# Mid-Term Project

## Exploring the Air Travel Industry

**Isabelle de Robert**

**& Maggie Fiander**

# Brief Agenda

- Project Overview
- Exploratory Data Analysis
  - Review of Tasks
- Modeling and Evaluations
  - Review of Process
  - Model Evaluation Metrics
  - Model Selection
- Conclusions and Next Steps

.

# The Data Structure

**passengers**

The observations of passenger totals – spanning 2015-2019

**fuel**

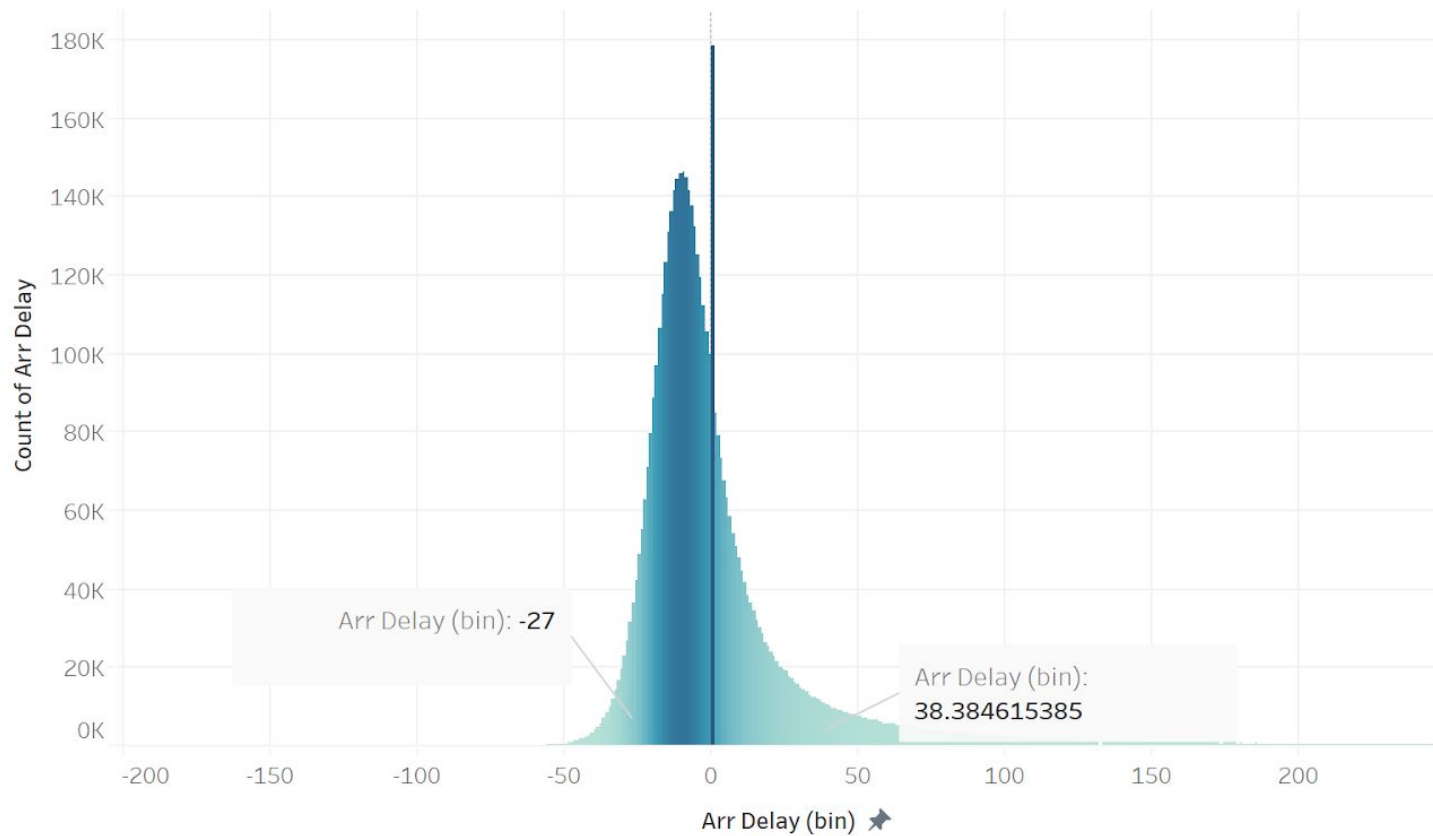The observations on fuel consumption – spanning 2015 – 2019

**flights (flights_test)**

The observations on flight arrival and departure patterns – data spanning 2018 and 2019
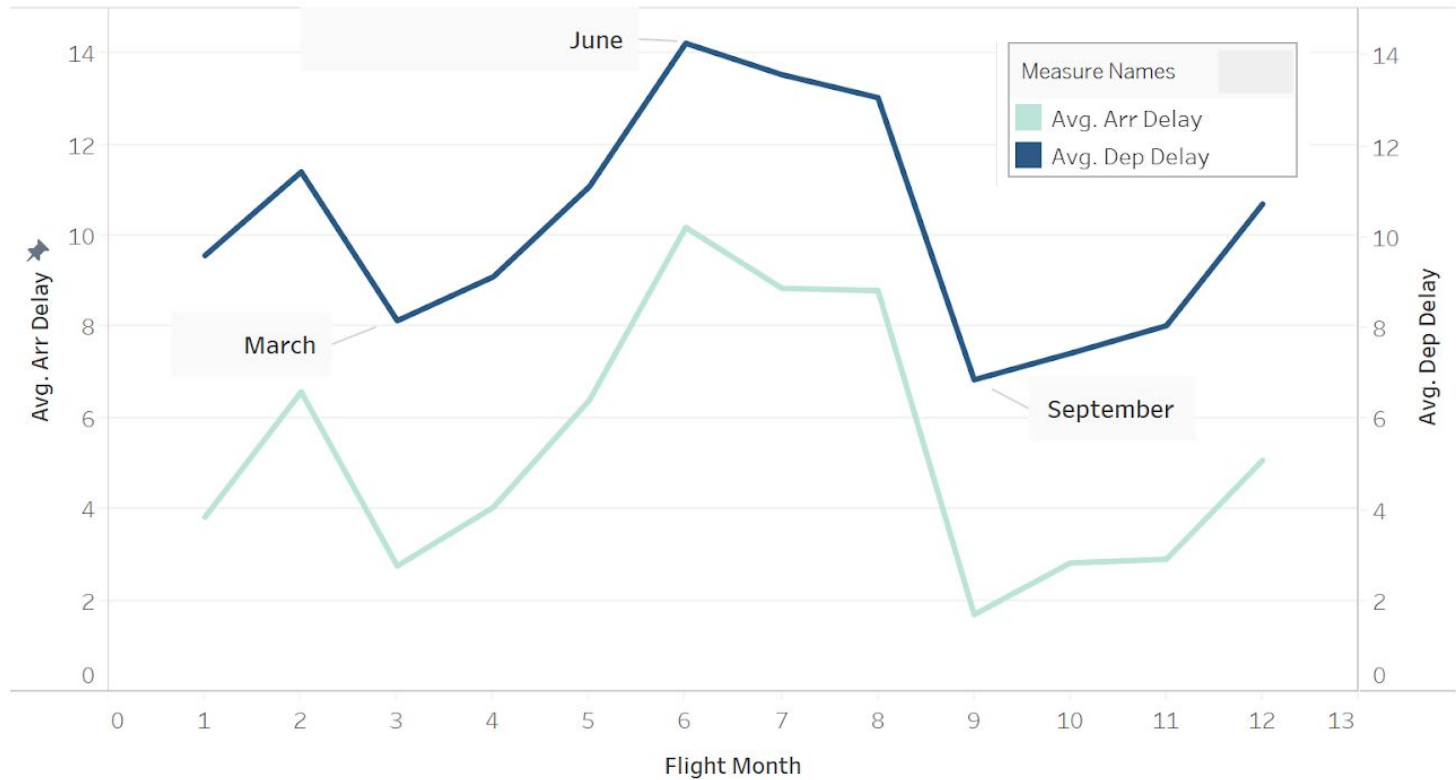
# 1

# Exploratory Data Analysis
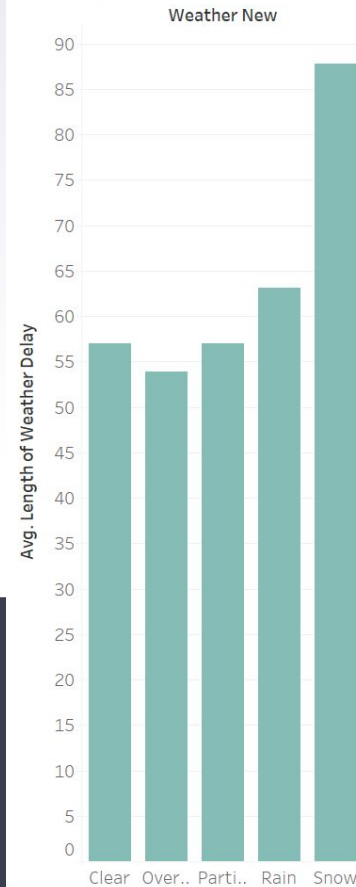
Arrival Delays (including outliers)

# Mean Monthly Arrival and Departure Delays

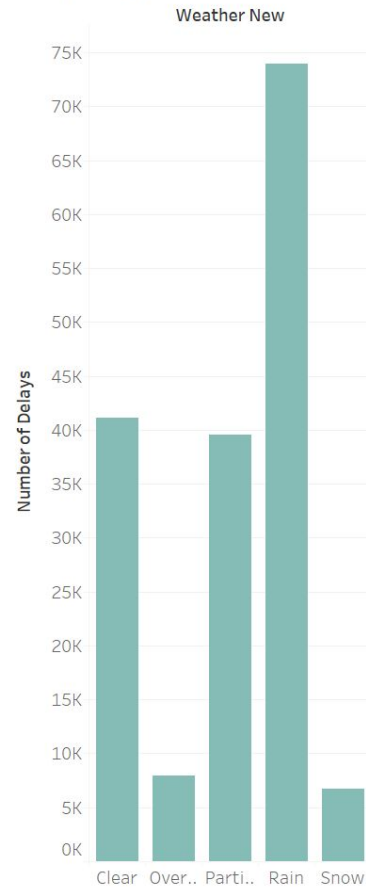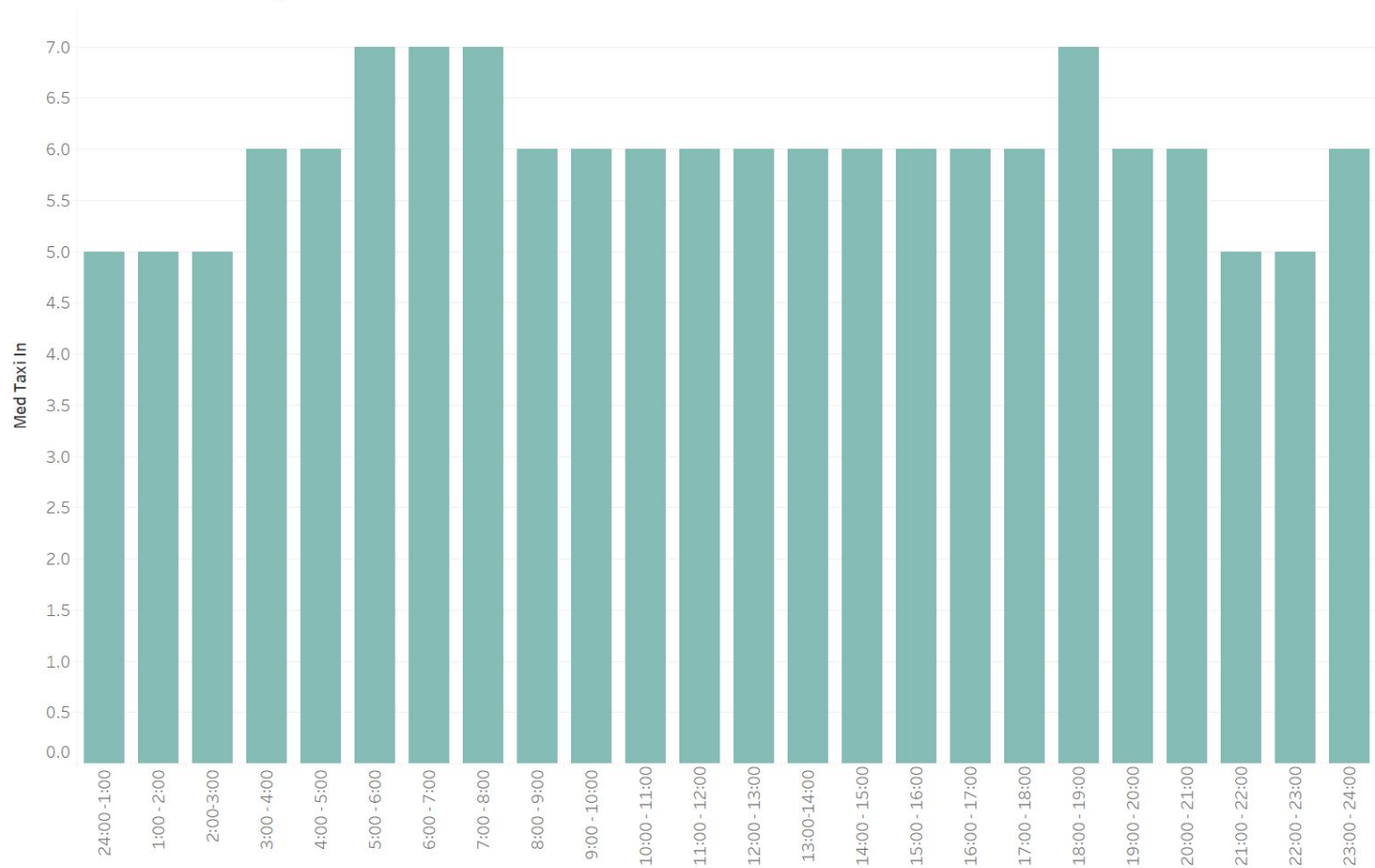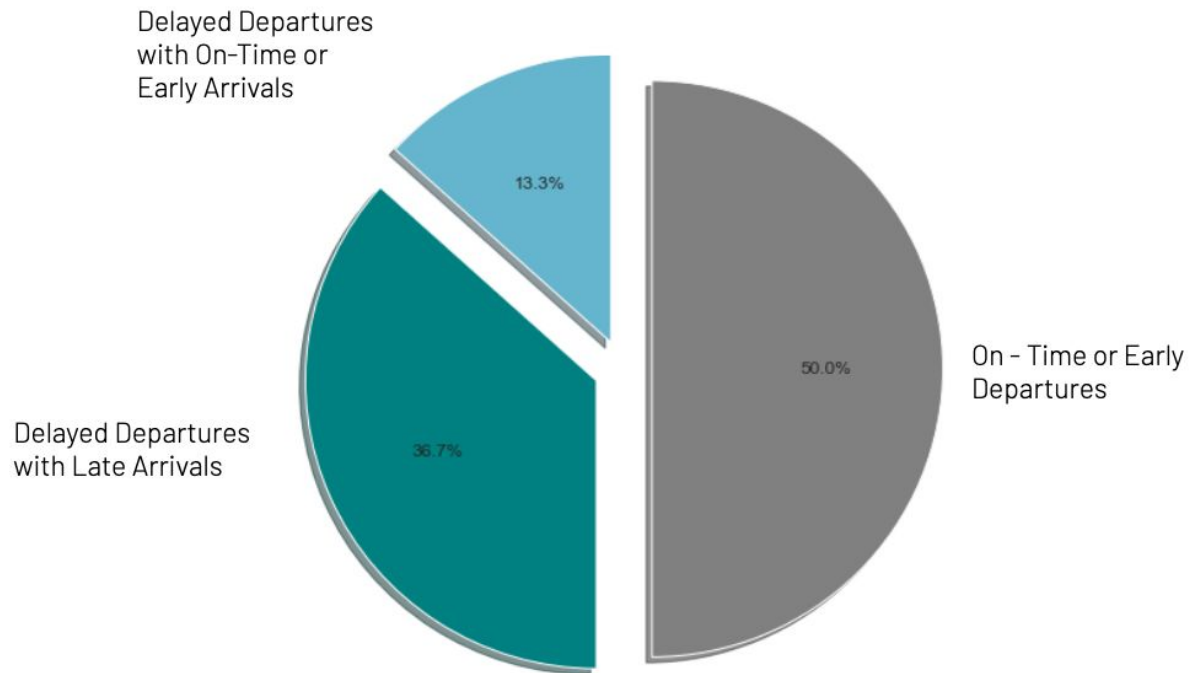# Delay Length Impacts due to Weather



Average Weather Delay by Type

Delay Frequency – Weather

Median Taxi in Time by Scheduled Arrival Time

Delayed Departures with On-Time or Early Arrivals — 13.3%

On - Time or Early Departures — 50.0%

Delayed Departures with Late Arrivals — 36.7%

# Percent of Air Traffic by State
## (using departing flights)



© 2020 Mapbox © OpenStreetMap

Comparision to Average Speed (with no departure delay)

Faster than Average 40.6%

59.4% Slower than Average

Breakdown of Haul Length

Long Haul Flights — 0.2%
Medium Haul Flights — 42.1%
Short Haul Flights — 57.7%

# Scheduled Hourly Departure Frequency by Haul Type

Airports by No of Passengers

Relationship between Arrival Delay and Fuel Consumption

Relationship between Arrival Delay And Fuel Consumption Per Mile

# 2 Modeling and Evaluation

# Process

On samples data     On entire data
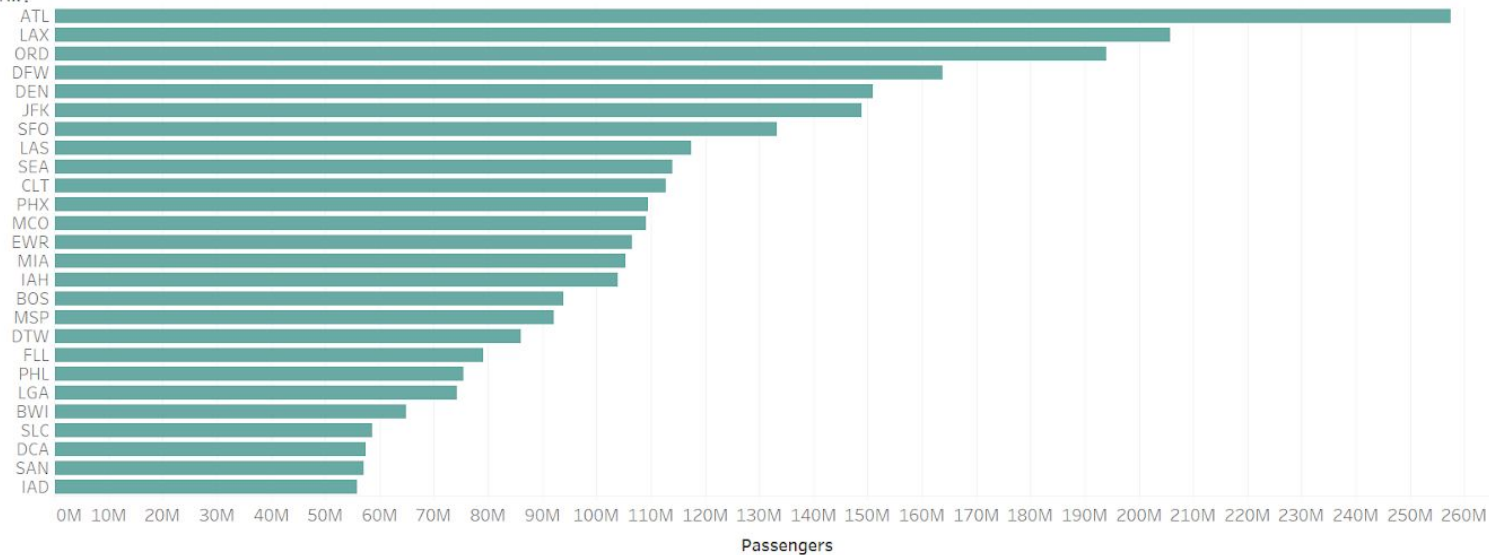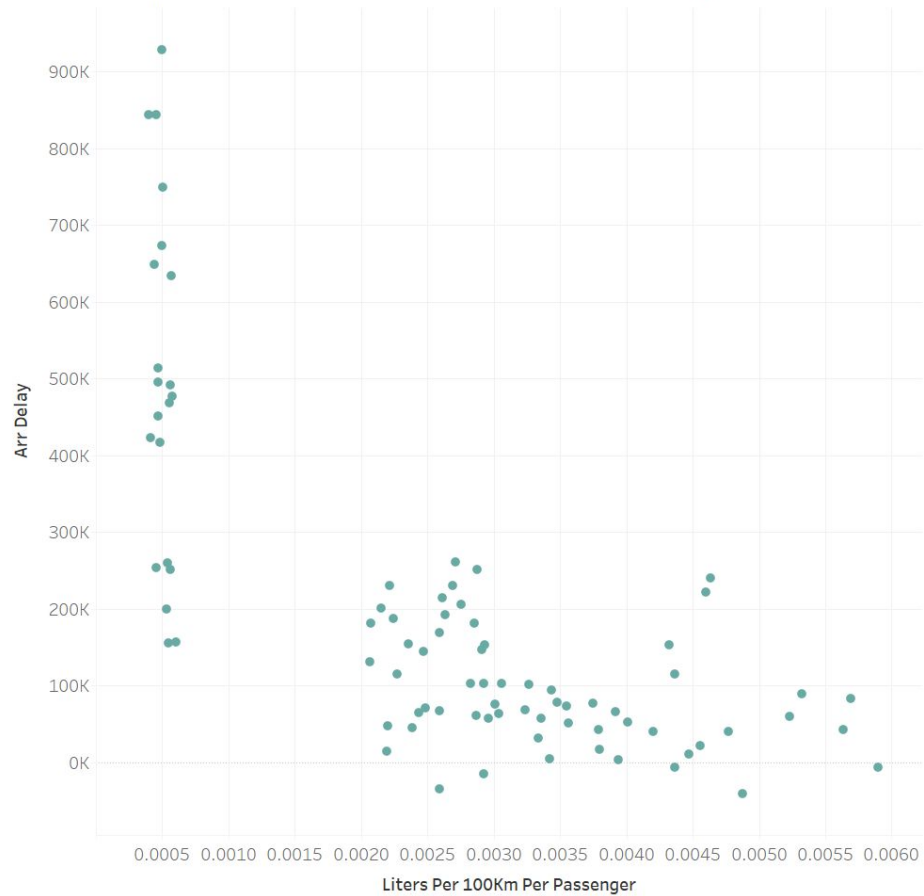
**Create samples**

**Prepare the Data**

**Feature selection and engineering**

**+**

Determine input data structure and targets

Determine Evaluation Metrics

Pick Models to compare

**→**

Spot check Models using Cross Validate

Try different sampling techniques and parameters

Tune Hyperparameters using Grid Search

Training and Testing

Model Deployment

# Feature Engineering

- Including an average measure delay from historical data - linked to flight month

- Dealing with Categorical Variables

- Time Related Data — Binning

- Dimensionality reduction (PCA)

# Evaluation Metrics

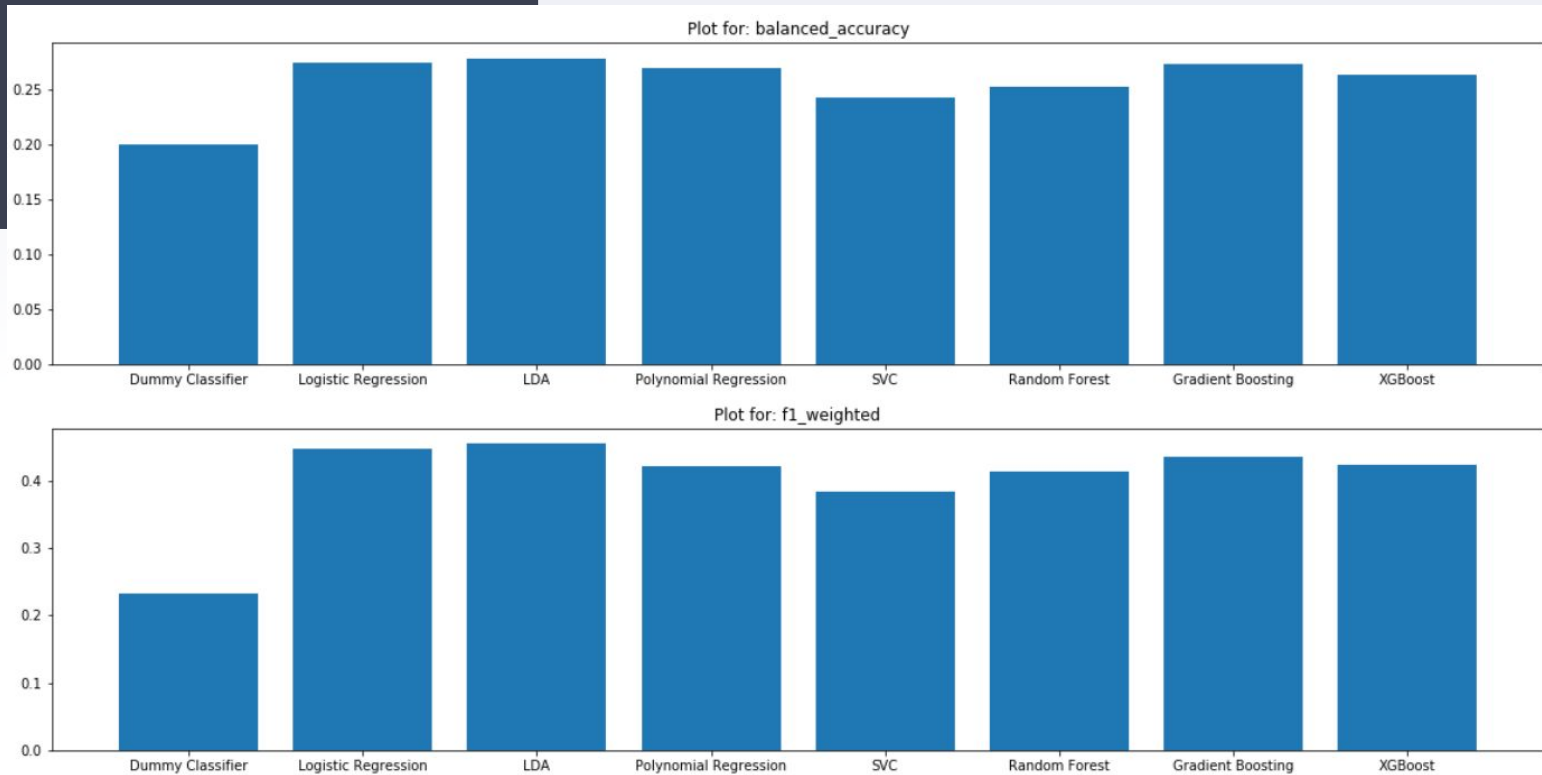▸ Regression : $R^2$ and MSE

▸ Multiclass Classification :

F1 Score (weighted), Accuracy (balanced) and AUC

▸ Binary Classification :

Brier Score, Precision-Recall AUC and F1 Score(weighted)

# Model Selection

Used scikit-learn's cross_validate function to compare algorithms
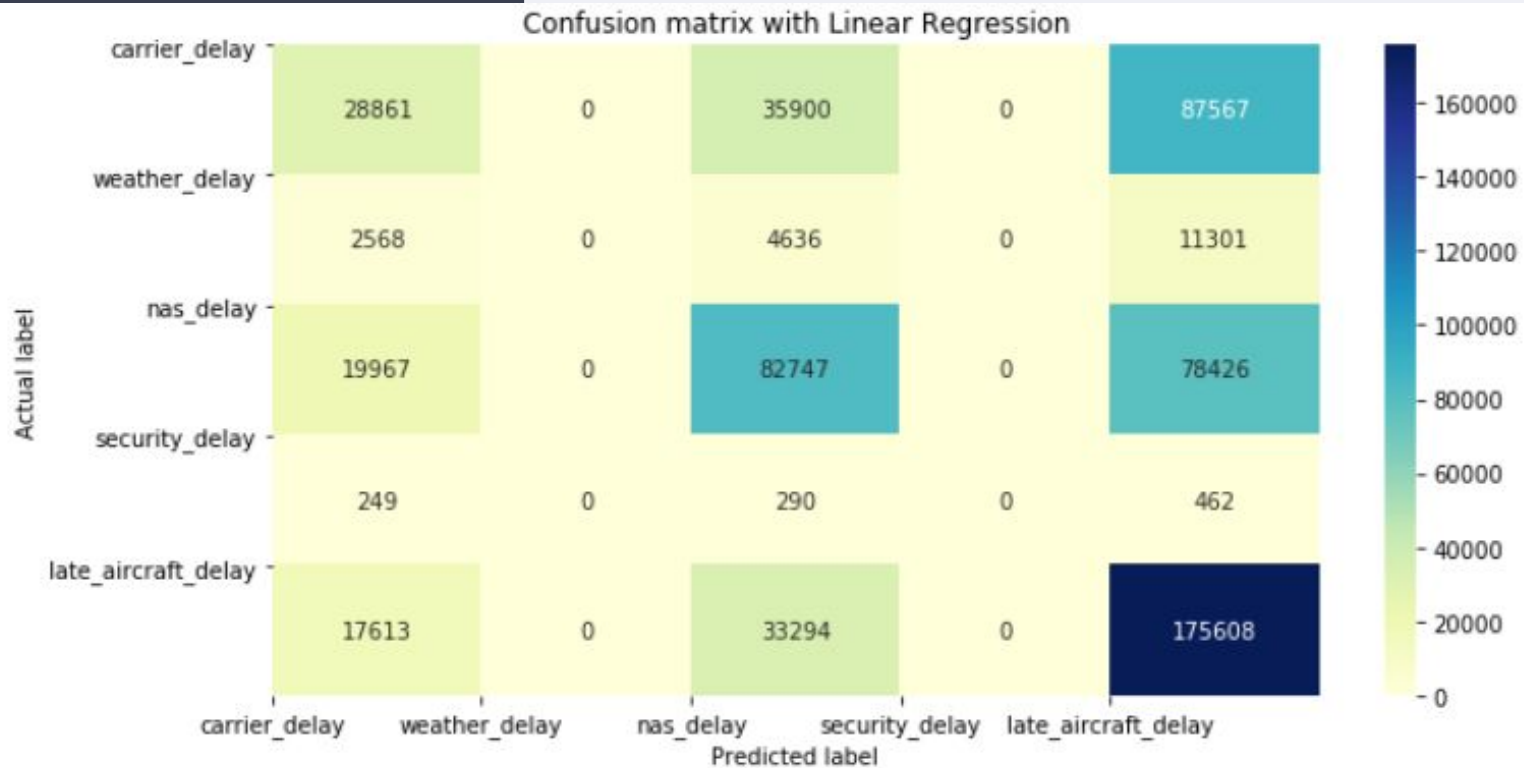
Tests performed with Linear, SVM and Ensemble Models

- ▸ Regression : XGBoost
- ▸ Multiclass Classification : LDA
- ▸ Binary Classification : Logistic Regression
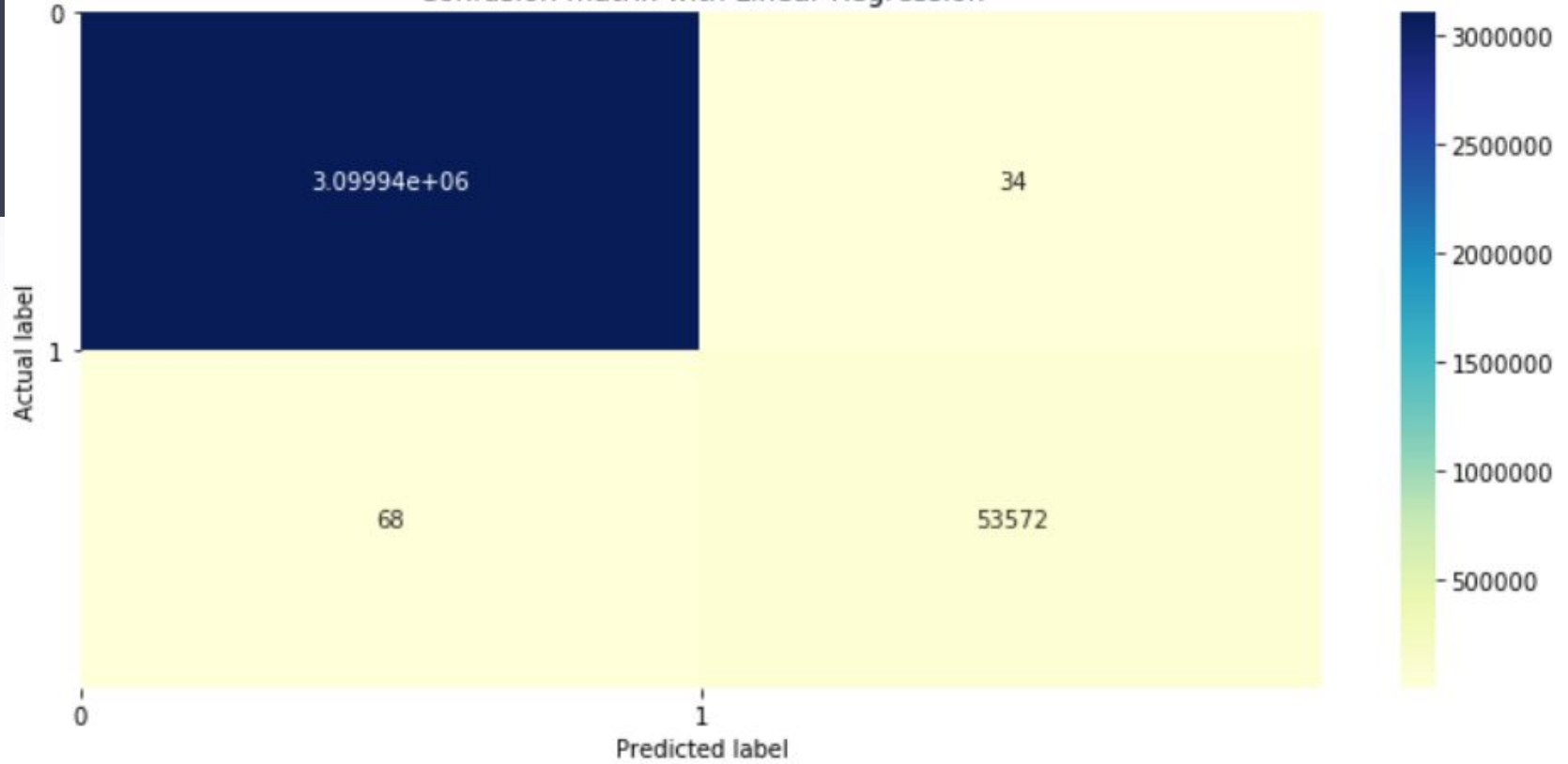
Plot for: balanced_accuracy

Plot for: f1_weighted

# Model Performance

Some visualizations of the output we obtained in our modeling process

Confusion matrix with Linear Regression

Confusion matrix with Linear Regression

# 3 Conclusions and Next Steps

# Challenges

- Remote Collaboration and Version Control
  - File Share (large .csv)
- Data Prep and Feature Engineering
  - Time Required for All Iterations
- Bug Resolution Progression
  - Propagated Errors

# Lessons

This project showed the value of:

- ▸ Prototyping and Working with Samples
- ▸ Establishing a Plan
- ▸ Clear and Reusable Code
- ▸ Understanding the Data (targets and observations)

# Next Steps

- Back to the sandbox on things we had particular challenges
- Prepare an "ideal" workflow for the changes in how we would approach the problem with what we know now
- Learning ways to integrate API data into model
- Learn more about scripts and packages

# Any questions?

You can find our repo at

- ▶ github.com/Isabelle-Dr/MidTerm-Project

# Mid-Term Project
## Exploring the Air Travel Industry

**Isabelle de Robert**

**& Maggie Fiander**