

An Investigation into

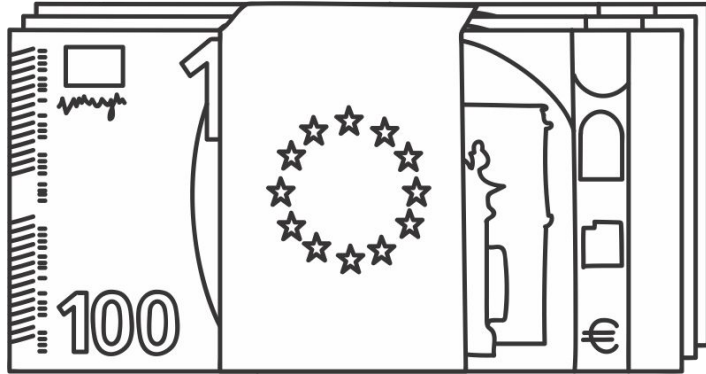
Flight Delays



Rachael Alexandroff, Sofia Pignataro, Racquel Fygenson, Ruxin Shen
Group 13

Flight delays are a hassle.

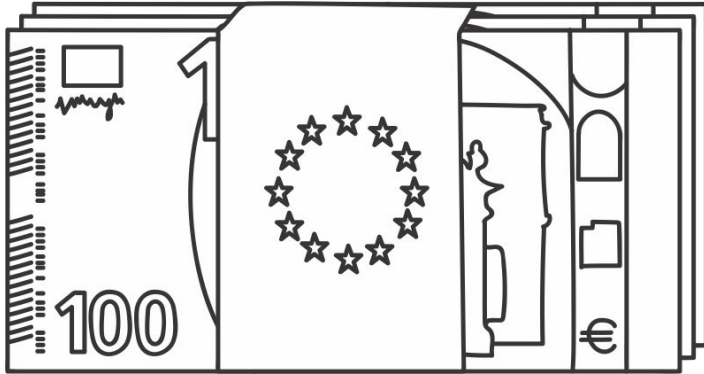
Flight delays are a hassle



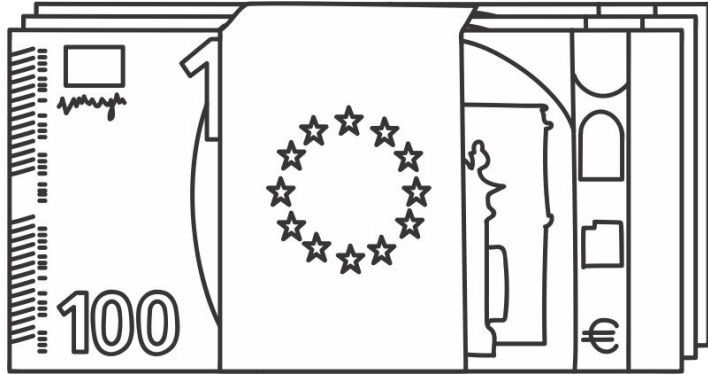
EU regulation 261/2004 requires
airlines to give you money
if your flight is
delayed > 3 hours!

Flight delays are a hassle

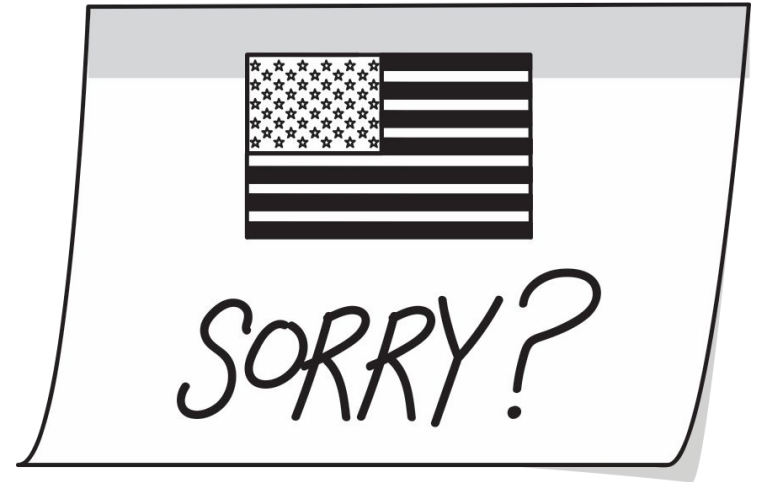
What about in the US?



Flight delays are a hassle



VS



Flight delays are a hassle

Are delayed flights
a problem in the US?

Are delayed flights
a problem in the US?

Heck, yes.

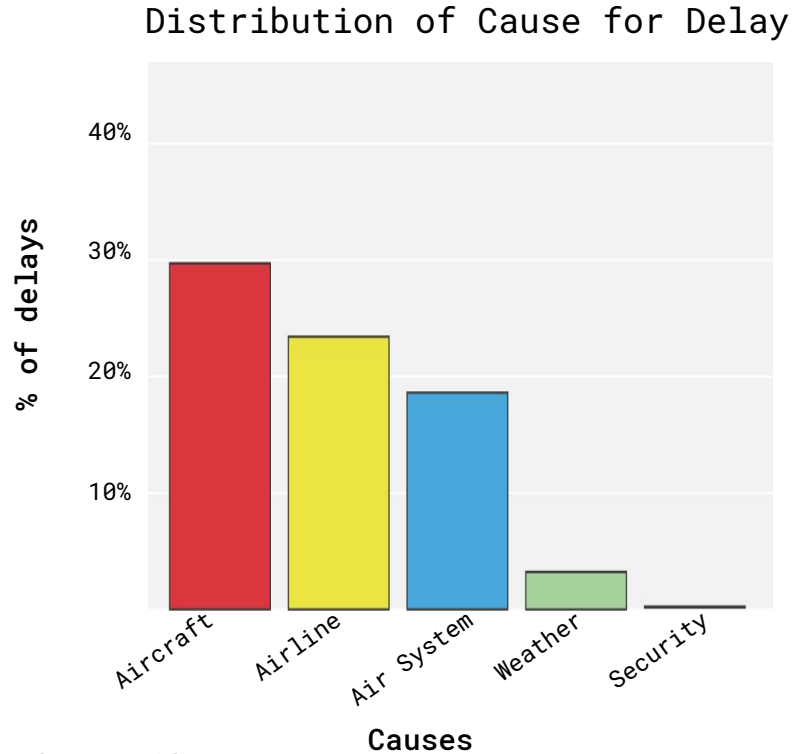
1 in 5 US flights
were delayed in 2017



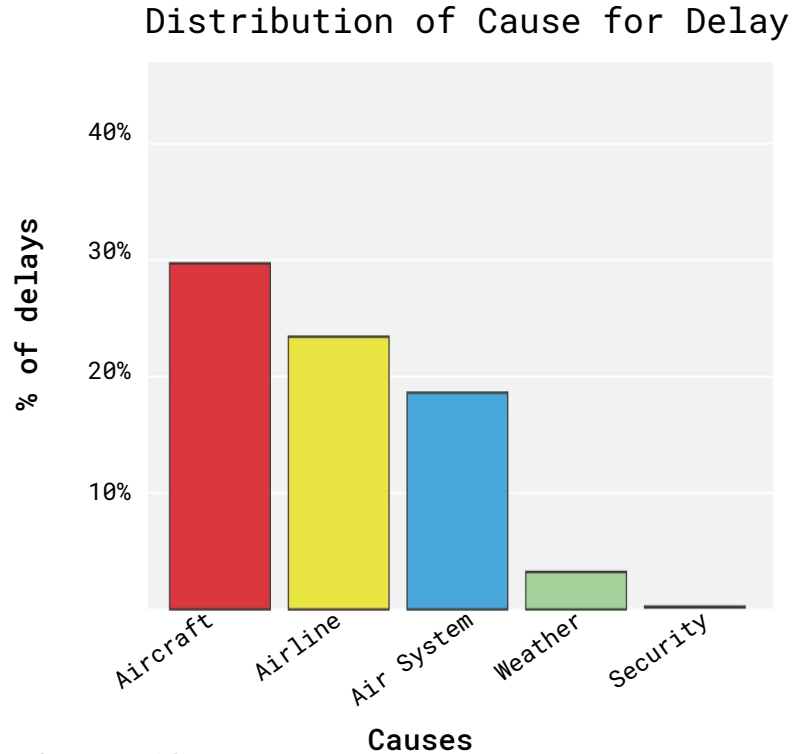


What factors correlate with delay?

What factors affect delay?

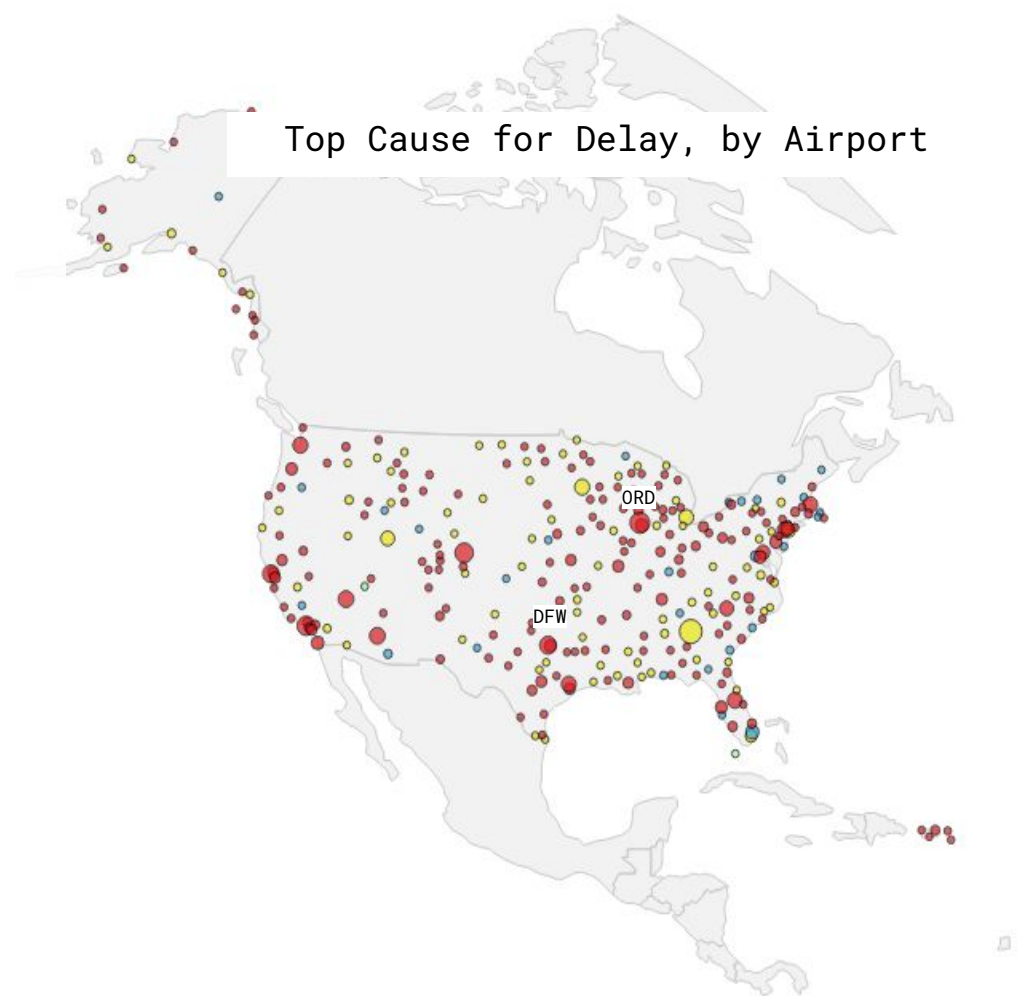


What factors affect delay?



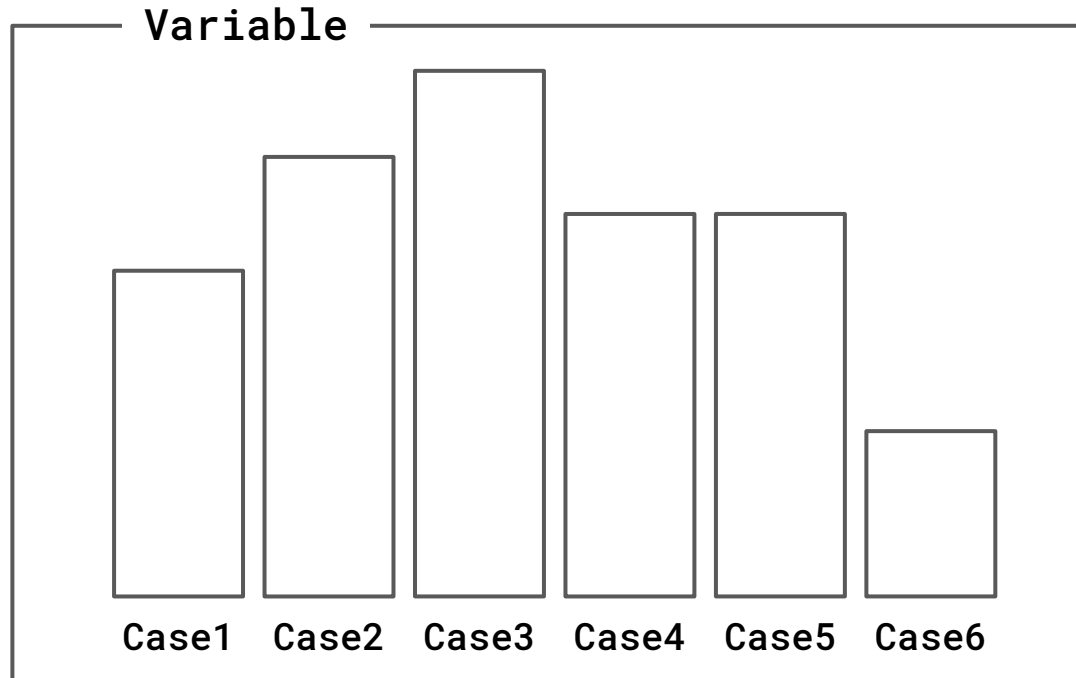
Group 13

Top Cause for Delay, by Airport



What factors affect delay?

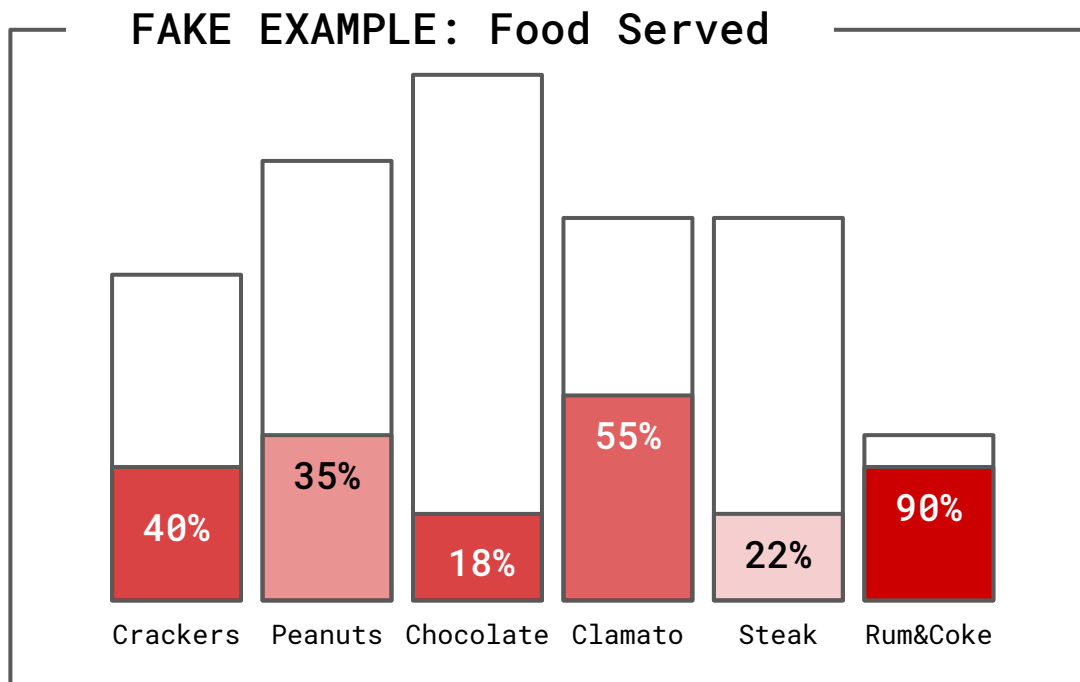
Proportional Summary*



Let's separate by variables we think might be relevant

What factors affect delay?

Proportional Summary*

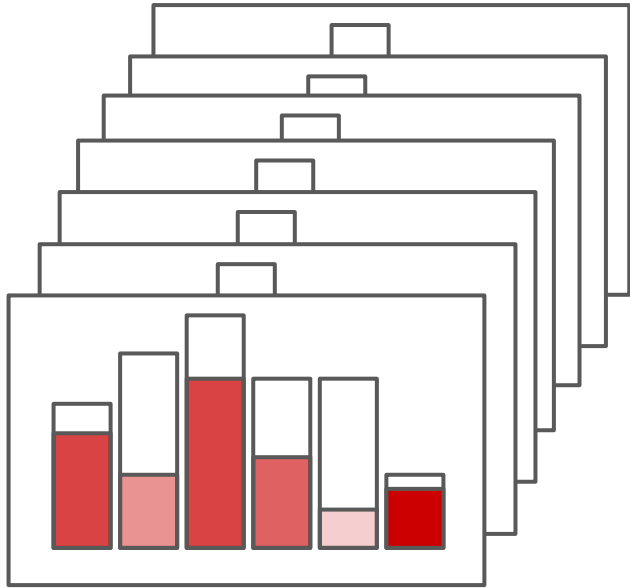


Let's separate by variables we think might be relevant

and look at the proportion of flights that are delayed

What factors affect delay?

Proportional Summary*

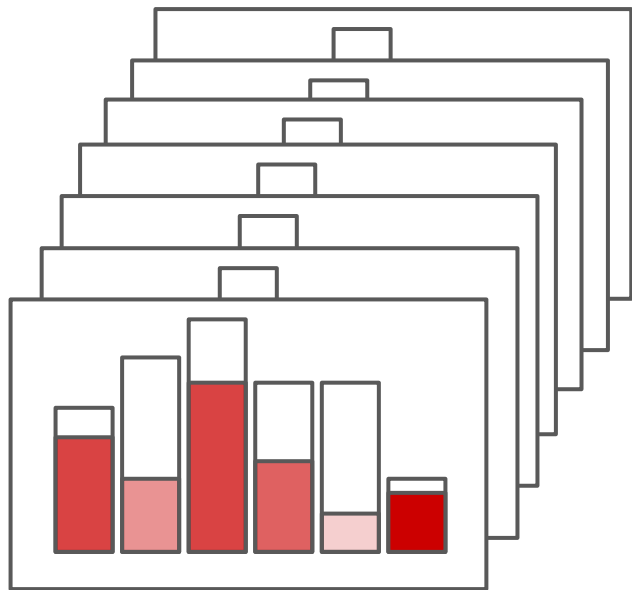


Variables we evaluated

- Day of the week
- Month of the year
- Time of day (4 buckets)
- Elapsed flight time
- Distance of flight
- Airlines
- Season

What factors affect delay?

Proportional Summary*



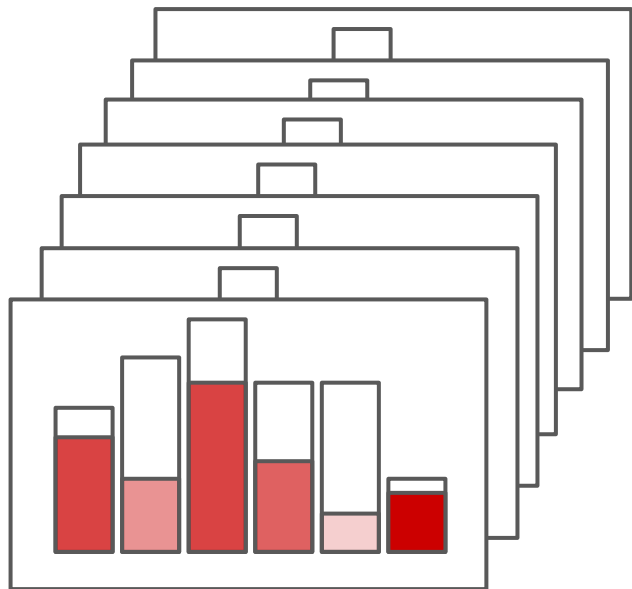
Group 13

Variables we evaluated

- Day of the week
- Month of the year
- Time of day (4 buckets)
- Elapsed flight time
- Distance of flight
- Airlines (2 buckets)
- Season

What factors affect delay?

Proportional Summary*



Data Used
Flight Traffic

Data Not Used
Weather
Fare
Event

Outside scope of question →

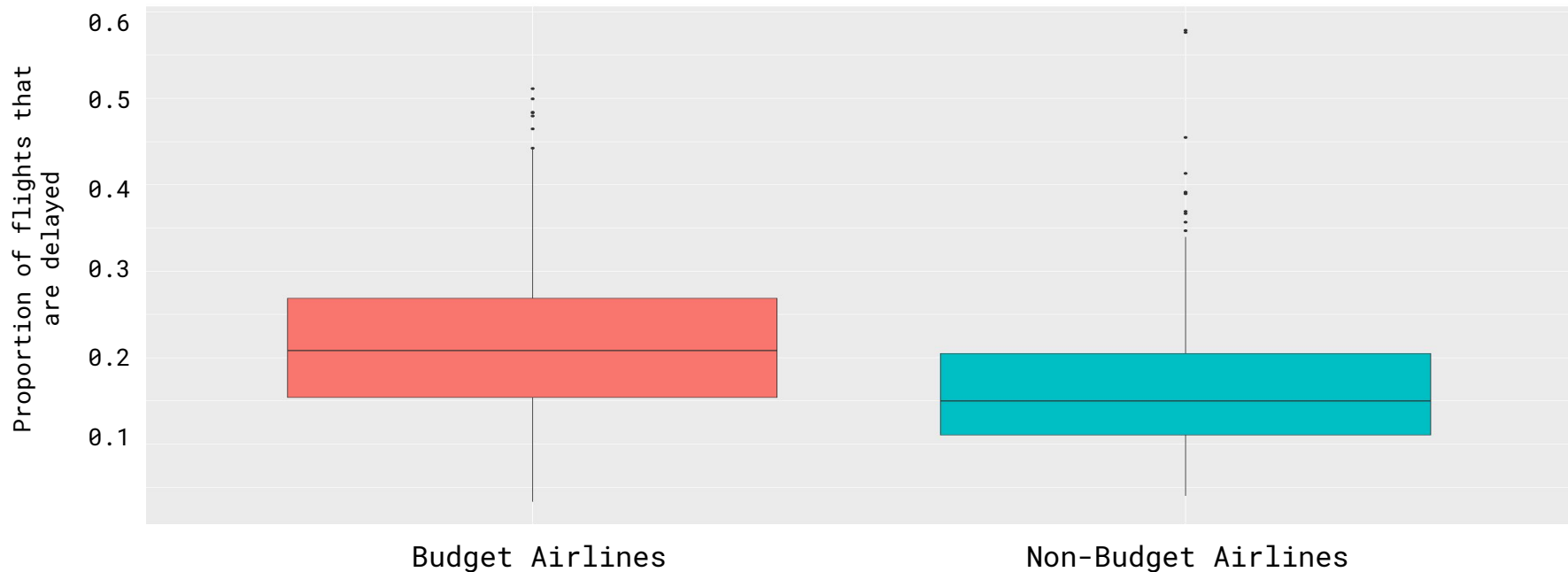
Do not contain important data points (like ORD or DFW) that are represented in our main dataset

Variables we evaluated

- Day of the week
- Month of the year
- Time of day (4 buckets)
- Elapsed flight time
- Distance of flight
- Airlines (2 buckets)
- Season

What factors affect delay?

Hypothesis Testing



What factors affect delay?

Hypothesis Testing



	Budget
Not Budget	S

S	Significant
NS	Not Significant

H_0 = Proportion of delayed flights are equal
 H_a = Proportion of delayed flights are not equal

"Budget" Airlines:

- Spirit
- JetBlue
- ExpressJet
- Frontier
- SkyWest
- Southwest
- Virgin

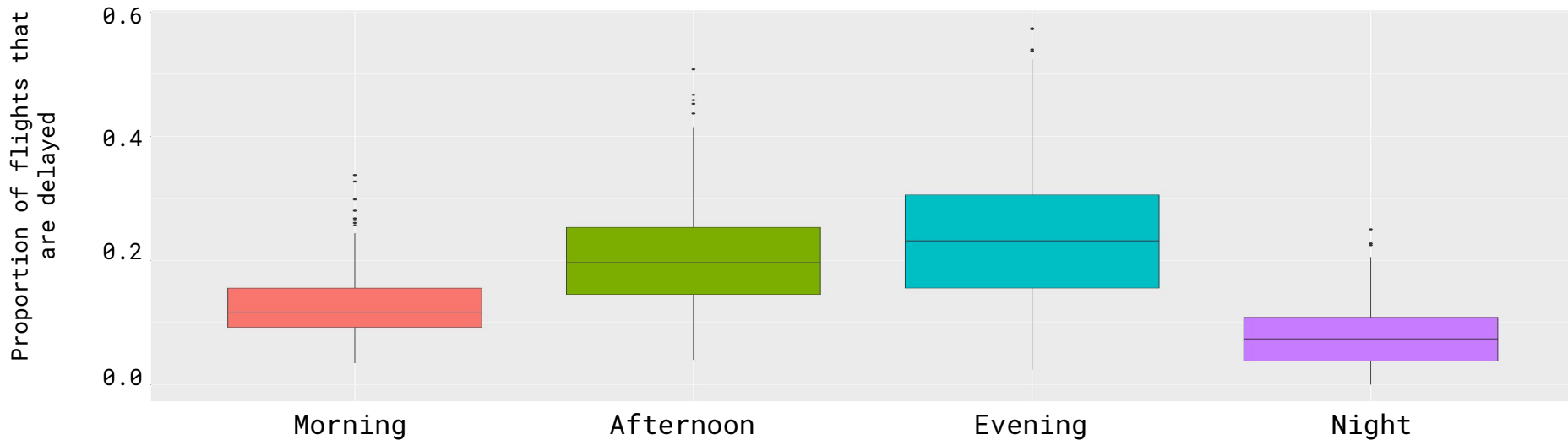
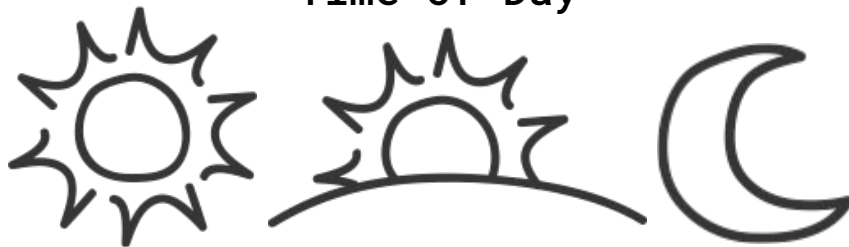
"Non-Budget" Airlines:

- American
- Delta
- Hawaiian Air
- United
- Alaska Air

What factors affect delay?

Hypothesis Testing

Time of Day



What factors affect delay?

Hypothesis Testing

Time of Day



	Morning	Afternoon	Evening
Afternoon	S		
Evening	S	S	
Night	S	S	S

H_0 = Proportion of delayed flights are equal

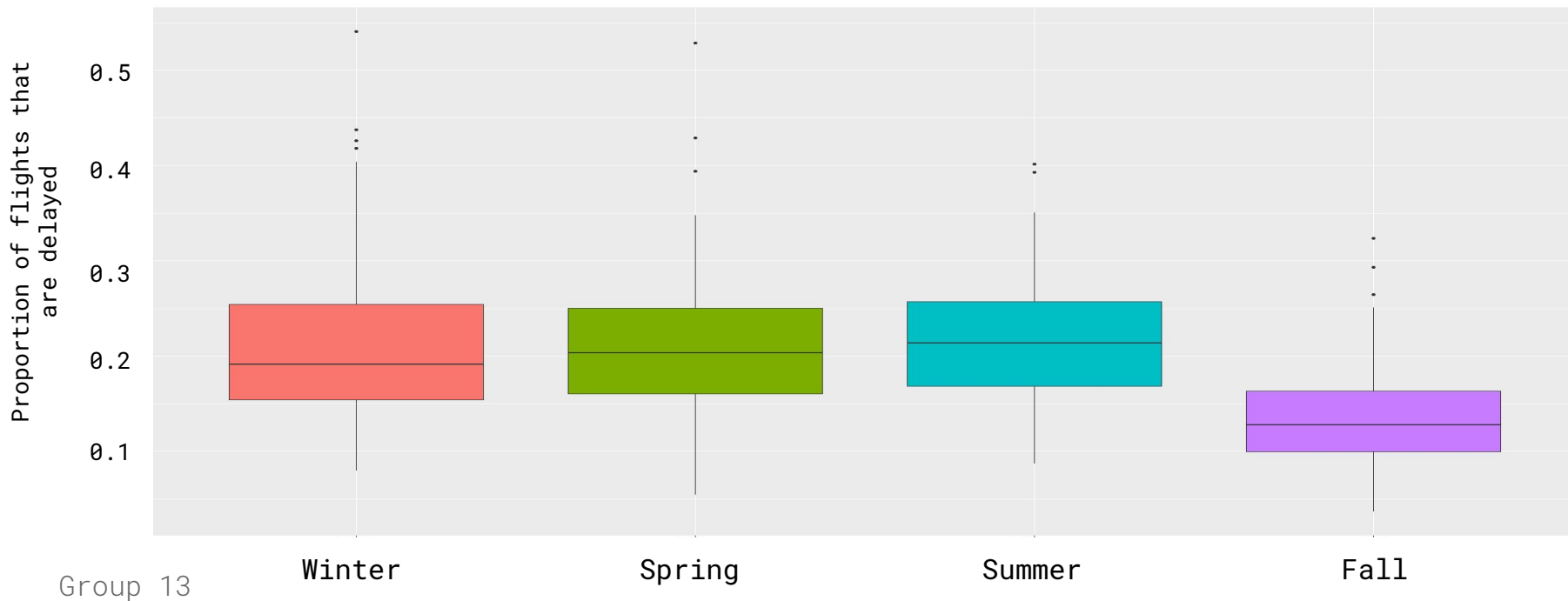
H_a = Proportion of delayed flights are not equal

S	Significant
NS	Not Significant

What factors affect delay?

Hypothesis Testing

Season



What factors affect delay?

Hypothesis Testing



Season

	Spring	Summer	Autumn
Summer	NS		
Autumn	S		
Winter	NS	NS	S

H_0 = Proportion of delayed flights are equal

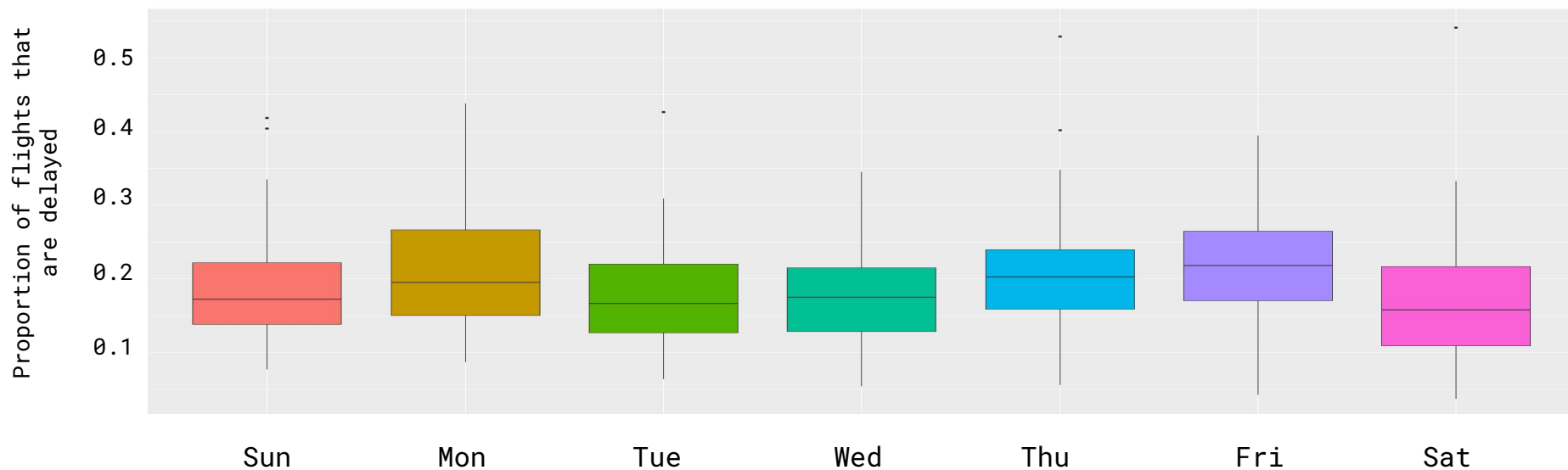
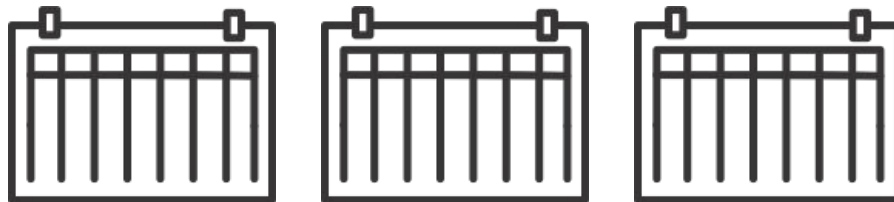
H_a = Proportion of delayed flights are not equal

S	Significant
NS	Not Significant

What factors affect delay?

Hypothesis Testing

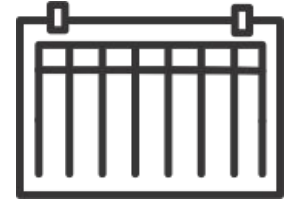
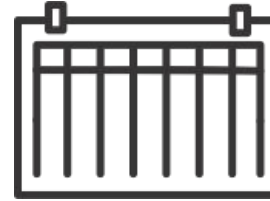
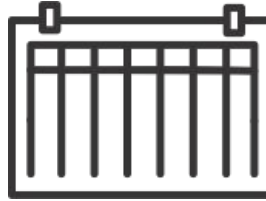
Day of the Week



What factors affect delay?

Hypothesis Testing

Day of the Week



	Sun	Mon	Tue	Wed	Thu	Fri
Mon	S					
Tue	NS	S				
Wed	NS	S	NS			
Thu	S	NS	S	S		
Fri	S	NS	S	S	S	
Sat	S	S	S	S	S	S

H_0 = Proportion of delayed flights are equal

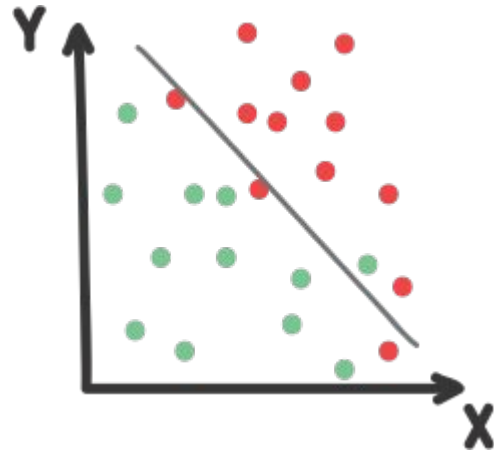
H_a = Proportion of delayed flights are not equal

S	Significant
NS	Not Significant

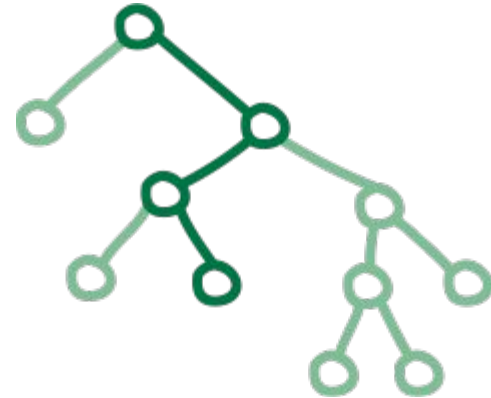
Can we predict delay?

Machine Learning

Benchmark:
Logistic Regression



Modeling:
Random Forest

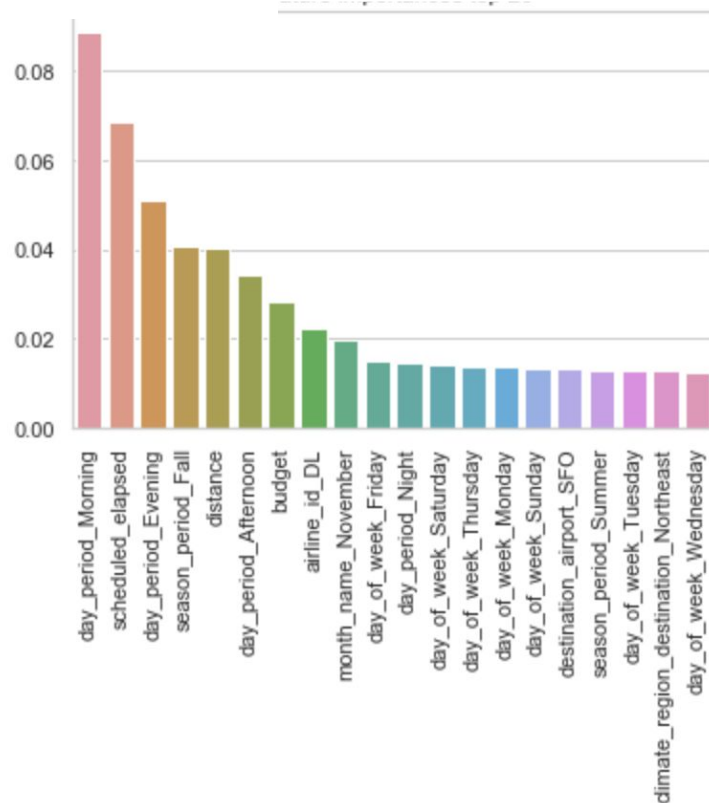


Can we predict instance of delay?

Random Forest

Delay: Y/N?

Features



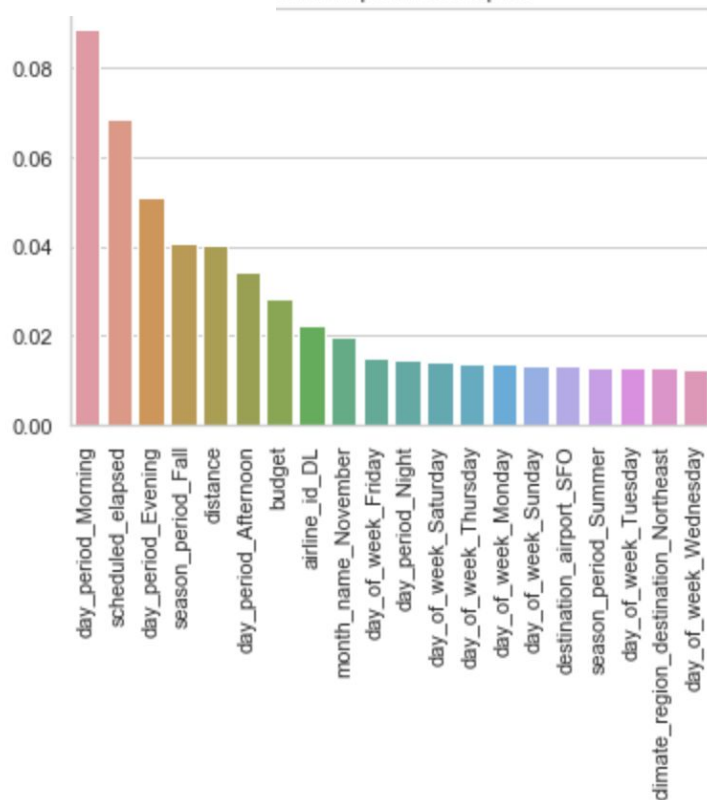
Can we predict instance of delay?

Random Forest

Delay: Y/N?

Time of day is an important feature

Features



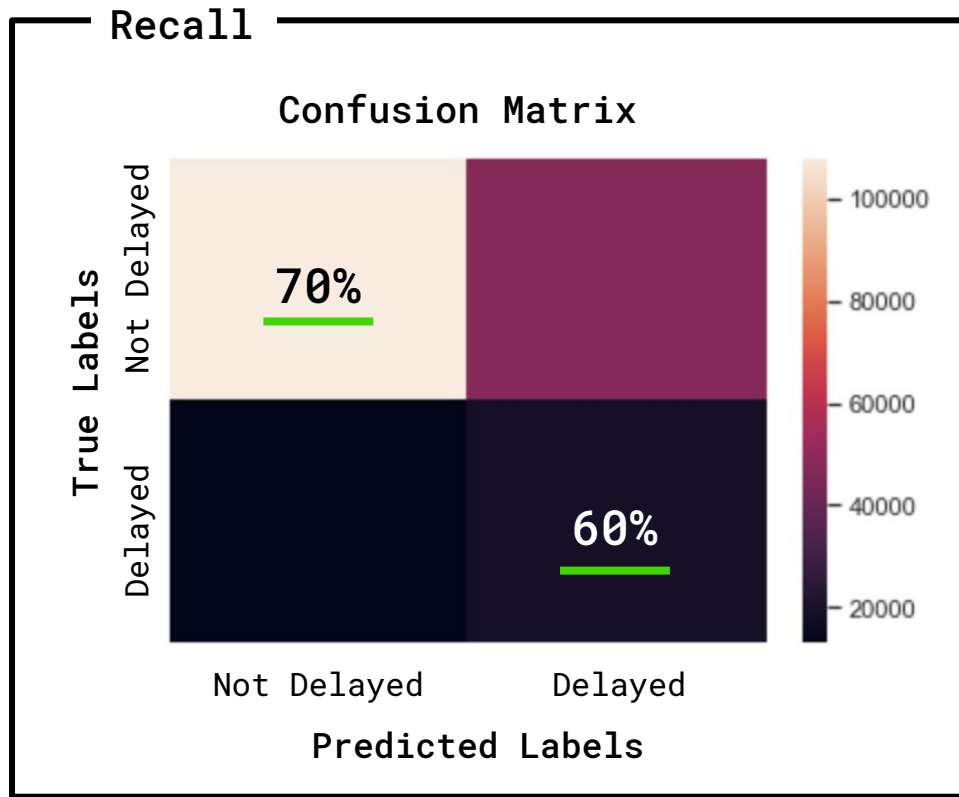
Can we predict instance of delay?

Random Forest

Delay: Y/N?

Right now, our
model has a lot of
false positives

Overall Accuracy: 68%



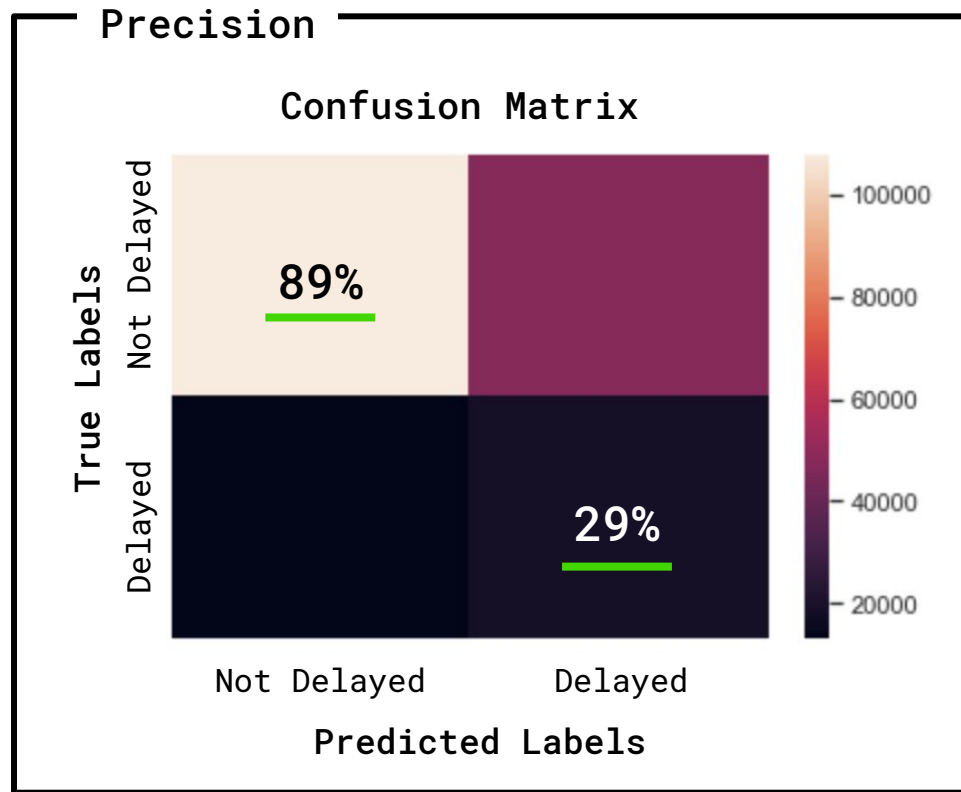
Can we predict instance of delay?

Random Forest

Delay: Y/N?

Right now, our
model has a lot of
false positives

Overall Accuracy: 68%

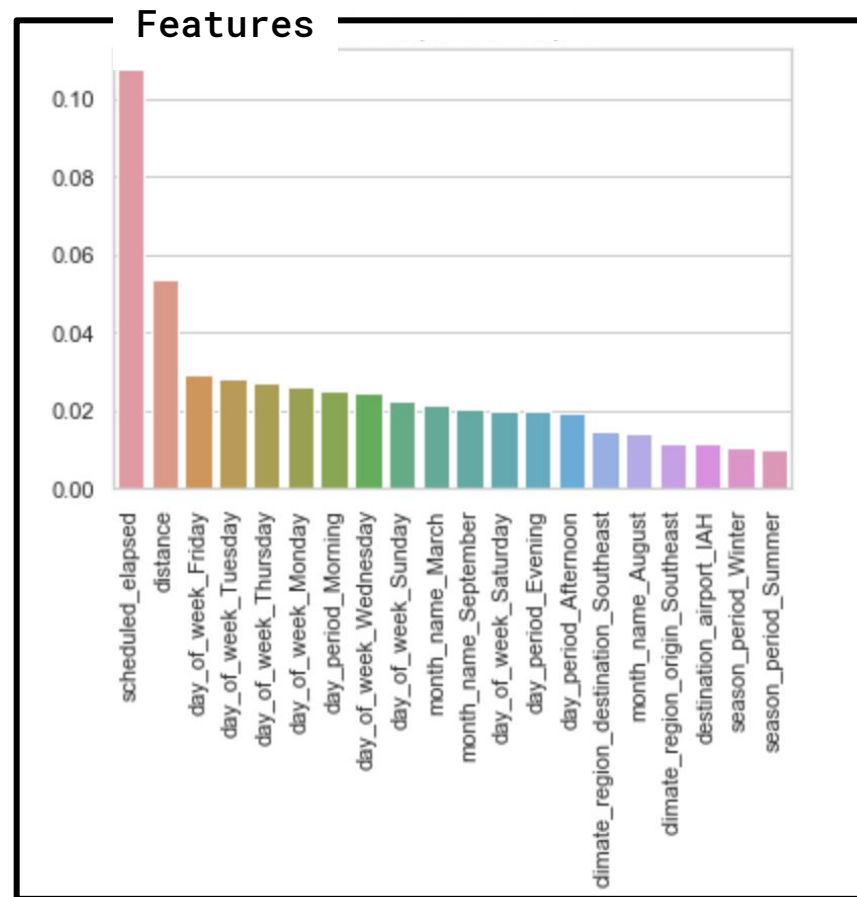


Can we predict length of delay?

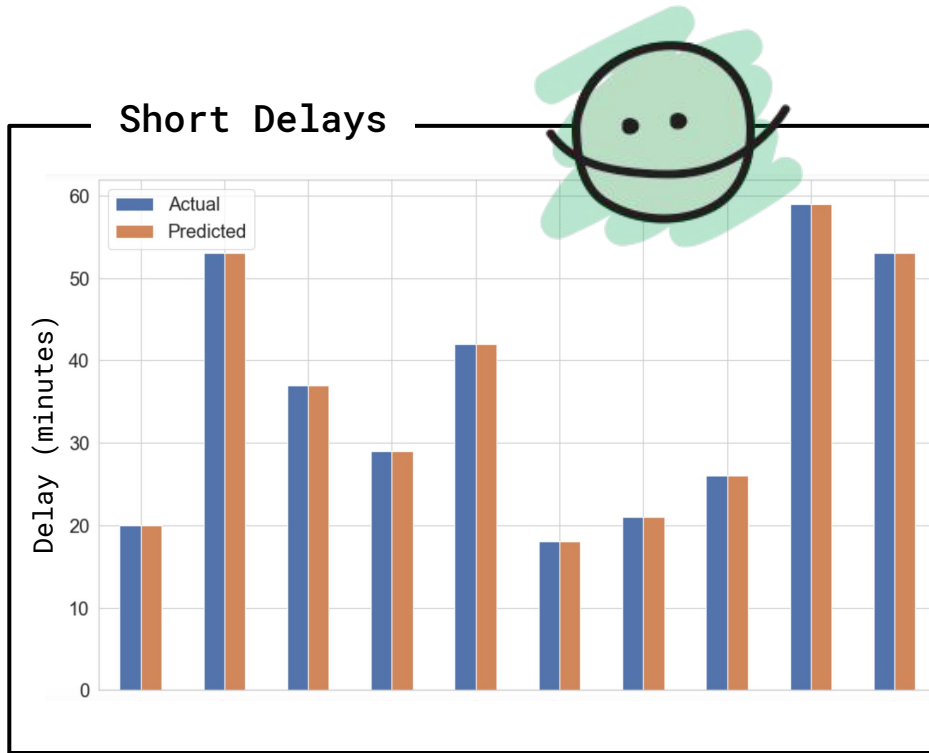
Random Forest

Length of Delay?

Flight **duration**,
flight **distance**,
day of week,
and **month** are
important features

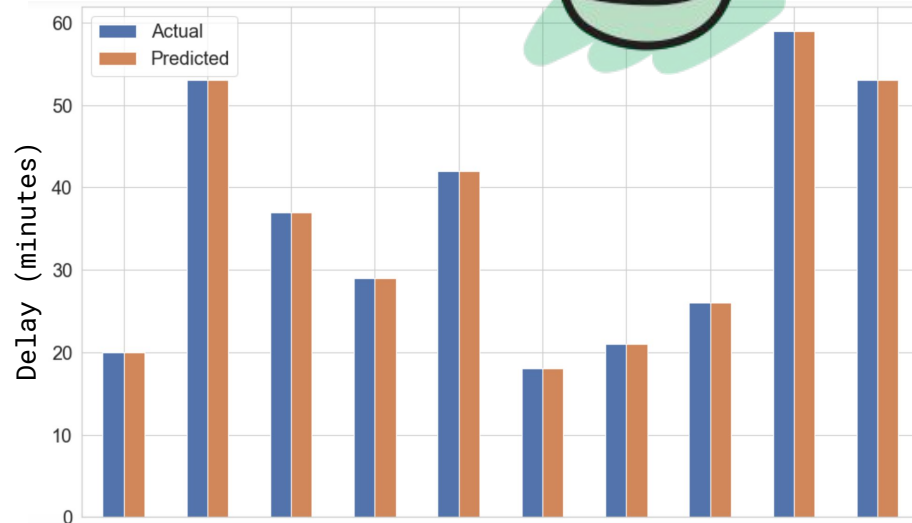
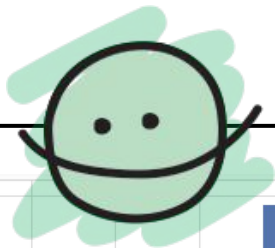


Can we predict length of delay?

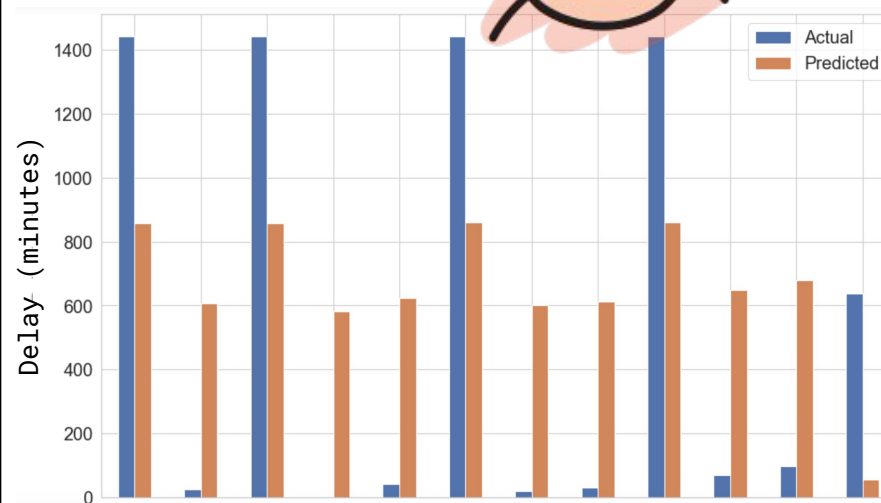
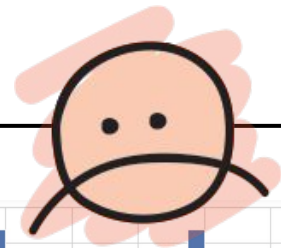


Can we predict length of delay?

Short Delays



Long Delays



Practical Application

(So What?)

Alert shoppers when a flight
is at risk of being delayed.

Practical Application

(So What?)

Alert shoppers when a flight
is at risk of being delayed.



Practical Application

(So What?)

Alert shoppers when a flight is at risk of being delayed.



One way ▾ 1 passenger ▾ Economy ▾



Trips



Explore



Flights



Hotels



Packages

○ Boston BOS



📍 Newark EWR

📅 Fri, Jul 10



Bags ▾

Stops ▾

Airlines ▾

Price ▾

Times X

Connecting airports ▾

More ▾

📈 Track prices ⓘ ☐

📅 Date grid

📊 Price graph

📍 Nearby airports

Best flights ⓘ

Total price includes taxes + fees for 1 adult. [Additional bag fees](#) and other fees may apply.

Sort by: ⬆️



12:00 PM – 1:38 PM

United

1h 38m

BOS-EWR

Nonstop

🛫 \$81



8:45 PM – 10:16 PM

United

1h 31m

BOS-EWR

Nonstop

🛫 \$81



1:26 PM – 2:59 PM

Del... · Operated by Republic Airways Delta Connecti... BOS-EWR

1h 33m

Nonstop

\$111



Practical Application

(So What?)

Alert shoppers when a flight is at risk of being delayed.



One way ▾ 1 passenger ▾ Economy ▾



Trips



Explore



Flights



Hotels



Packages

○ Boston BOS



📍 Newark EWR

📅 Fri, Jul 10



Bags ▾

Stops ▾

Airlines ▾

Price ▾

Times X

Connecting airports ▾

More ▾

📈 Track prices ⓘ ☐

📅 Date grid

📊 Price graph

📍 Nearby airports

Best flights ⓘ

Total price includes taxes + fees for 1 adult. [Additional bag fees](#) and other fees may apply.

Sort by: ⬆️



12:00 PM – 1:38 PM

United

1h 38m

BOS-EWR

Nonstop



\$81



8:45 PM – 10:16 PM

United

1h 31m

BOS-EWR

Nonstop



\$81



1:26 PM – 2:59 PM

Del... · Operated by Republic Airways Delta Connecti... BOS-EWR

1h 33m

Nonstop

\$111



Practical Application

(So What?)

Alert shoppers when a flight is at risk of being delayed.

Google

One way ▾ 1 passenger ▾ Economy ▾

○ Boston BOS ↔ Newark EWR

Fri, Jul 10

Bags ▾ Stops ▾ Airlines ▾ Price ▾ Times X Connecting airports ▾ More ▾

Track prices ⓘ ☐

Best flights ⓘ

Total price includes taxes + fees for 1 adult. [Additional baggage fees](#)

	12:00 PM – 1:38 PM	United	BOS-EWR		\$81	▾	
	8:45 PM – 10:16 PM	United	1h 31m BOS-EWR	Nonstop		\$81	▾
	1:26 PM – 2:59 PM	Del... · Operated by Republic Airways Delta Connecti...	1h 33m BOS-EWR	Nonstop		\$111	▾

Nearby airports

Sort by: ↑↓

According to our predictions, this flight has a **70% chance of being delayed**. Consider flying on **Saturday morning** instead to **reduce this risk to 30%**.

Things we could do better

Get our hands on more data



Optimize feature selection

Bucket delays by type:

- >3hrs ("Catastrophic")
- <3hrs ("Not Catastrophic")
- Cancelled



Bucket delays by cause:

- Airline
- Weather
- Air System
- Aircraft
- Security



Thank you

A quick overview...

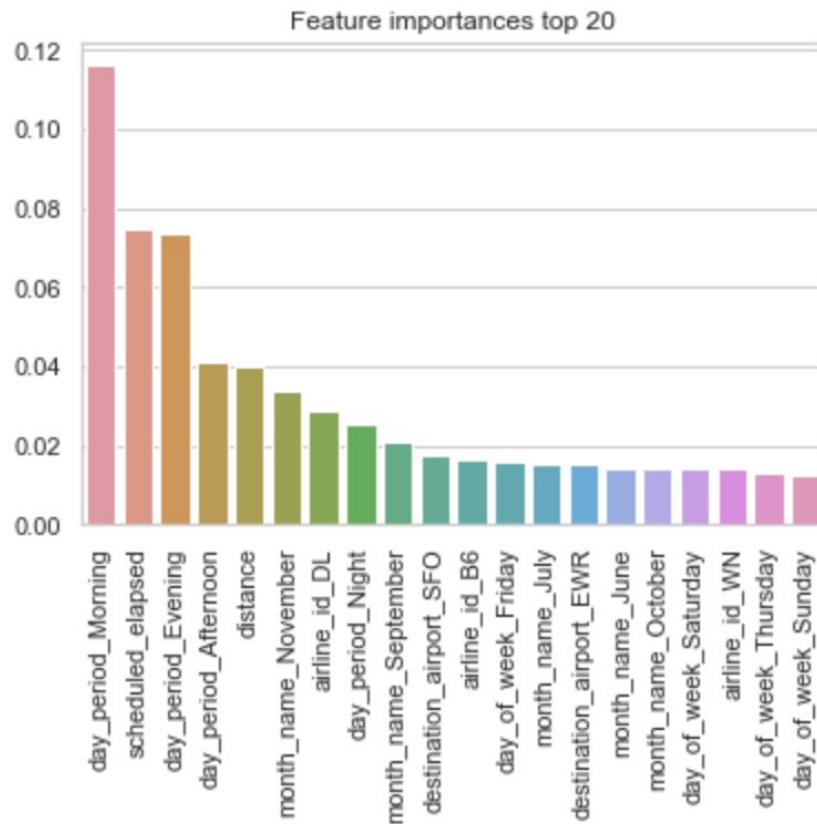
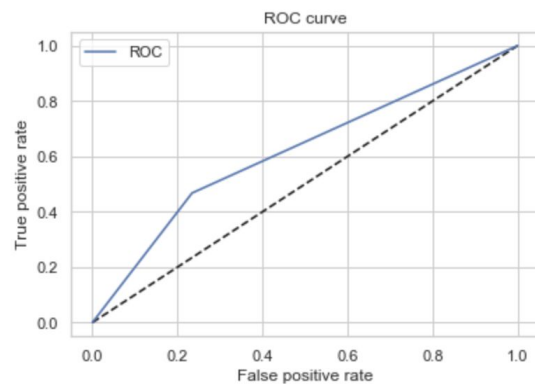
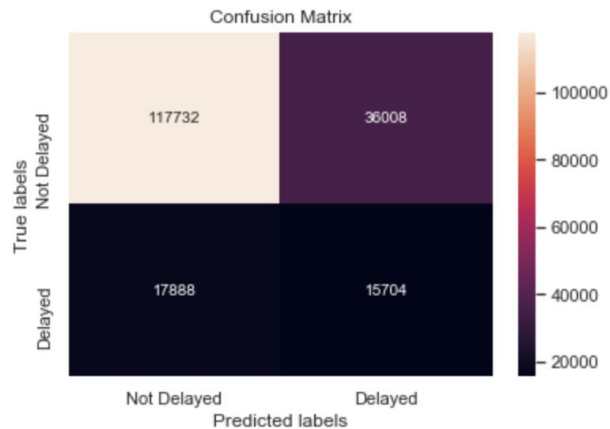
Time of day and **length of flight** most affect whether a flight will be delayed.

Our model recall is **60% for delayed** flights and **70% for non-delayed** flights.

To improve we would use **more data**, optimize **feature selection** and investigate differences between **types of delays**.

... to open for questions :)

Outputs :\



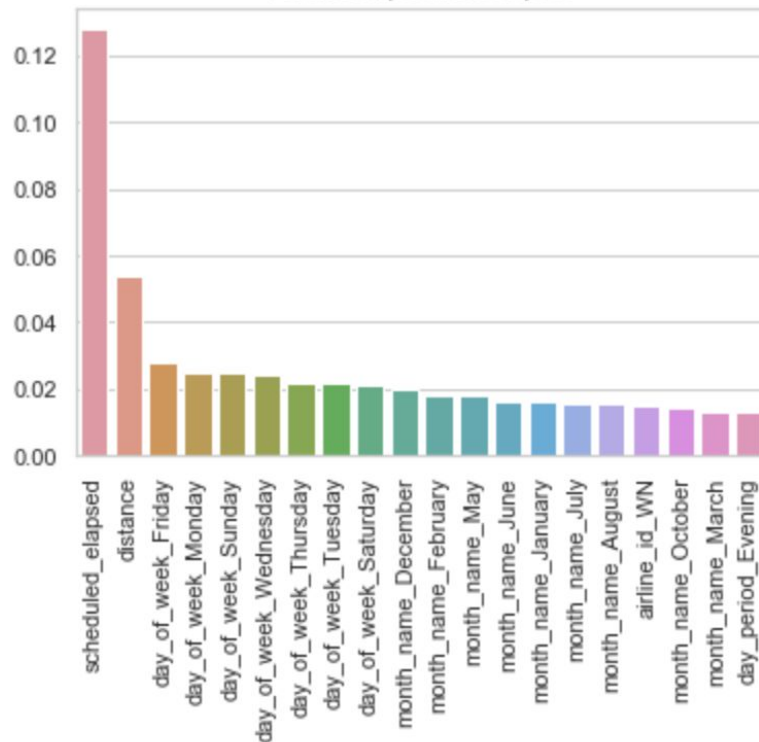
How Long would the Delay be ?

scheduled_elapsed
distance
day_of_week_Friday
day_of_week_Monday
day_of_week_Sunday
day_of_week_Wednesday
day_of_week_Thursday
day_of_week_Tuesday
day_of_week_Saturday
month_name_December
month_name_February
month_name_May
month_name_June
month_name_January
month_name_July
month_name_August
airline_id_WN
month_name_October
month_name_March
day_period_Evening

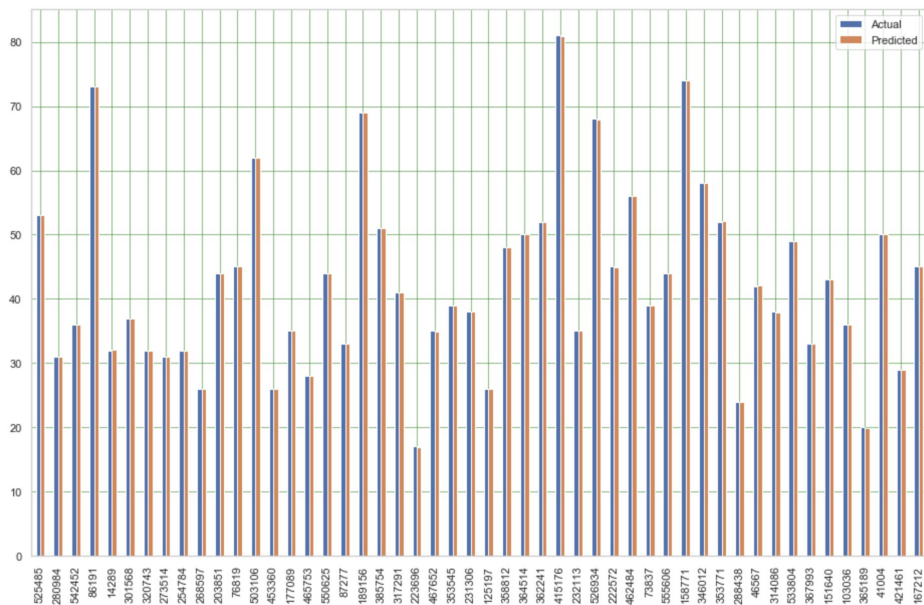
Important Factors

Scheduled duration
Distance
Day of week
Month

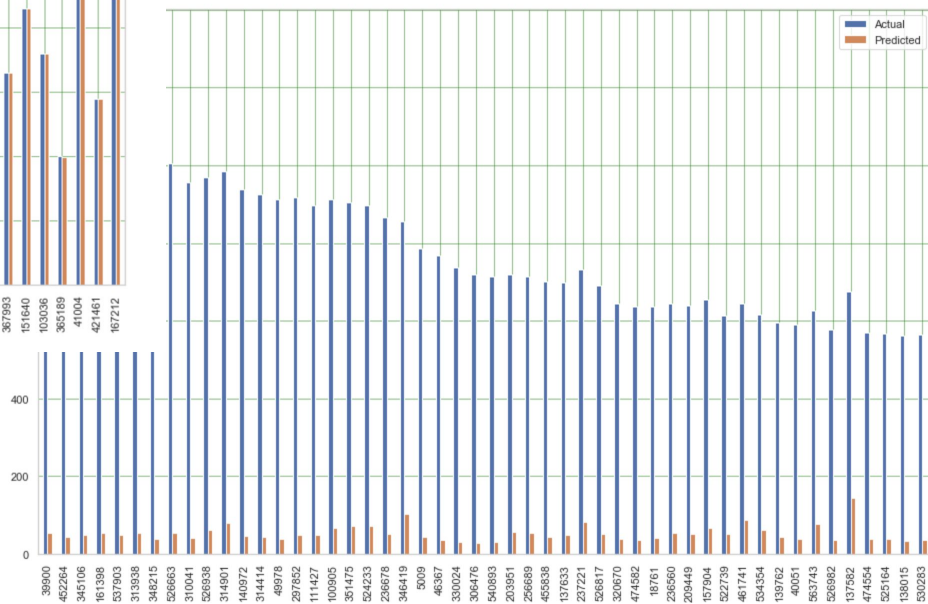
Feature importances top 20



Good At predicting short-term delay



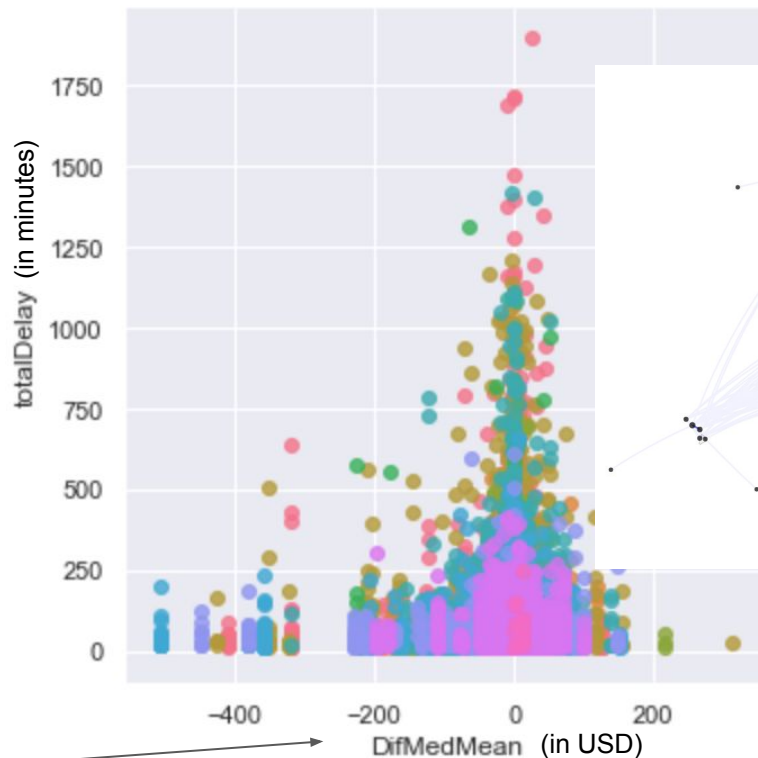
Not so well when the delay is long



Recommendation: Alternative Airports?

Run same on airports within a given distance-- still working on this code

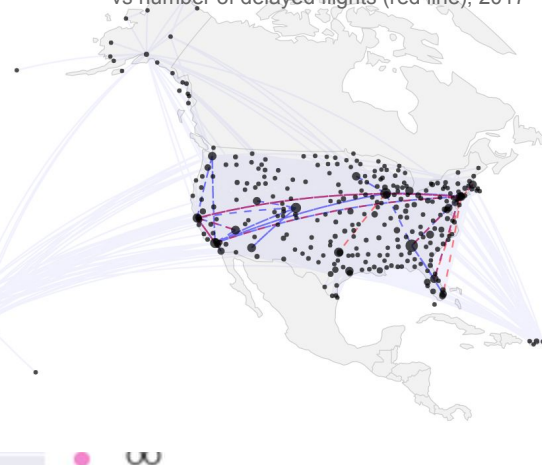
Route Comparison: Fare vs. Delay



We came up with a pricing metric that calculates the median price for a single flight and then compares it against the average median price for that flight path across airlines

Do some airlines
tend to have
more delays?
What about more

Map with main routes in number of flights (bold blue line)
vs number of delayed flights (red line), 2017



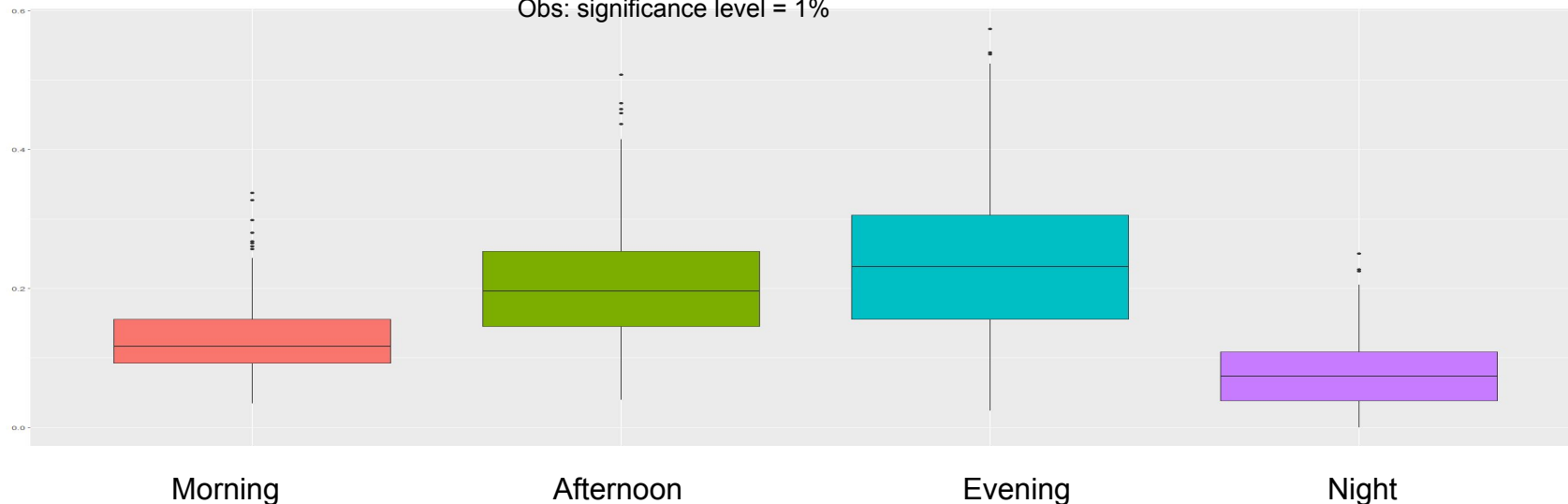
What factors affect delay?

Hypothesis Testing

In hypothesis testing with pairwise comparison procedures, we found out that all the periods of the day presented statistically significant difference in their proportions of delayed/cancelled flights.

	Morning	Afternoon	Evening
Afternoon	S		
Evening	S	S	
Night	S	S	S

Obs: significance level = 1%

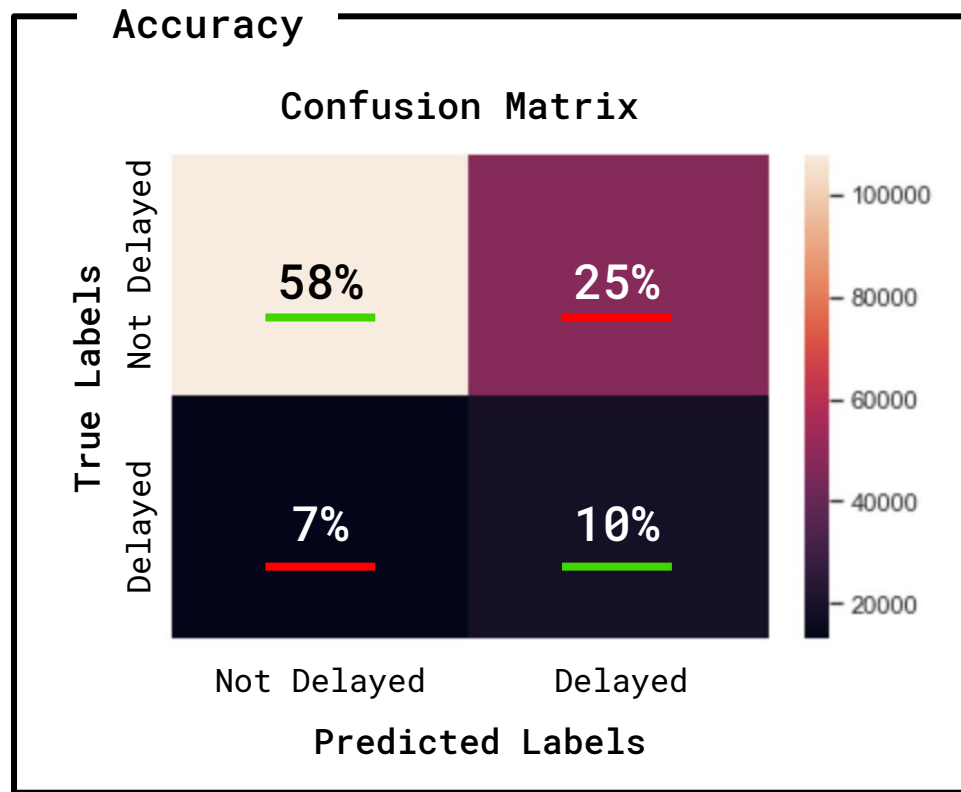


Can we predict instance of delay?

Random Forest

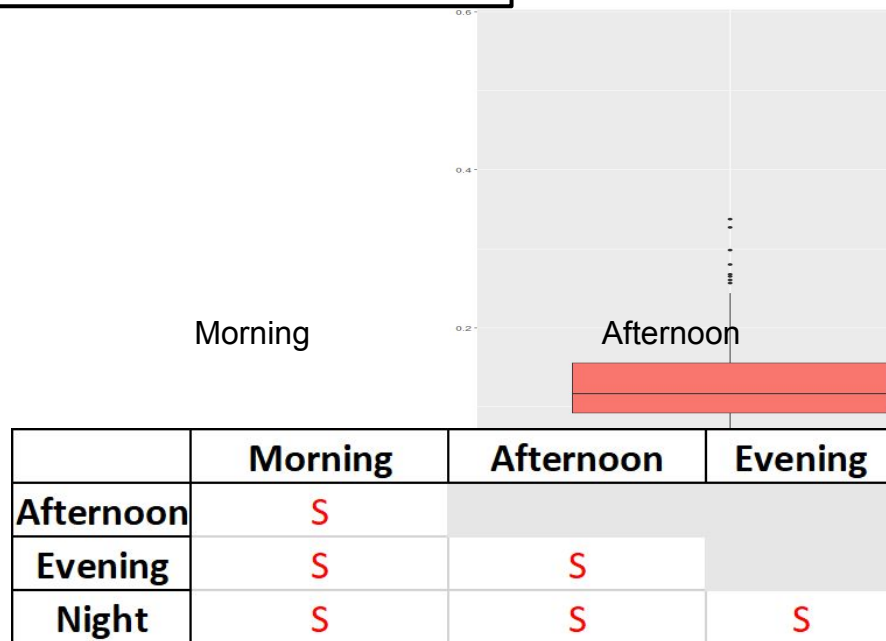
Delay: Y/N?

Right now, our
model has a lot of
false positives :/



What factors affect delay?

Hypothesis Testing



Evening

Night

presented statistically significant difference in their proportions of delayed/cancelled flights.

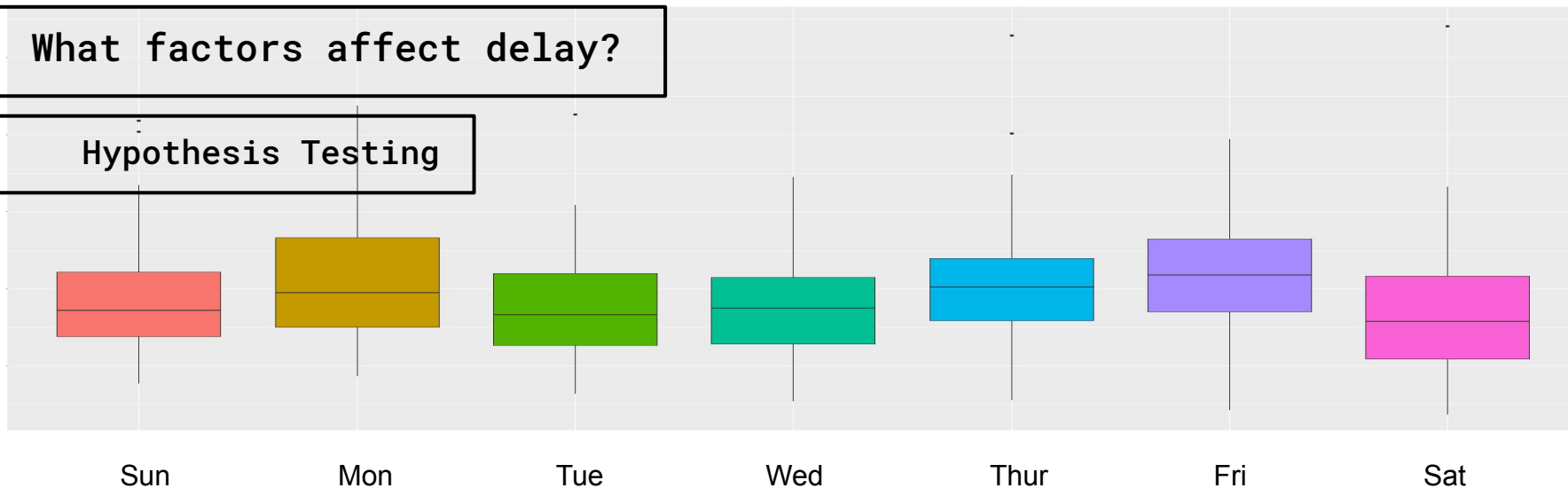
Obs: significance level = 1%

S: significant

NS: not significant

What factors affect delay?

Hypothesis Testing



	Sun	Mon	Tue	Wed	Thur	Fri
Mon	S					
Tue	NS	S				
Wed	NS	S	NS			
Thur	S	NS	S	S		
Fri	S	NS	S	S	S	
Sat	S	S	S	S	S	S

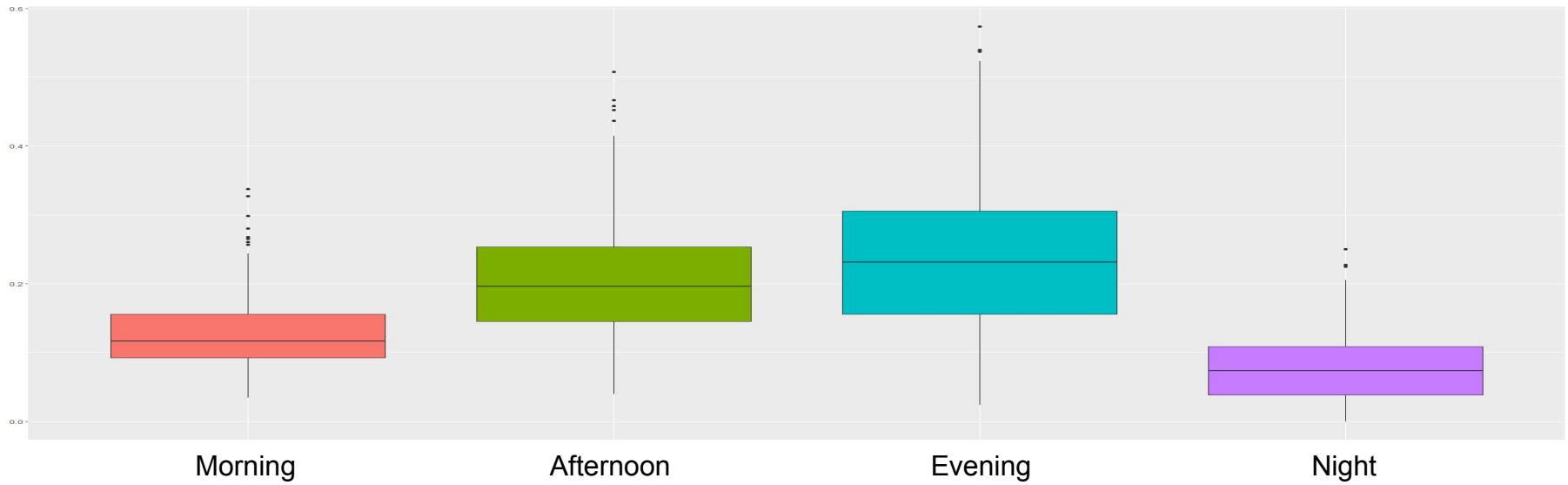
S: significant

NS: not significant

In hypothesis testing with pairwise comparison procedures, we found out that the days of the week presented statistically significant difference in their proportions of delayed/cancelled flights, except for:

- Sunday versus both Tuesday and Wednesday
- Monday versus Thursday and Friday
- Tuesday versus Wednesday

Obs: significance level = 1%



	Morning	Afternoon	Evening
Afternoon	S		
Evening	S	S	
Night	S	S	S

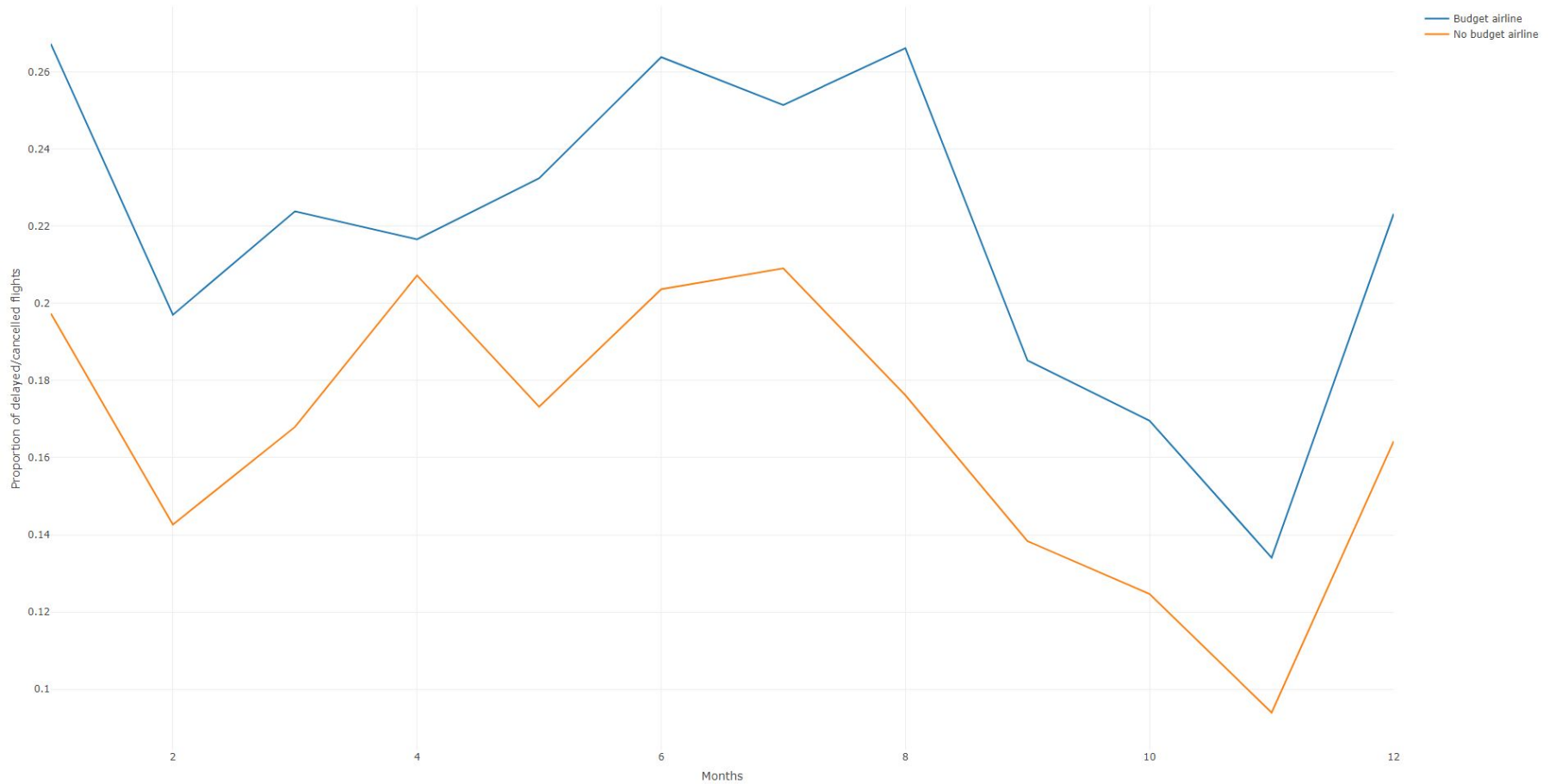
S: significant

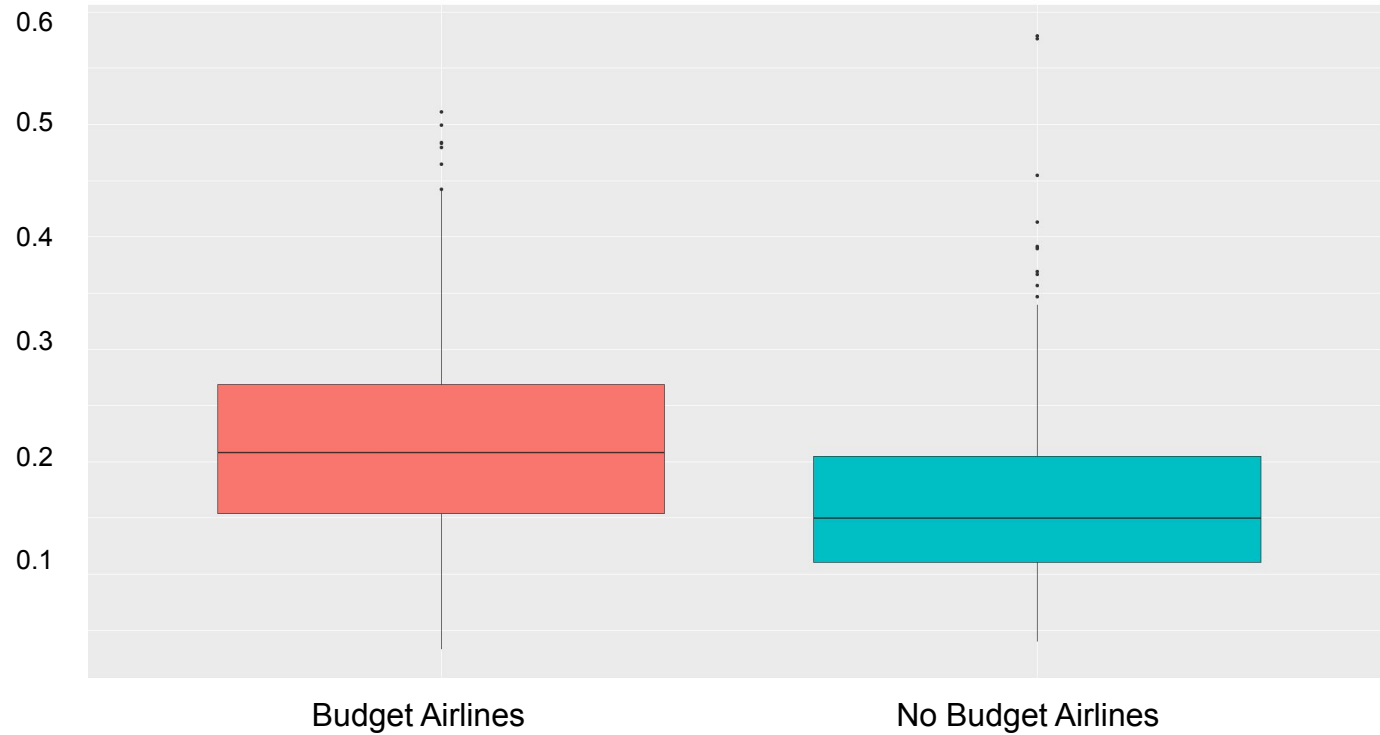
NS: not significant

Group 13

In hypothesis testing with pairwise comparison procedures, we found out that all the periods of the day presented statistically significant difference in their proportions of delayed/cancelled flights.

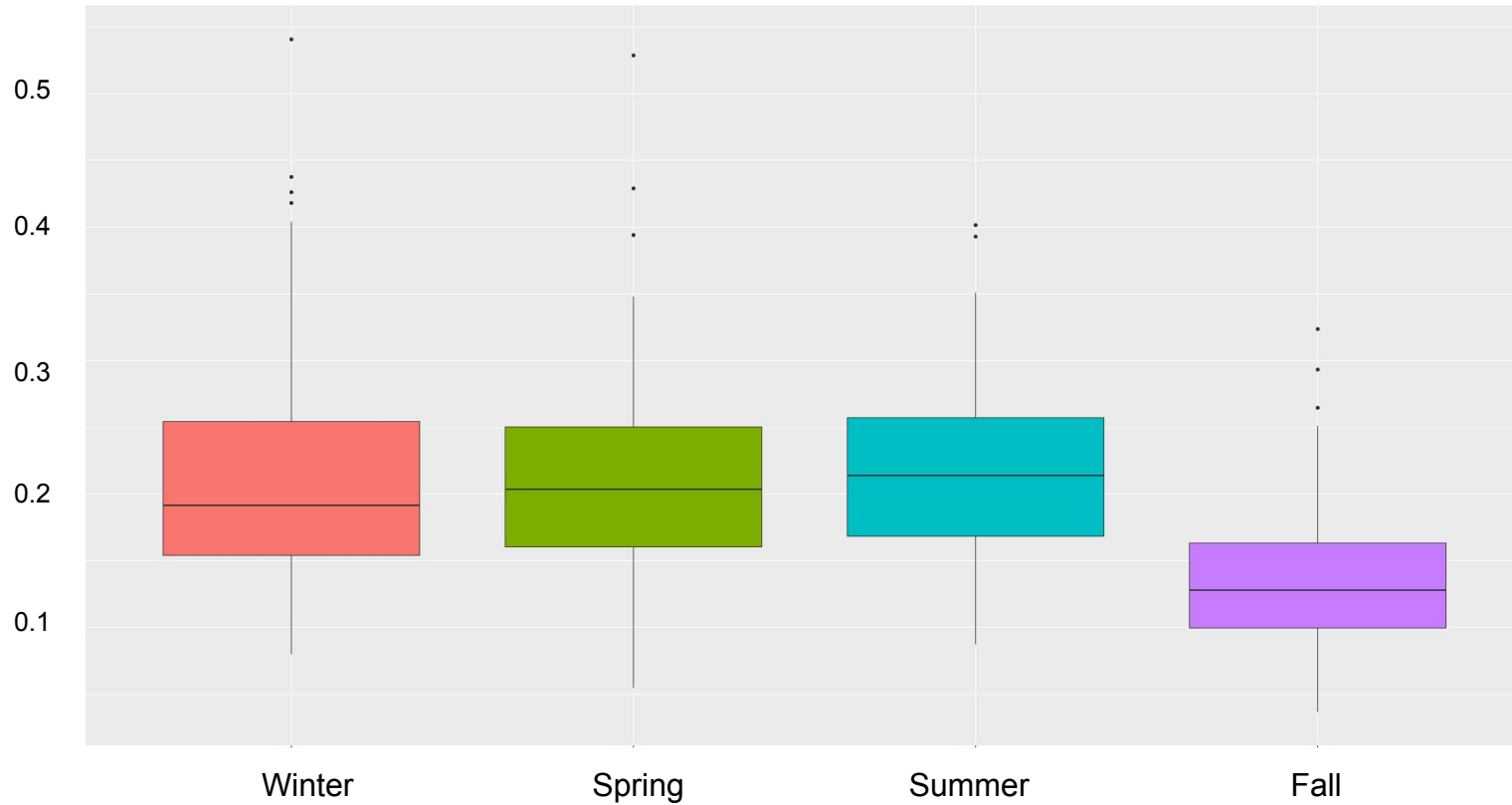
Obs: significance level = 1%





In hypothesis testing, we found out that budget airlines presented statistically significant difference in their proportions of delayed/cancelled flights.

Obs: significance level = 1%



In hypothesis testing with pairwise comparison procedures, we found out that Fall season flights presented statistically significant difference in their proportions of delayed/cancelled flights.

Obs: significance level = 1%

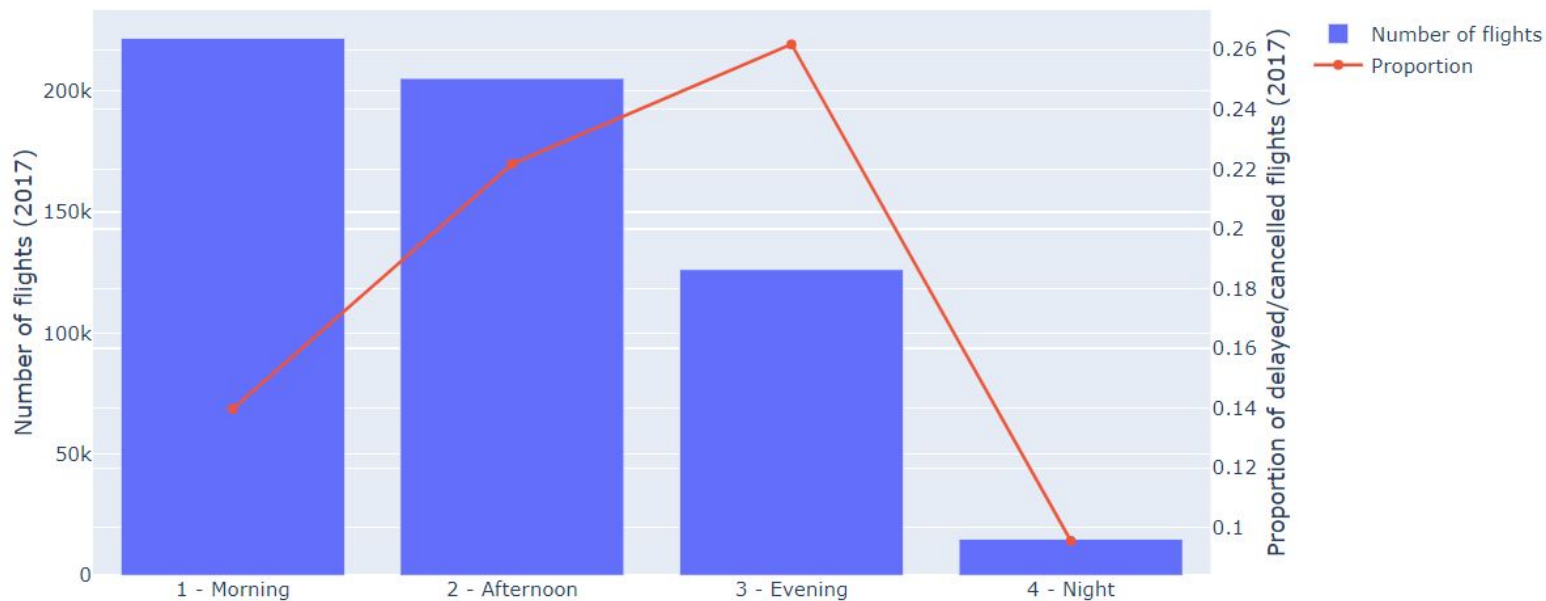
Create your own delay-adjusted fare

Some model for how much an individual is willing to pay to avoid/minute of delay

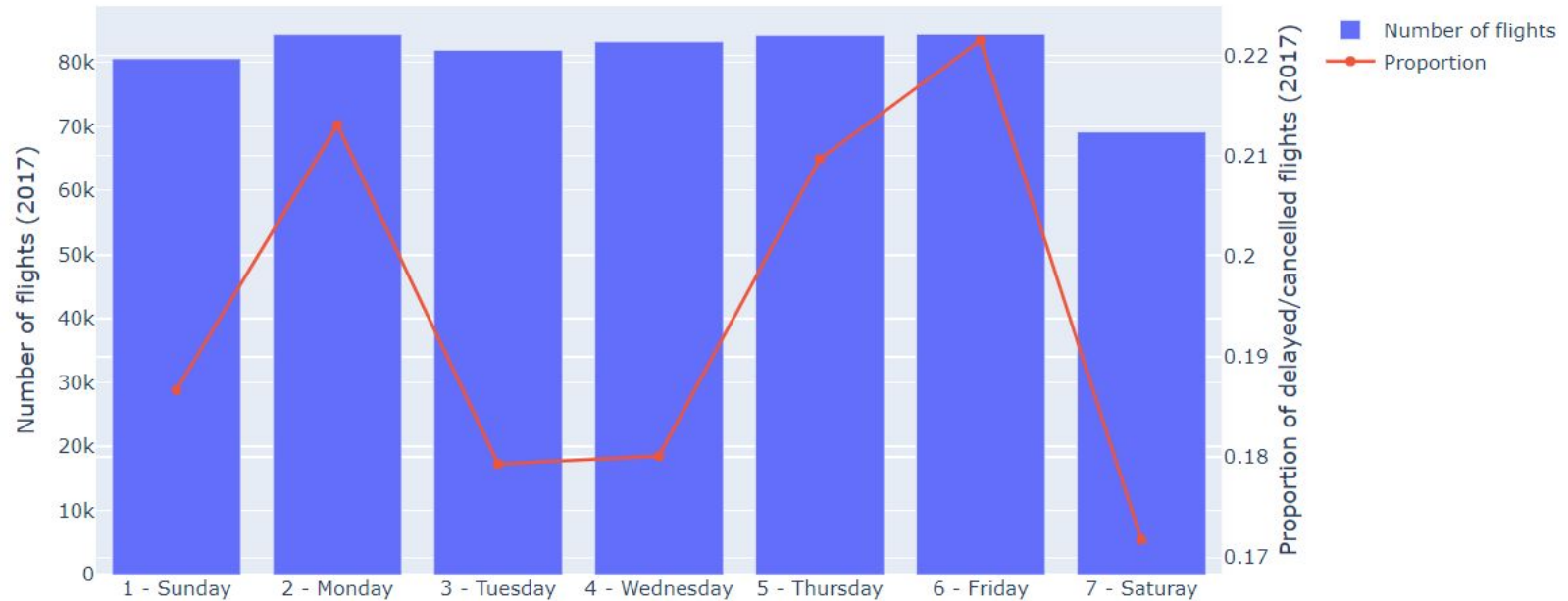
Model as a polynomial or exponential with some assumptions, user input

Exploratory Analysis

Proportional Summary*



Delay Comparison: What factors most affect delay?



Can we predict instance of delay?

Random Forest

Delay: Y/N?

Right now, our
model has a lot of
false positives :/

