# Predicting US Flight Cancellations

By: Jason Kwan and Kyle Leung (Lecture 2 Group C)

# Table of Contents

## 01
### Introduction
Flight context and data set overview

## 02
### Methodology
Data cleaning and modeling

## 03
### Results & Discussion
Final constructed model analysis

## 04
### Limitations & Conclusions
Setbacks, assumptions and final words

# 01

# Introduction

Flight context and data
set overview

# The Flight Cancellation Problem



Flight cancellations pose quite a costly problem (both in $ and time) with both passengers and airlines

Especially, with the COVID-19 pandemic, predicting flight cancellations would be an incredibly help in mitigating this annoyance.

# Flights Data Set

## 69225

### Observations

Each observation represents a domestic flight operated by large air carriers in the US

## 45

### Variables

Detailed information recorded with each flight (ex: destination airport, day, and scheduled flight time)

# 02

# Methodology

Data cleaning and
modeling

# The Process

**Clean Data**

Deal with NA's and multicollinearity

01

**Model Data**

Create several models from our data for predictions
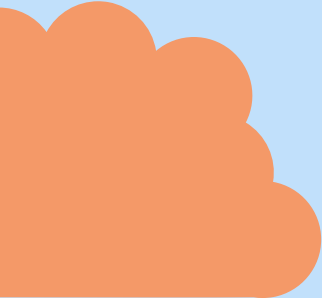
02

**Analyze Models**

Find misclassification rates from testing

03

**Compare Models**

Choose model with best MCR and simplicity

04

# Data Cleaning: Dealing with NA's

First, we needed to see if there were any predictors with NA values, and if there were, what was the proportion of NA's that each predictor has.

## Predictors with >80% NA values

- Air system delay
- Security delay
- Airline delay
- Late aircraft delay
- Weather delay
- Share White/Black/Native/Asian/Hispanic
- Median Income
- Poverty Rate
- Percent Completed HS

## Predictors with <10% NA values

- Tail number
- Pass traffic
- Aircraft Movement

Concerning the predictors with >80% NA values, it would not be beneficial to try to impute these values since we are already missing so many already. So, we deleted these predictors entirely. We also removed the predictors with <10% NA values because they ended up being highly correlated with other predictors or had massive amounts of levels

# Data Cleaning: Dealing with Multicollinearity

We found 3 main groupings of highly correlated variables. Afterwards, we chose and removed variables, accordingly.

## Time

**SCHEDULED_DEPARTURE**
SCHEDULED_TIME
**SCHEDULED_ARRIVAL**

## Origin

O.City
**O.State**
Origin_airport
Origin_city
Org_airport_lat
Org_airport_long

## Destination

**Destination_airport**
Destination_city
Destination_population
Dest_airport_lat
Dest_airport_long
Rank
Average.Passengers
Rank.Status

# Cleaned Flights Data Set

## 69225

### Observations

We were able to keep 100% of the original observations from the original data set

## 16

### Variables

We were able to reduce the 45 original variables down to only 16, including cancellations

We were now able to move on to the modeling phase and try to create the best model at predicting flight cancellations

# Data Modeling: Choosing our model

### Logistic Regression

Pros:
Simple
Fast
Interpretable

Cons:
Inaccurate
Handles Poorly
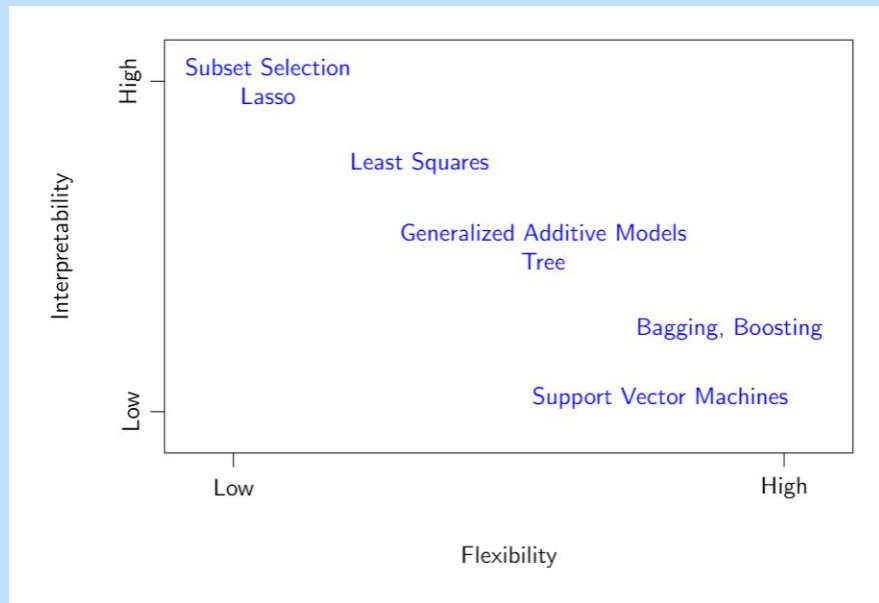Non-Flexible

### Random Forests/Boosting

Pros:
Accurate
Handles Well
Flexible

Cons:
Complicated
Slow
Non-Interpretable



After splitting the data into training and testing sets (70/30), we decided to try both methods to see which one works better.

# Logistic Regression Model

## Method

We constructed a logistic regression model on predicting flight cancellations (YES for cancelled, NO for not cancelled) with the other 15 predictor variables on the training data:

Destination_airport, O.State, Passengers, Seats, Flights, Distance, Origin_Population, MONTH, DAY, DAY_OF_WEEK, AIRLINE, FLIGHT_NUMBER, SCHEDULED_DEPARTURE, SCHEDULED_ARRIVAL, and DIVERTED
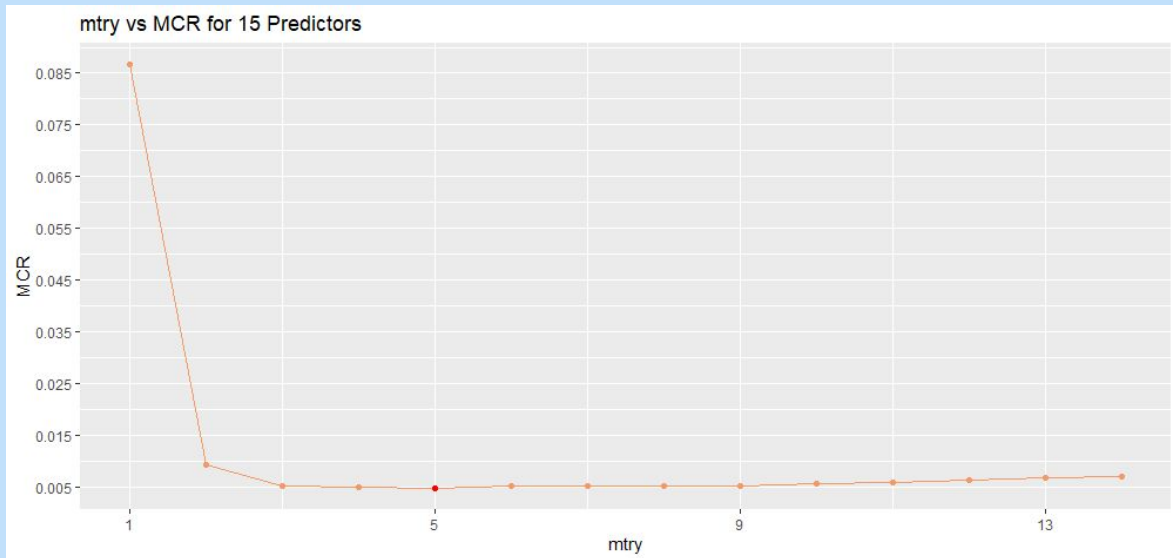
We then tried to predict cancellations on the testing data using the model.

## Test Data Confusion Matrix (X: Predicted vs Y: Actual)

|      | NO    | YES  |
|------|-------|------|
| NO   | 11859 | 3806 |
| YES  | 1950  | 3152 |

We got a misclassification rate of 27.72%, which isn't bad, but let's see if forests can improve this!

# Random Forest Model (Full)
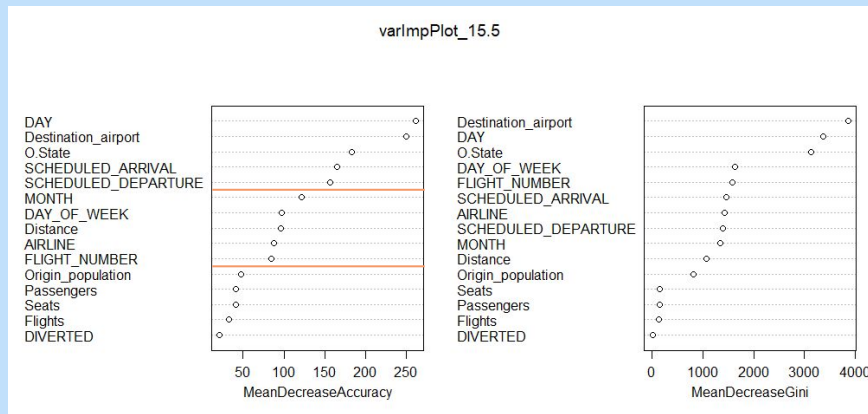


mtry vs MCR for 15 Predictors

## Method

We constructed random forest models using all 15 predictors from training. However, we varied the mtry parameter from 1-14 to see which mtry works the best.

We then took all 14 models and tried to predict the testing data and compared the misclassification rates.

We can see that the MCR converges to about 0.58% with the minimum occurring at mtry = 5.

# Random Forest Model (Full)

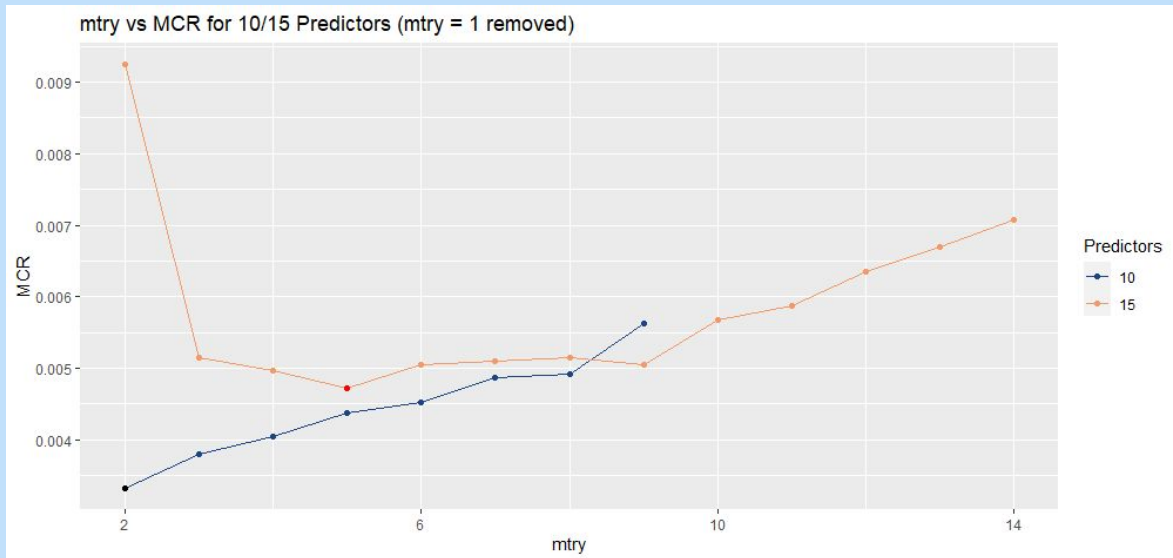## Variable Importance Plot
## 15 Predictors, mtry = 5



After looking through all of the variable importance plots, we noticed that the same predictors were grouped like above. Let's try simplifying our model by removing the predictor group with the lowest importance.

## Test Data Confusion Matrix
## (X: Predicted vs Y: Actual)

|      | NO    | YES  |
|------|-------|------|
| NO   | 13735 | 24   |
| YES  | 74    | 6934 |

We got a misclassification rate of 0.47%, which is a massive improvement from logistic, but let's see if we can get a similar rate but with fewer variables!

# Random Forest Model (10P's)



mtry vs MCR for 10/15 Predictors (mtry = 1 removed)

## Method

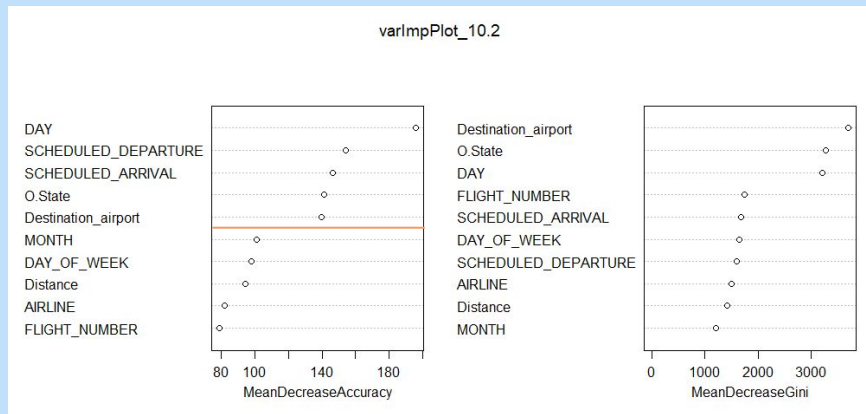We removed the variable group with the lowest importance:
Origin_population, Passengers, Seats, Flights, DIVERTED
Then, we repeated what we did before and varied mtry, now from 1-9 and compared the resulting misclassification rates.

Surprisingly, our new MCR's were actually better than the MCR's from our full model! (Minimum occurred at mtry = 2)

# Random Forest Model (10P's)

## Variable Importance Plot
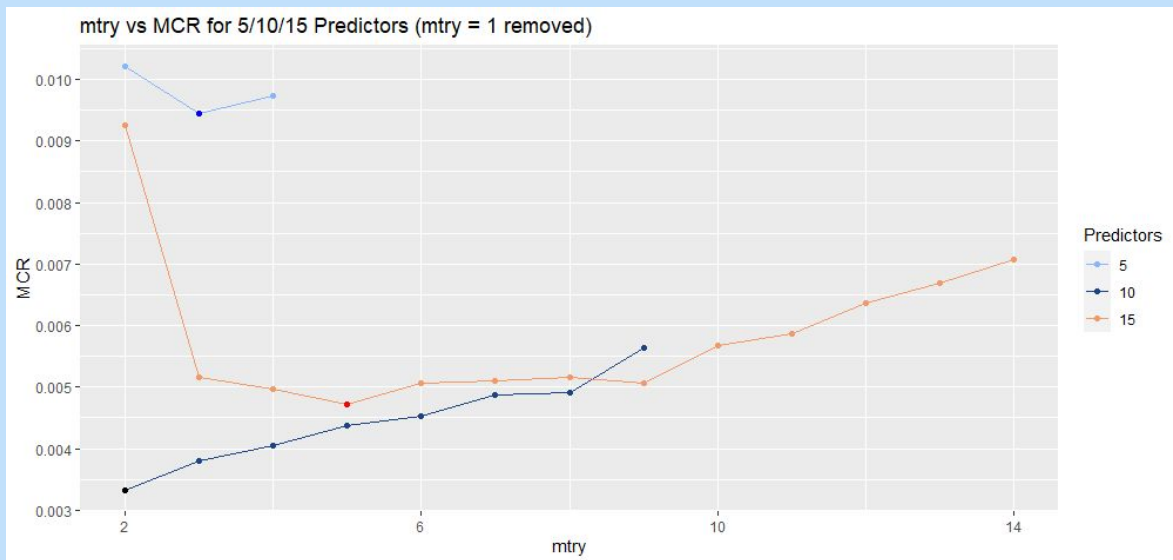## 10 Predictors, mtry = 2



After looking through all of the variable importance plots, we noticed that the same predictors were grouped like last time. Let's try to simplify our model even more again.

## Test Data Confusion Matrix
## (X: Predicted vs Y: Actual)

|       | NO    | YES  |
|-------|-------|------|
| NO    | 13765 | 25   |
| YES   | 44    | 6933 |

We got a misclassification rate of 0.33%, which is an improvement from last time but with even fewer variables! Let's see if we can do it again with even fewer variables.

# Random Forest Model (5P's)



mtry vs MCR for 5/10/15 Predictors (mtry = 1 removed)

## Method

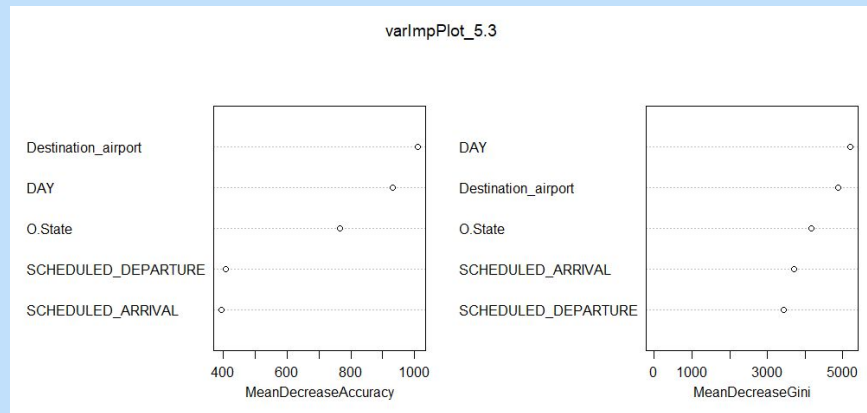We removed the variable group with the next lowest importance:
MONTH, DAY_OF_WEEK, Distance, AIRLINE, FLIGHT_NUMBER
Then, we repeated what we did before and varied mtry, now from 1-4 and compared the resulting misclassification rates.

Sadly, these new MCR's did worse and were actually far larger than even our full model MCR's. (Minimum occurred at mtry = 3)

# Random Forest Model (5P's)

## Variable Importance Plot
## 5 Predictors, mtry = 3



Even though, we can see an obvious separation here, we did not proceed with removing predictors further due to the large jump in MCR's.

## Test Data Confusion Matrix
## (X: Predicted vs Y: Actual)

|      | NO    | YES  |
|------|-------|------|
| NO   | 13643 | 30   |
| YES  | 166   | 6928 |

We got a misclassification rate of 0.94%, which was not an improvement from last time. However, we still have one more technique to try!

# Boosting Model

## Method

We constructed a boosting model with all 15 predictors. In terms of accuracy, random forests and boosting should theoretically be some of the best techniques at classification. So, we wanted to see if boosting does just as, or even better, than random forests.

After constructing the model from the training data, we then proceeded to, as always, try to predict cancellations on the testing data.

## Test Data Confusion Matrix (X: Predicted vs Y: Actual)

|       | NO    | YES  |
|-------|-------|------|
| NO    | 13451 | 358  |
| YES   | 48    | 6910 |

We got a misclassification rate of 1.96%, which is actually larger than all the MCR's from random forest. We then concluded our modeling phase and now we had to choose which model to use.

# Final Model

## Choosing the Model

When choosing the best model, we ideally wanted the model with the lowest MCR and is the most simple. However, realistically, this is not always possible and so we have to balance and weigh our options.

From seeing our model's MCRs and simplicity, we decided to choose our random forest model (mtry = 2) with 10 predictors.

It had the lowest MCR as well as it being "moderate" in terms of simplicity among our options.

### Misclassification Rates

| In % | Logis. | RF 15 | RF 10 | RF 5 | Boost |
|------|--------|-------|-------|------|-------|
| MCR  | 27.72  | 0.47  | **0.33** | 0.94 | 1.96 |

### Simplicity

| | Logis. | RF 15 | RF 10 | RF 5 | Boost |
|------|--------|-------|-------|------|-------|
| Rank | 1 | 4 | **3** | 2 | 5 |

03

# Results & Discussion

Final constructed model analysis

# Analysis: Key Parts

**Model**

Random Forest (mtry: 2)

**Observations**

69225 Flights

**Predictors**

10 Flight Predictors

**Simplicity**

Moderate

**MCR**

0.25%

**Rank**

Top 15

# Discussion: The Important Predictors

**Location**
Origin State, Destination
Airport, and Distance
01

**Date**
Month, Day, and
Day of the Week
02

03
**Flight Characteristics**
Airliner Name and
Flight Number

04
**Time**
Scheduled Arrival
and Departure

# Discussion: The MOST Important Predictors



### Origin State

### Day

### Destination Airport

From our model, we were able to find exactly what predictors were the most important in predicting flight cancellations. For the most part, if you know where the flight came from, where its going to, and what day it will happen, you should still be confident in being able to accurately predict whether the flight will be cancelled or not using our model, and be able to plan accordingly!

# 04

# Limitations & Conclusions

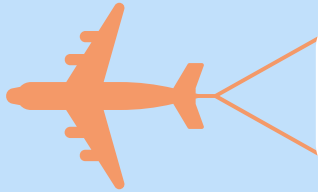Setbacks, assumptions and final words

# Limitations

**01**

Data Cleaning
May have chosen
weaker variables

**02**

Modeling
Was not able to use
regular trees

**03**

Assumptions
Fortunately, there are
no RF assumptions

# Conclusion

We get very good prediction using our model, however, we get poor understandings of our causes. The typical variables weren't included in our dataset (ex. Weather)

Further research would be necessary to delve deeper into this problem

But we believe passengers and airlines would still benefit from our model!
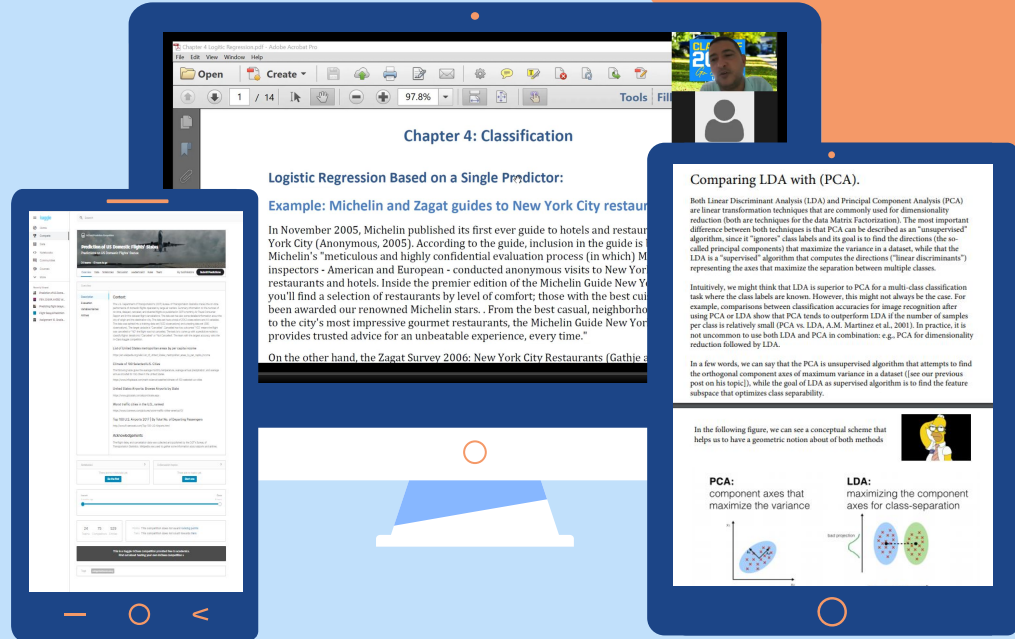
# References

Special Thanks To:

Almohalwas and Wang's Stats 101C Lectures and Discussions

Towards Data Science Inc. https://towardsdatascience.com/

Biodiversity and Climate Change Virtual Laboratory https://bccvl.org.au/

Thank you for traveling with us!
And send us your questions!

# THANKS

Merry Christmas and happy travels!

And always remember:
1. Don't drink and drive!
2. Don't shoot a mosquito with a missile!
3. And always know who the daddy is!