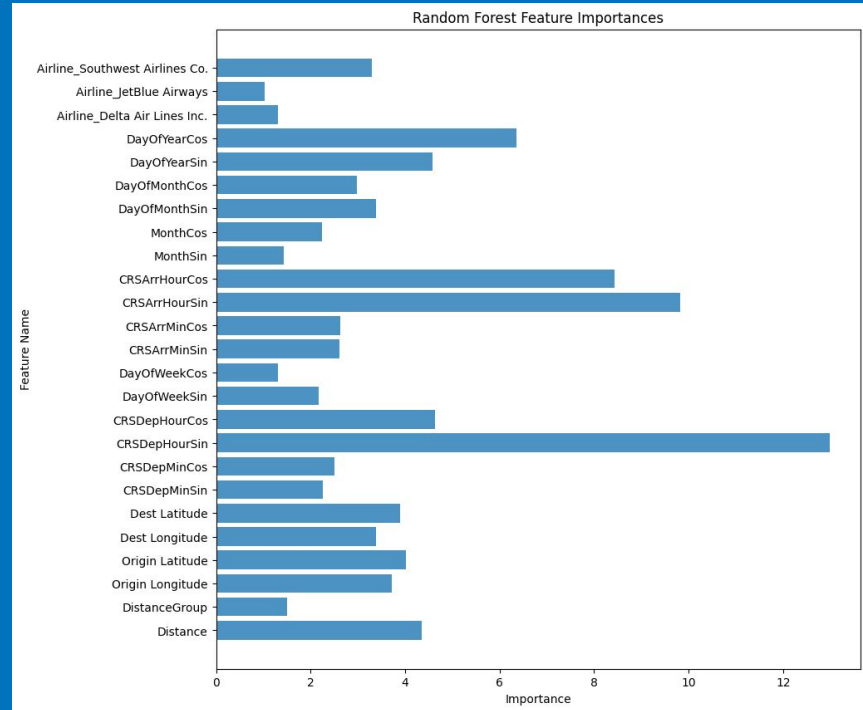
Testing Recall

**26.81%**

Testing Precision

# Feature Importance

- Hour of scheduled flight departure and arrival
- Day of the year





# Regression

# Linear Regression

Training

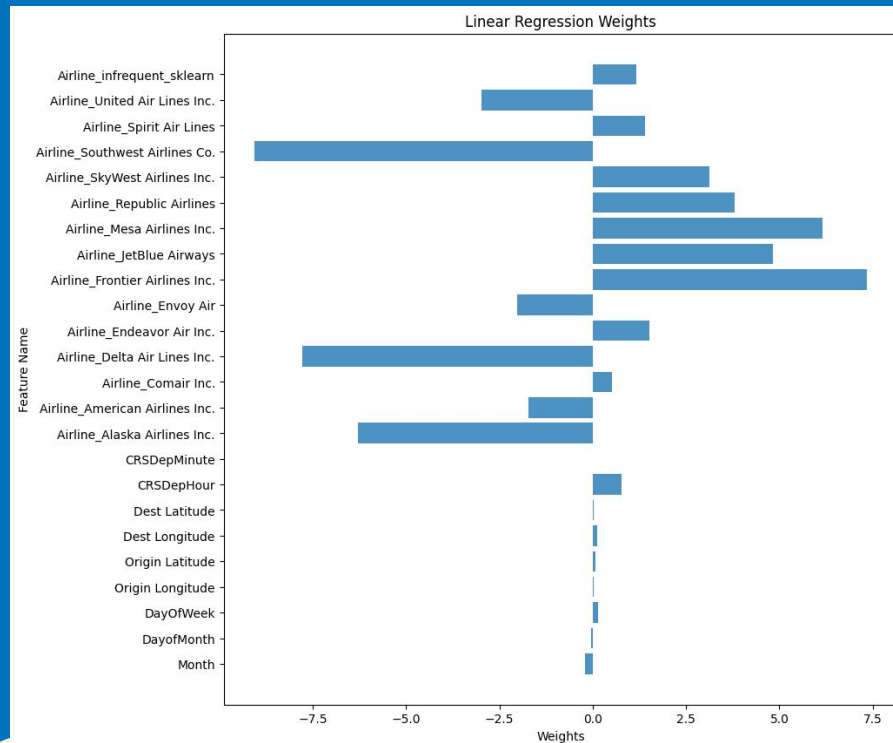
MSE: 1299.71

R2: 0.04

Test

MSE: 1299.71

R2: 0.04



# Optimized Linear Regression

Training

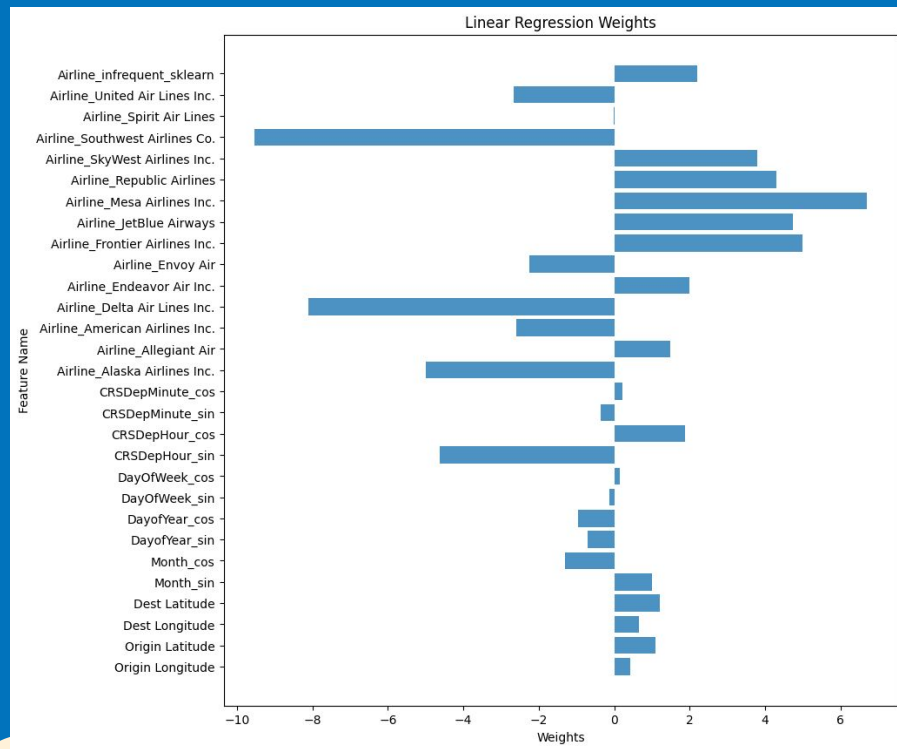
MSE: 1317.68

R2: 0.04

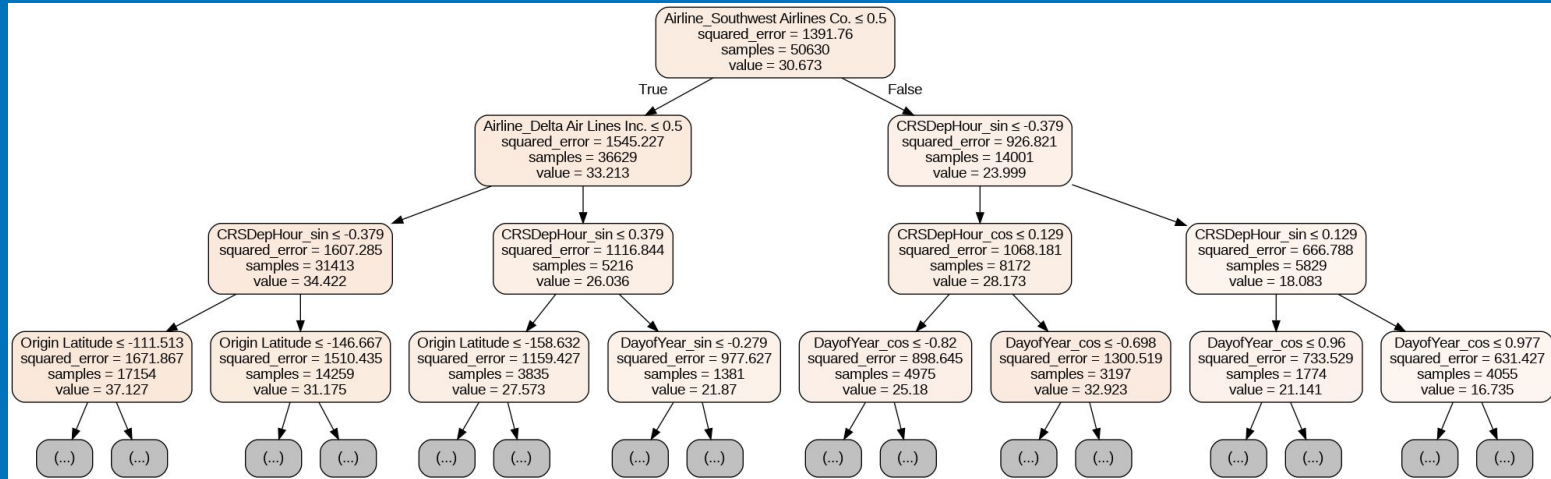
Test

MSE: 1317.68

R2: 0.04



# Random Forest + Bayes Search



Training MSE: 1263.55


Training R2: 0.08

Test MSE: 1263.55

Test R2: 0.05



# Implications and Insights

- Most important factors
    - Time of day
    - Day of the year
    - Airline
  - These three factors are already major factors individuals consider when purchasing a plane ticket.
    - Information from our model can supplement what we already know.
- 

# Limitations and Future Work

- Could not utilize entire dataset
  - Required us to sample small fraction
- Limited number of features
  - Even after encoding
- Classifier for “Cancelled” flights
- Address overfitting in Random Forest
- Add more features:
  - Airline controversy
  - Make/model of the plane: Boeing vs Airbus
  - Ticket price
  - Number of seats sold





# Reflections and Challenges

- Limited feature availability for our model
  - Even after encoding and feature engineering, our model effectively used airline, flight date, and airport location
- Should've spent more time on the EDA to identify more features
- Good practice of reviewing the entire course intimately
- Enjoyed data visualization

Thanks!

