An Investigation into
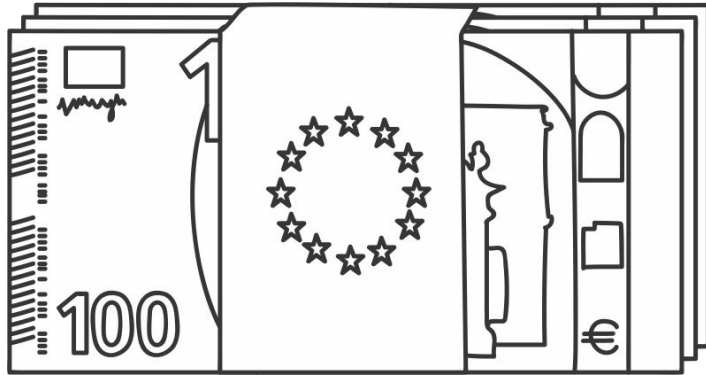
# Flight Delays

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

Flight delays are a hassle.

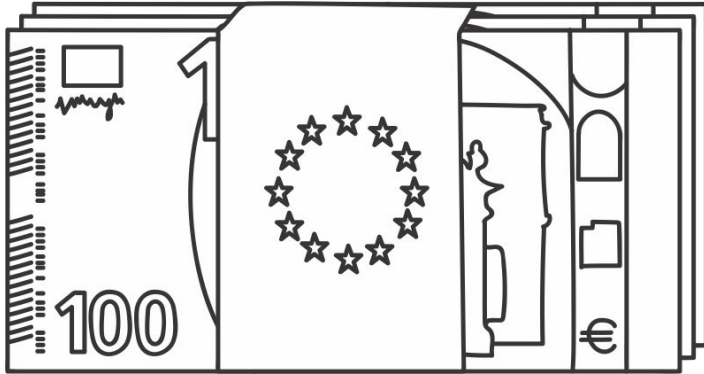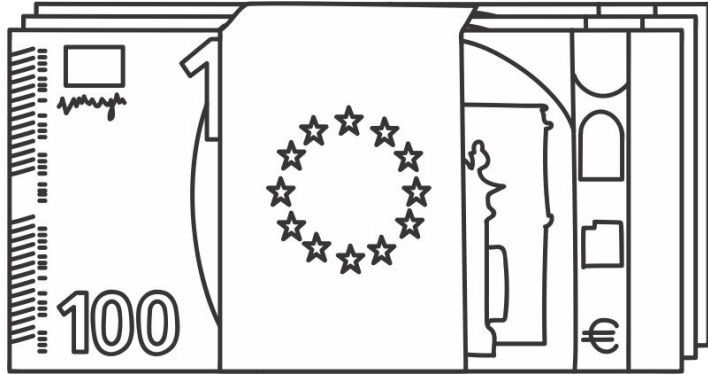EU regulation 261/2004 requires
**airlines to give you money**
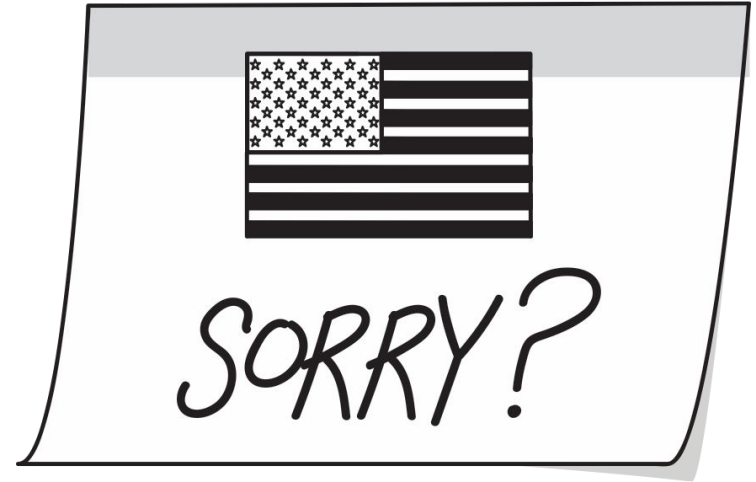if your flight is
delayed > 3 hours!

**Flight delays are a hassle**

What about in the US?

**Flight delays are a hassle**

VS

SORRY?

Group 13

Are delayed flights
a problem in the US?
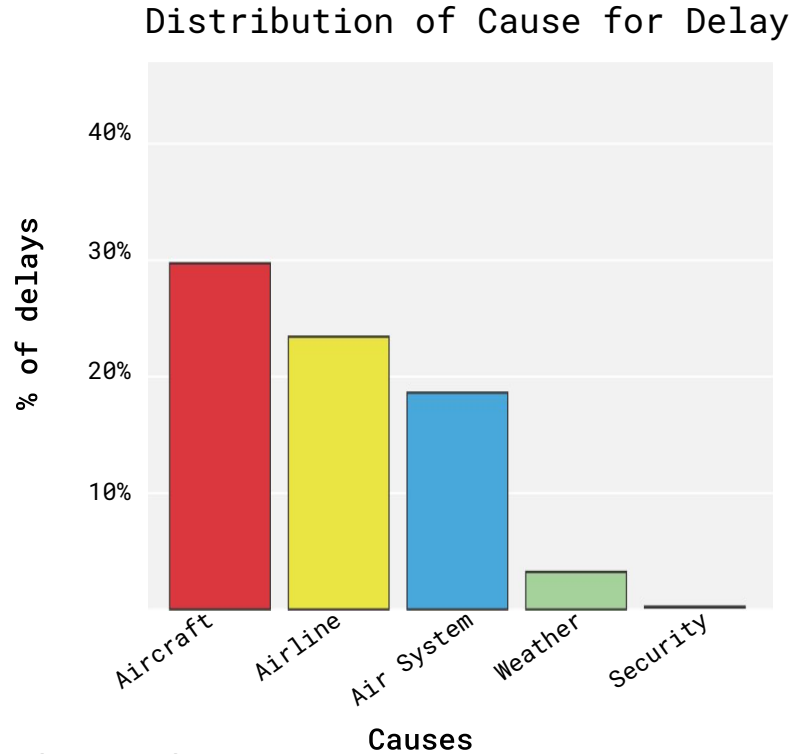
Are delayed flights
a problem in the US?

Heck, yes.

1 in 5 US flights
were delayed in 2017

Group 13

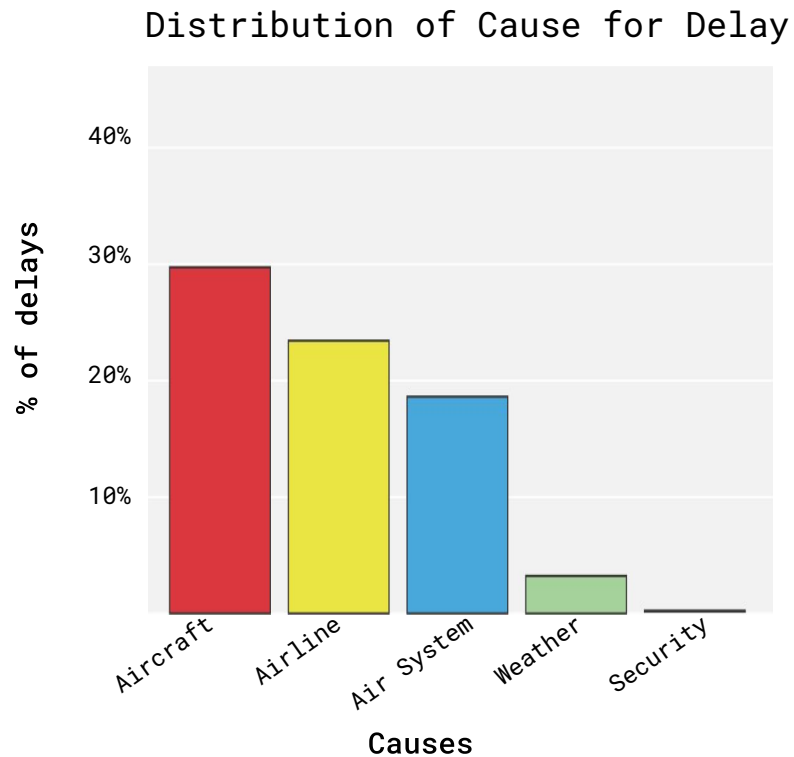# What factors correlate with delay?

# What factors affect delay?

Distribution of Cause for Delay



Group 13

# What factors affect delay?

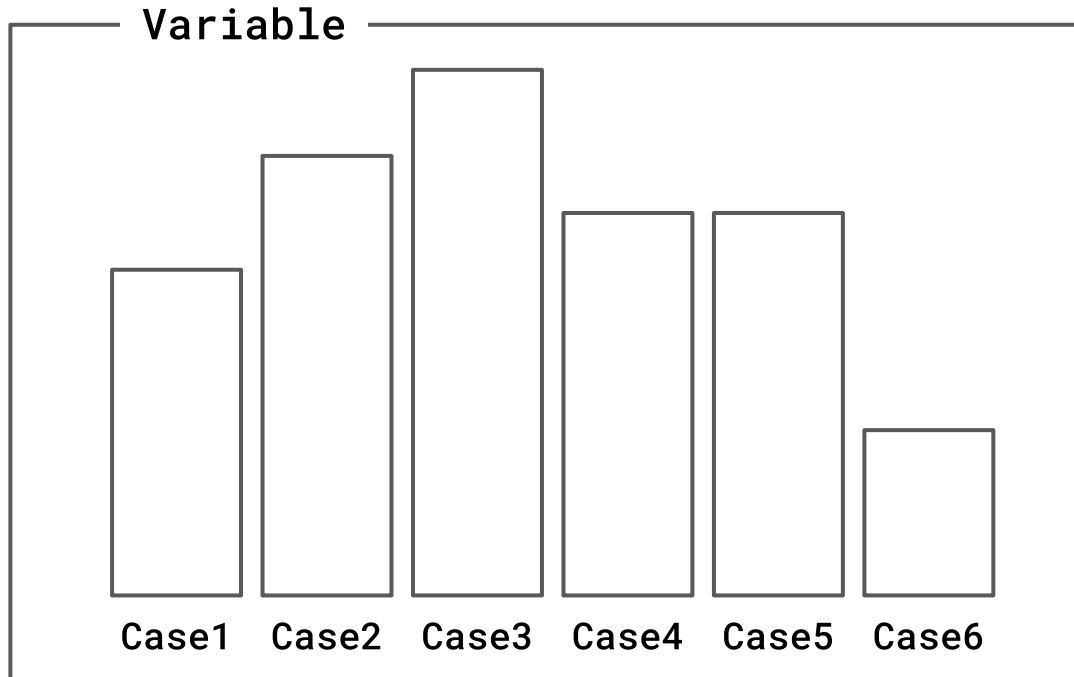## Distribution of Cause for Delay



Top Cause for Delay, by Airport



Weather database: missing airports for important cities, e.g. Chicago (ORD) and Dallas (DFW)
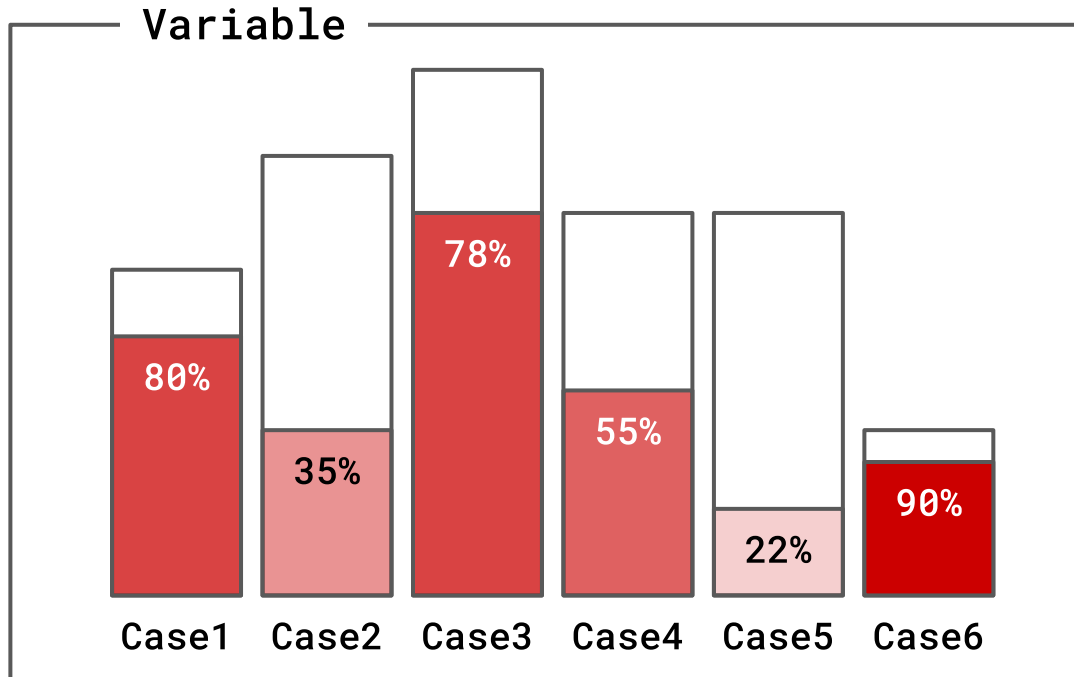
Group 13

# What factors affect delay?

## Proportional Summary*

**Variable**



Case1  Case2  Case3  Case4  Case5  Case6

Let's **separate by variables** we think might be relevant

# What factors affect delay?
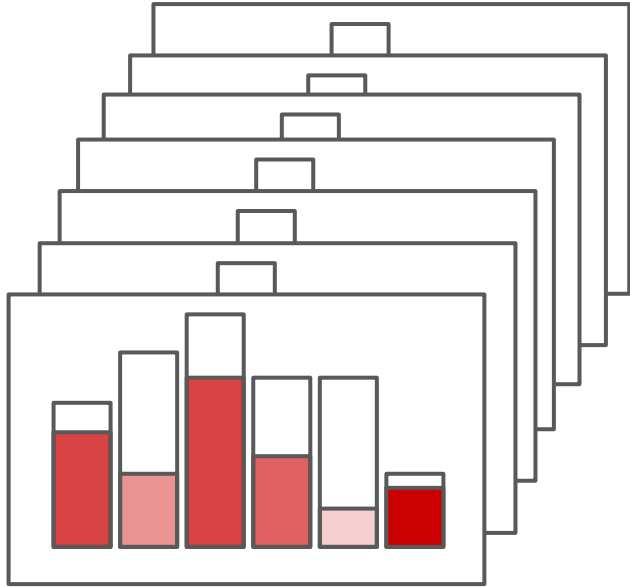
## Proportional Summary*



## Variables we evaluated

➜ Day of the week

➜ Month of the year

➜ Time of day (4 buckets)

➜ Elapsed flight time

➜ Distance of flight

➜ Airlines

➜ Season
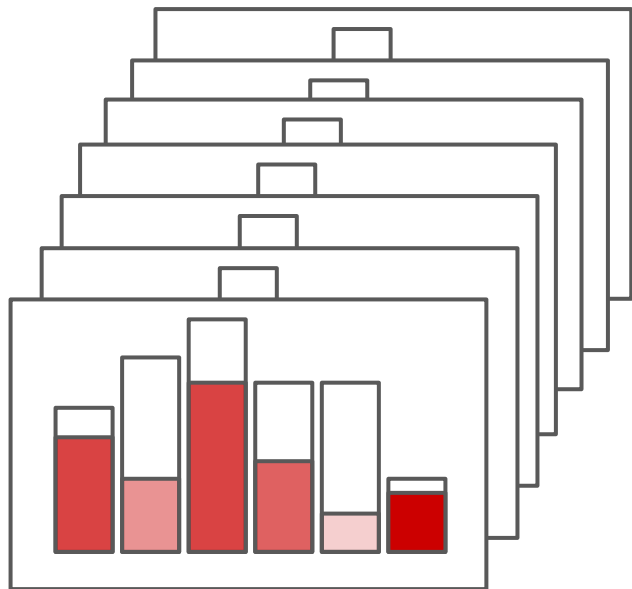
# What factors affect delay?

## Proportional Summary*



**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

Group 13

# What factors affect delay?

## Proportional Summary*

**Data Used**

Flight Traffic

**Data Not Used**

Outside → scope of question

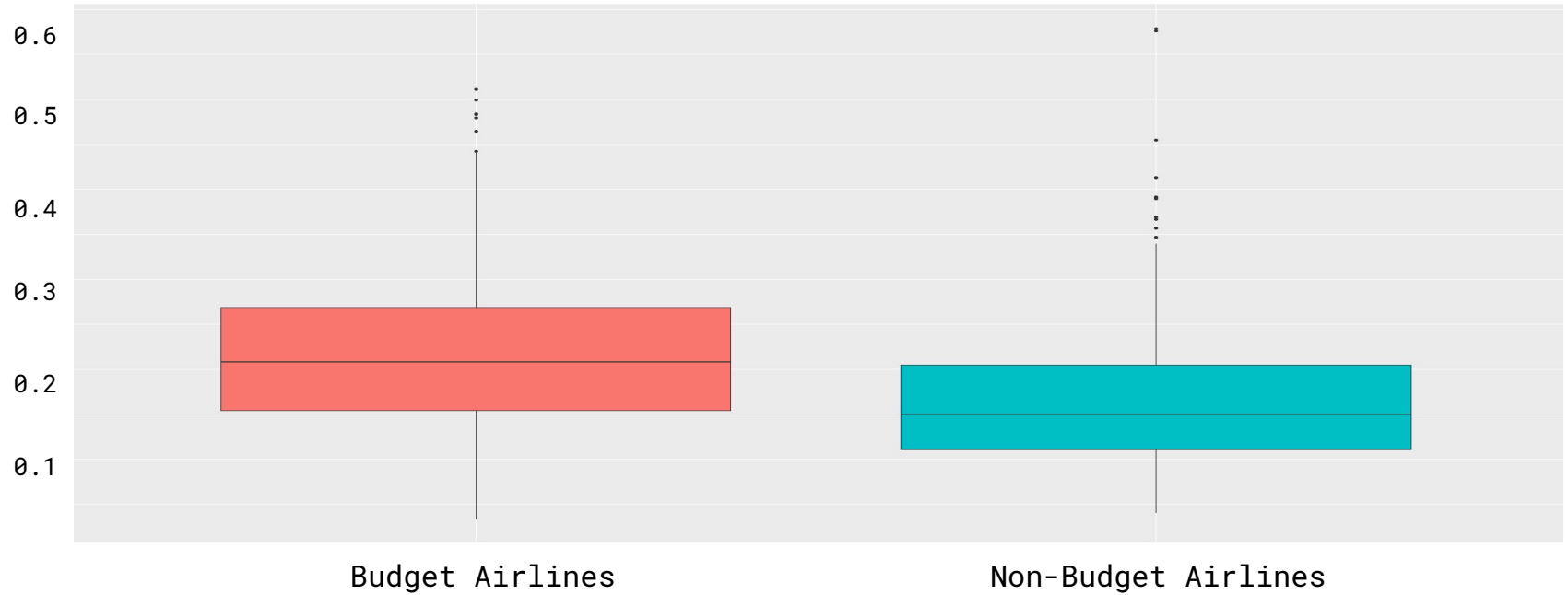Weather ← Retain important data points that were not represented in these dataframes

Fare

Event ←

## 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

## Airline Type



| | 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
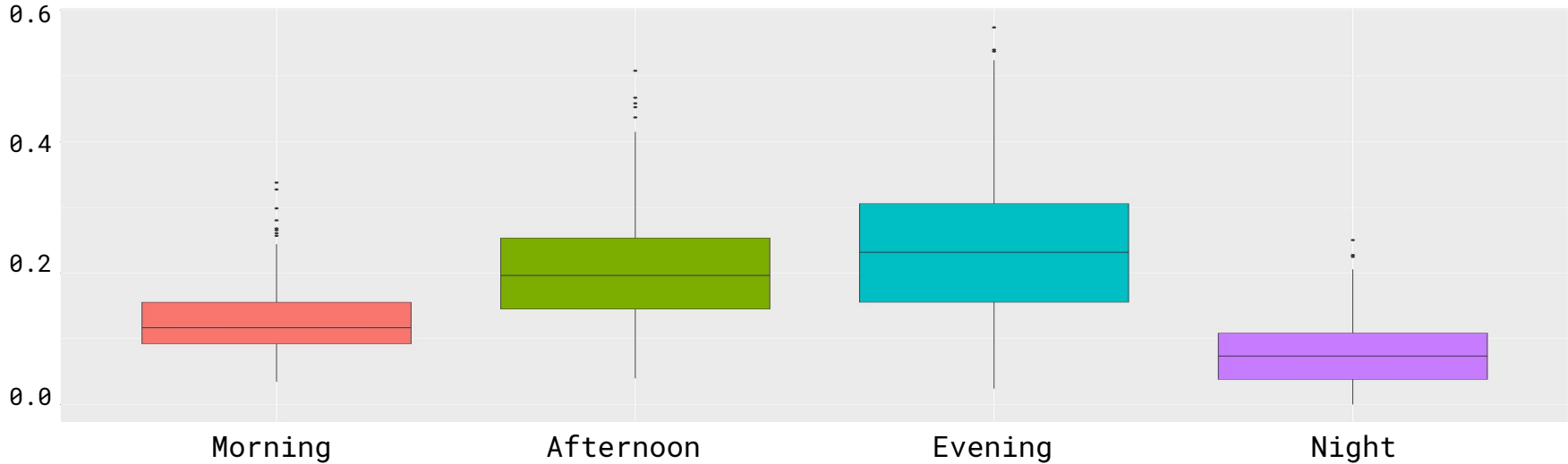➔ ExpressJet
➔ Frontier
➔ SkyWest
➔ Southwest
➔ Virgin

### "Non-Budget" Airlines:

➔ American
➔ Delta
➔ Hawaiian Air
➔ United
➔ Alaska Air

What factors affect delay?

Hypothesis Testing

Time of Day

Group 13

# 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**



|  | Spring | Summer | Autumn |
|---|---|---|---|
| Summer | NS |  |  |
| Autumn | S | 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 |

Group 13

# What factors affect delay?

## Hypothesis Testing

## Day of the Week



Group 13

# 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 **instance** of delay?

Random Forest

## Delay: Y/N?

Features

Random Forest

# Length of Delay?

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



Features

# Can we predict **length** of delay?
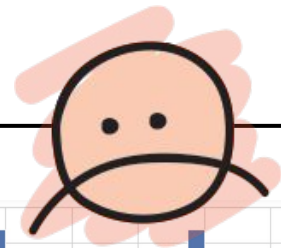


Short Delays

# Can we predict **length** of delay?



Short Delays

Long Delays

Group 13

## Practical Application

(So What?)

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

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.



Group 13

Practical Application

(So What?)

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

Group 13

# 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 :)