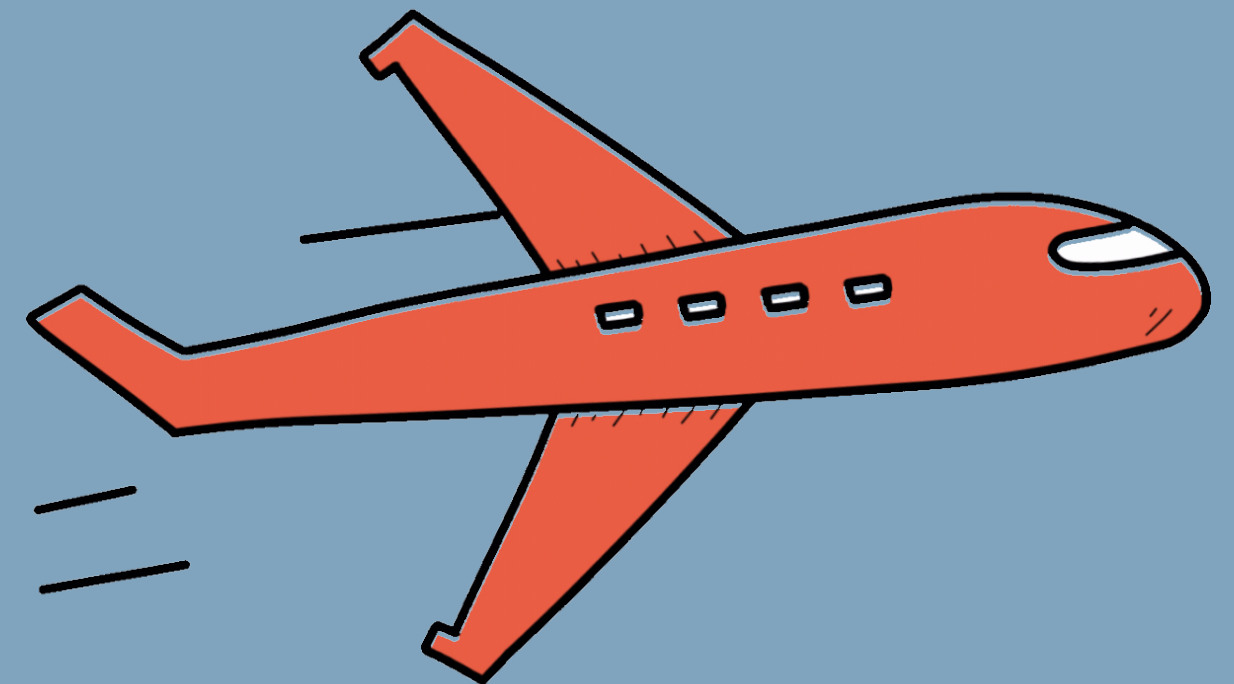


# Predicting Flight Delays

Caleb Ward and  
Jessica Moloney

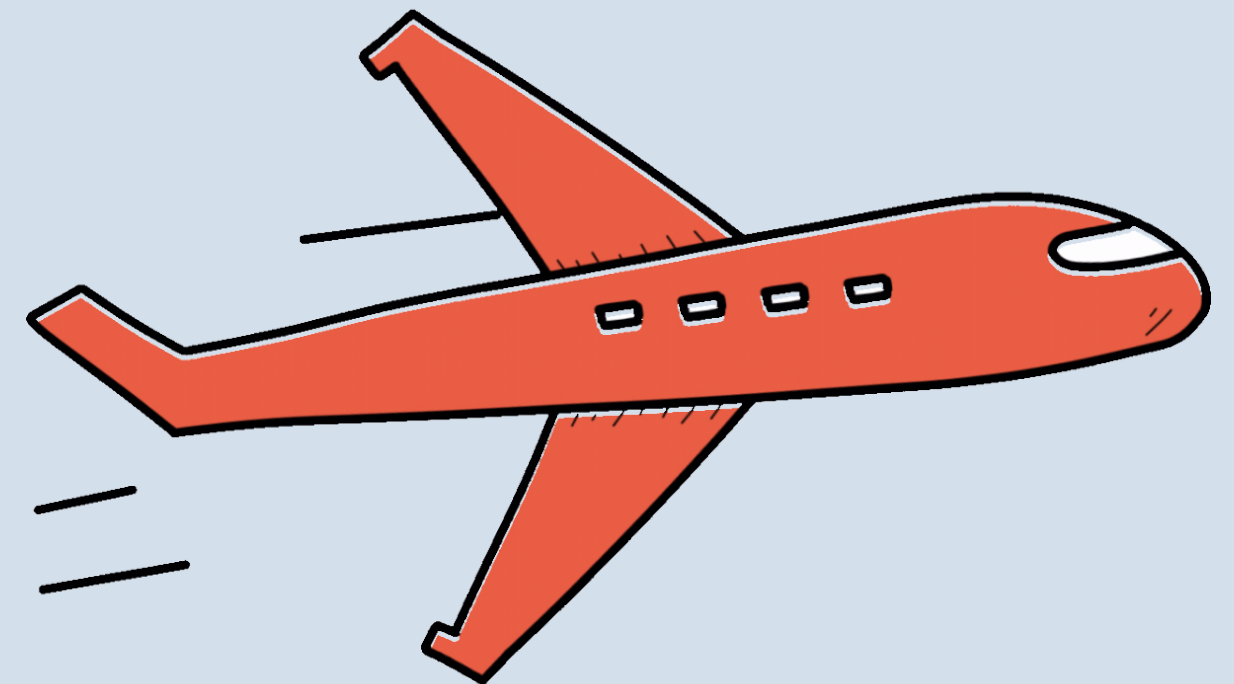
# Problem

Can we predict how delayed flights will be  
one week in advance?



# Data

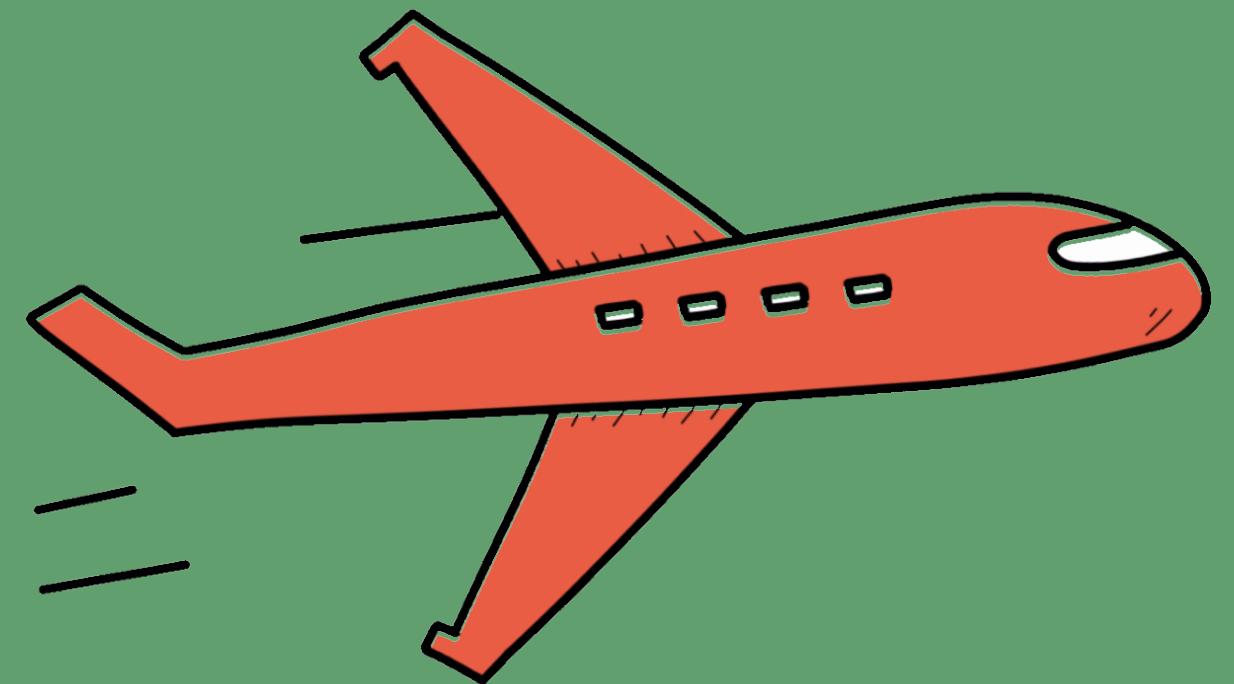
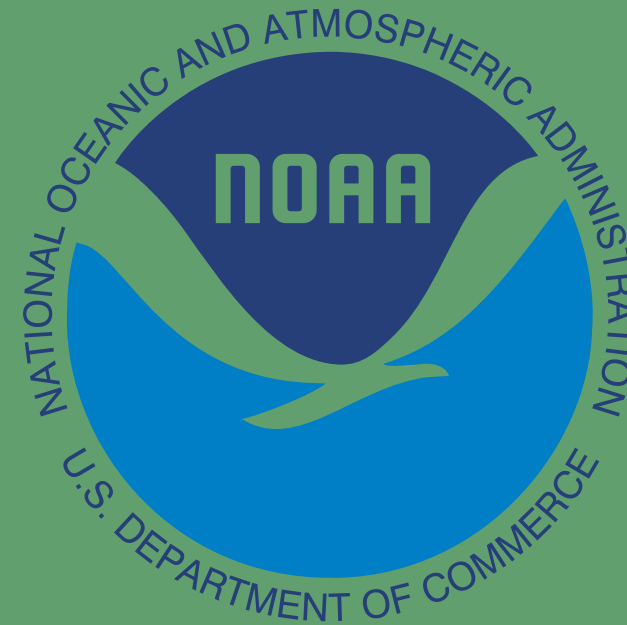
- 15 million records of US flights in 2018 and 2019
- Information about carrier, flight number, tail number, planned departure / arrival time, origin / destination, actual delays





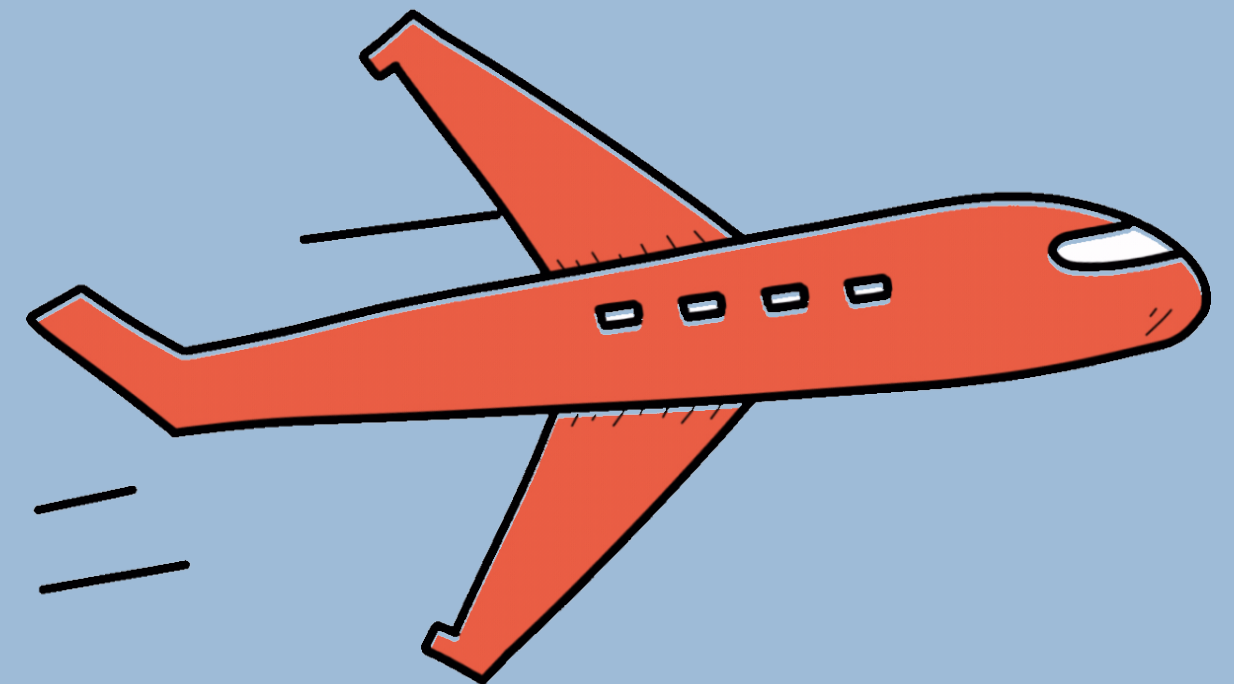
# Additional Sources

- NOAA - NWS Climate Prediction Center
- Government Source
- Historical forecasts, not historical weather records
- Exhaustive coverage
- Easily accessible and free

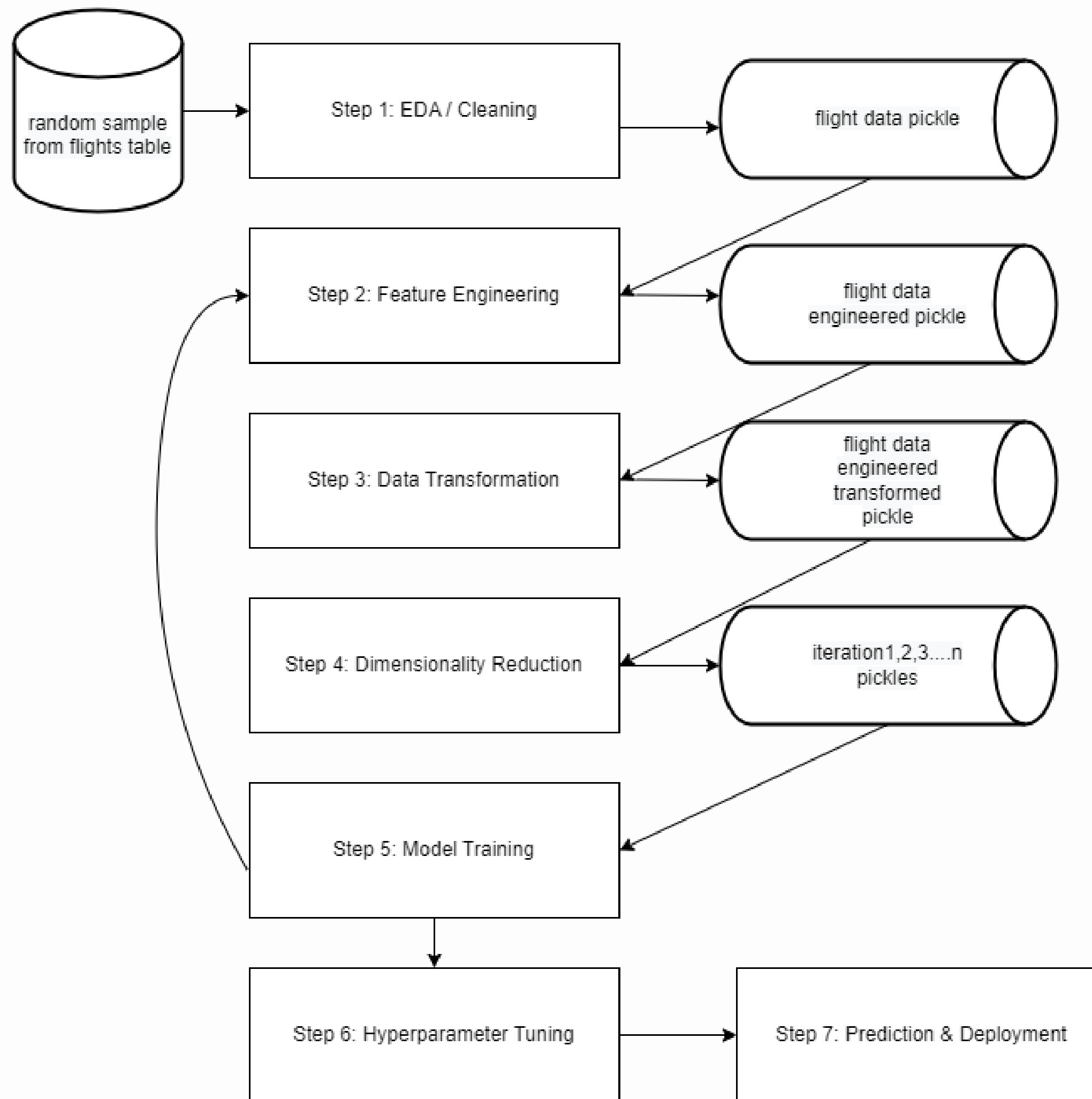


# Goal

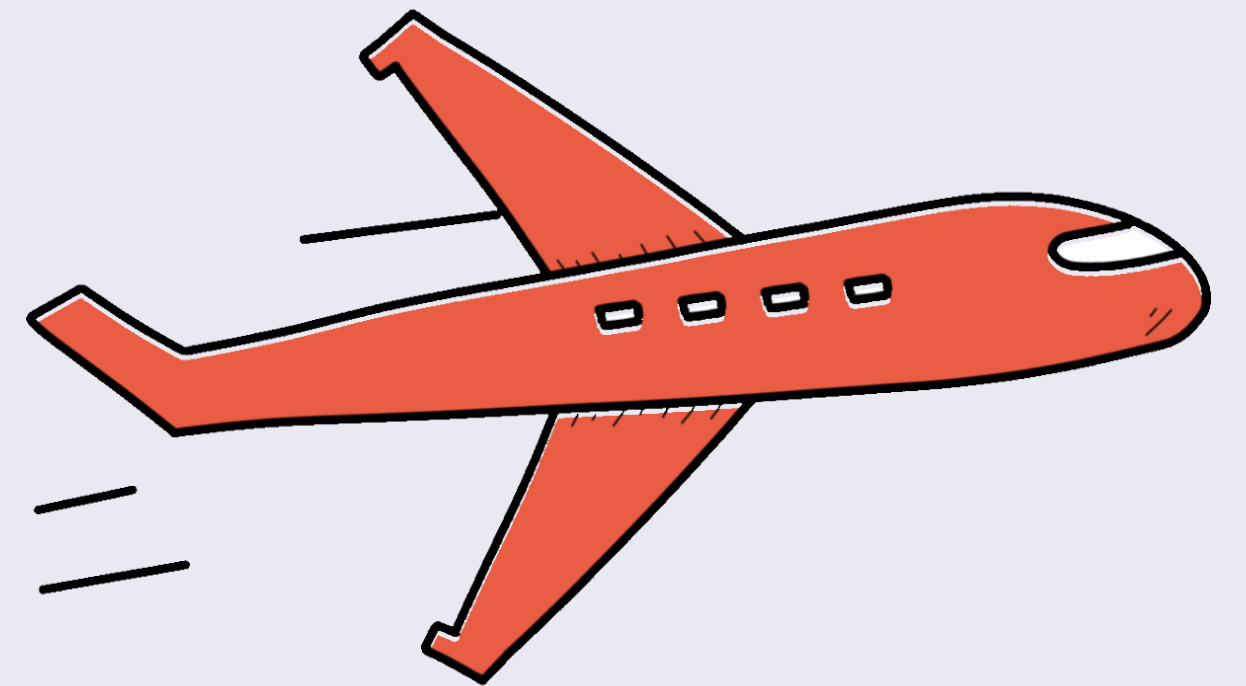
- Use machine learning to find patterns in the data and predict future delays



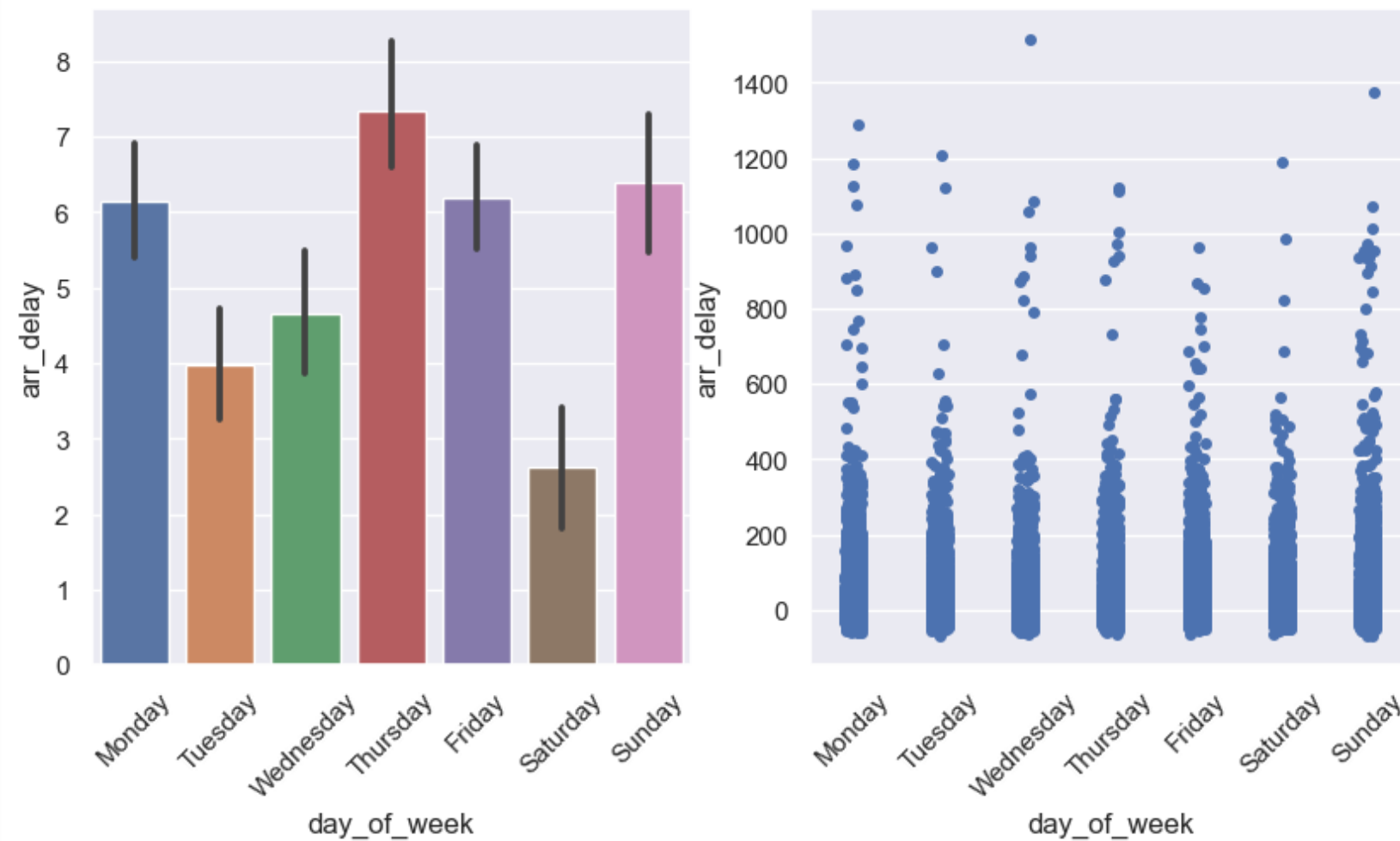
# Approach



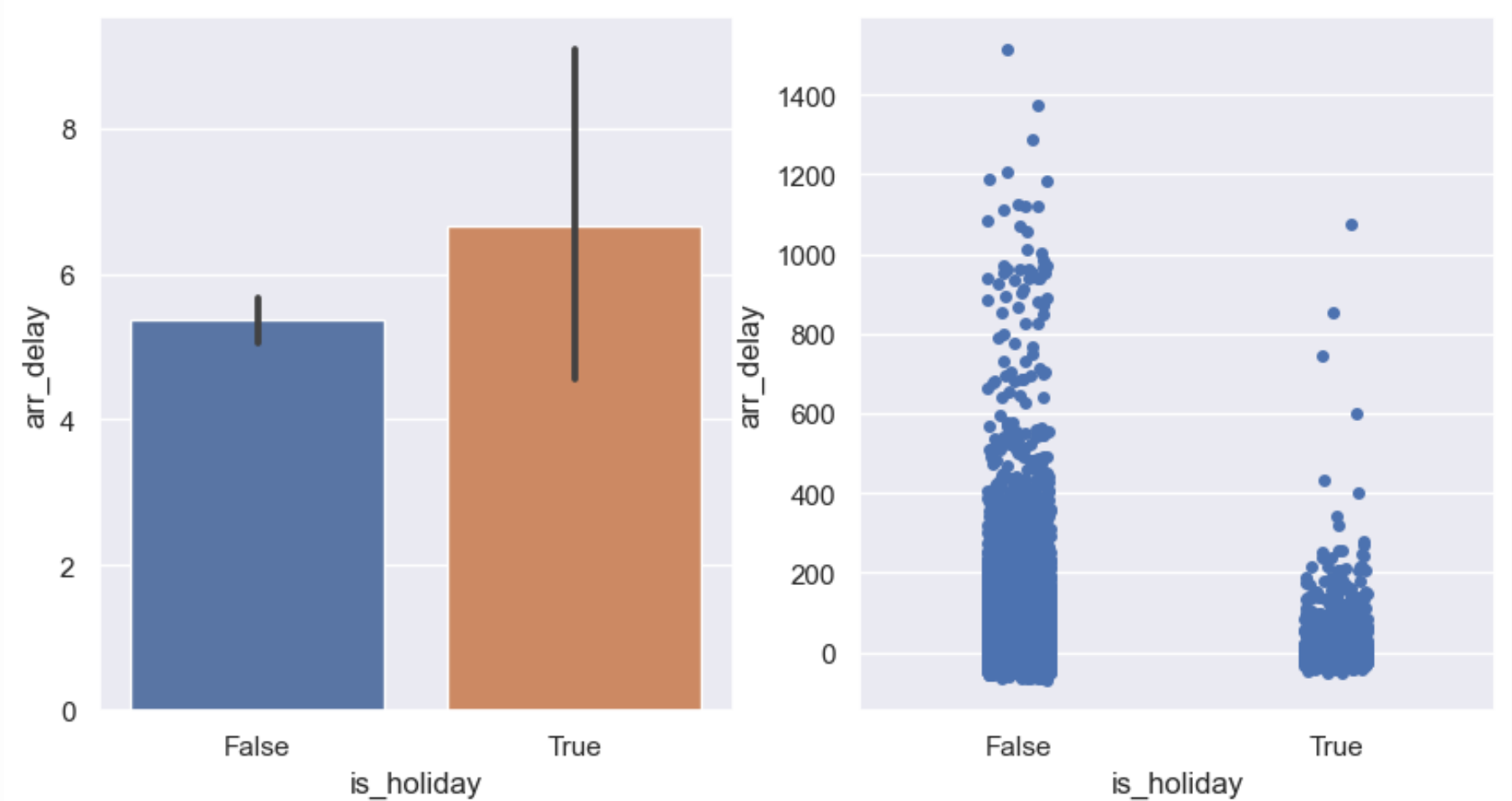
# Features based on Flight Date



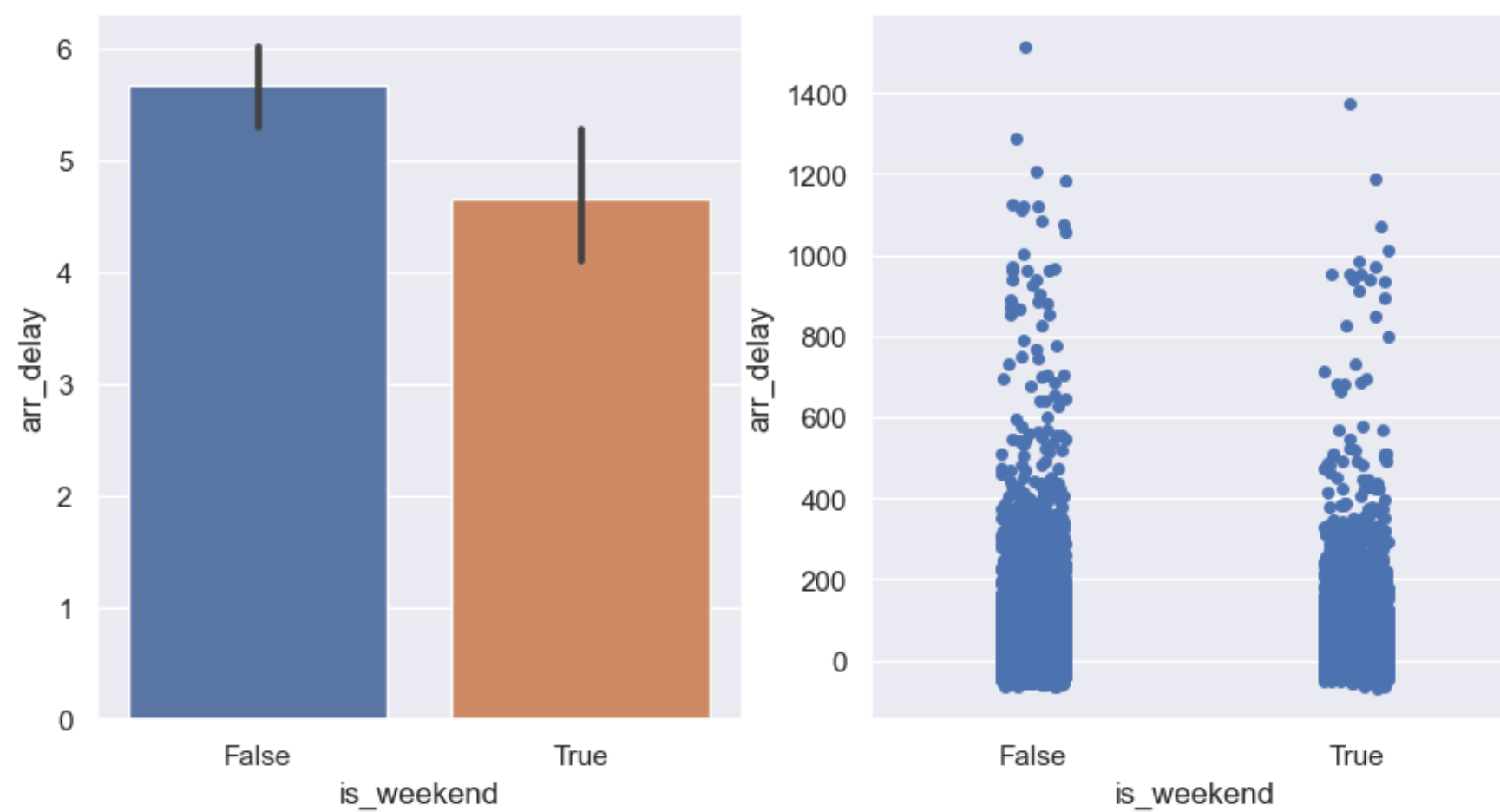
Arrival Delay by Day of Week



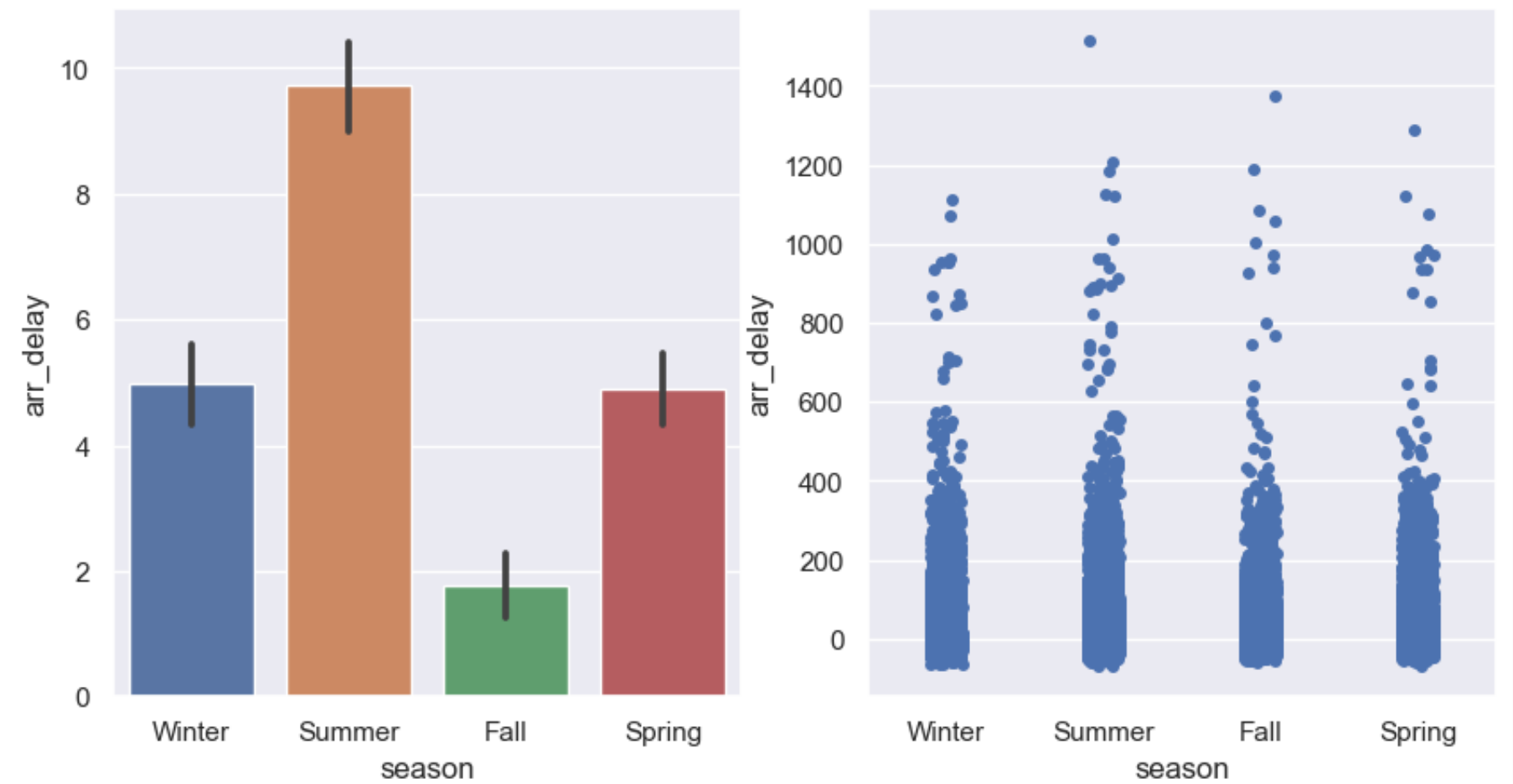
Arrival Delay by Normal Day vs. Holiday



Arrival Delay by Weekday vs. Weekend

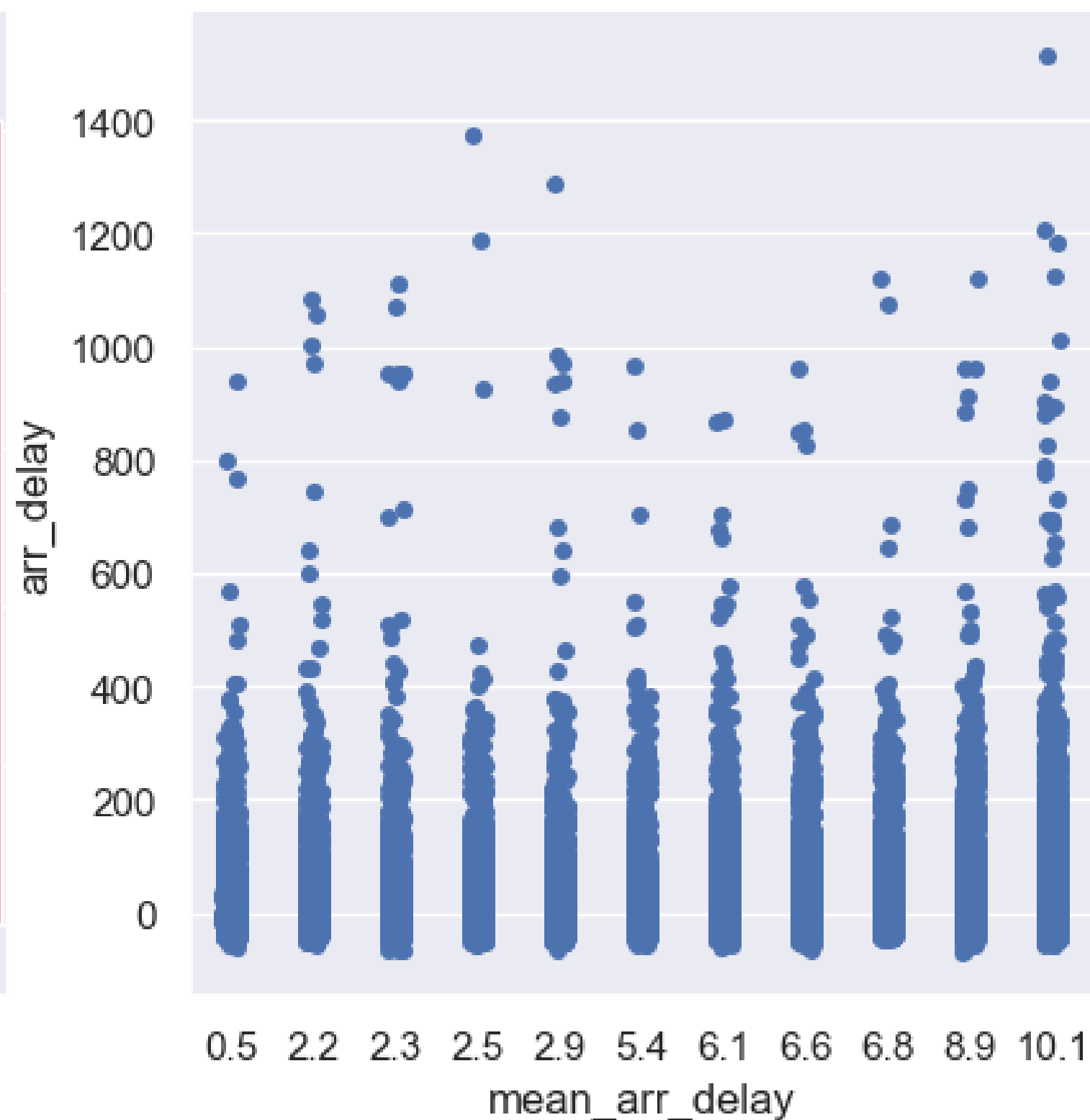
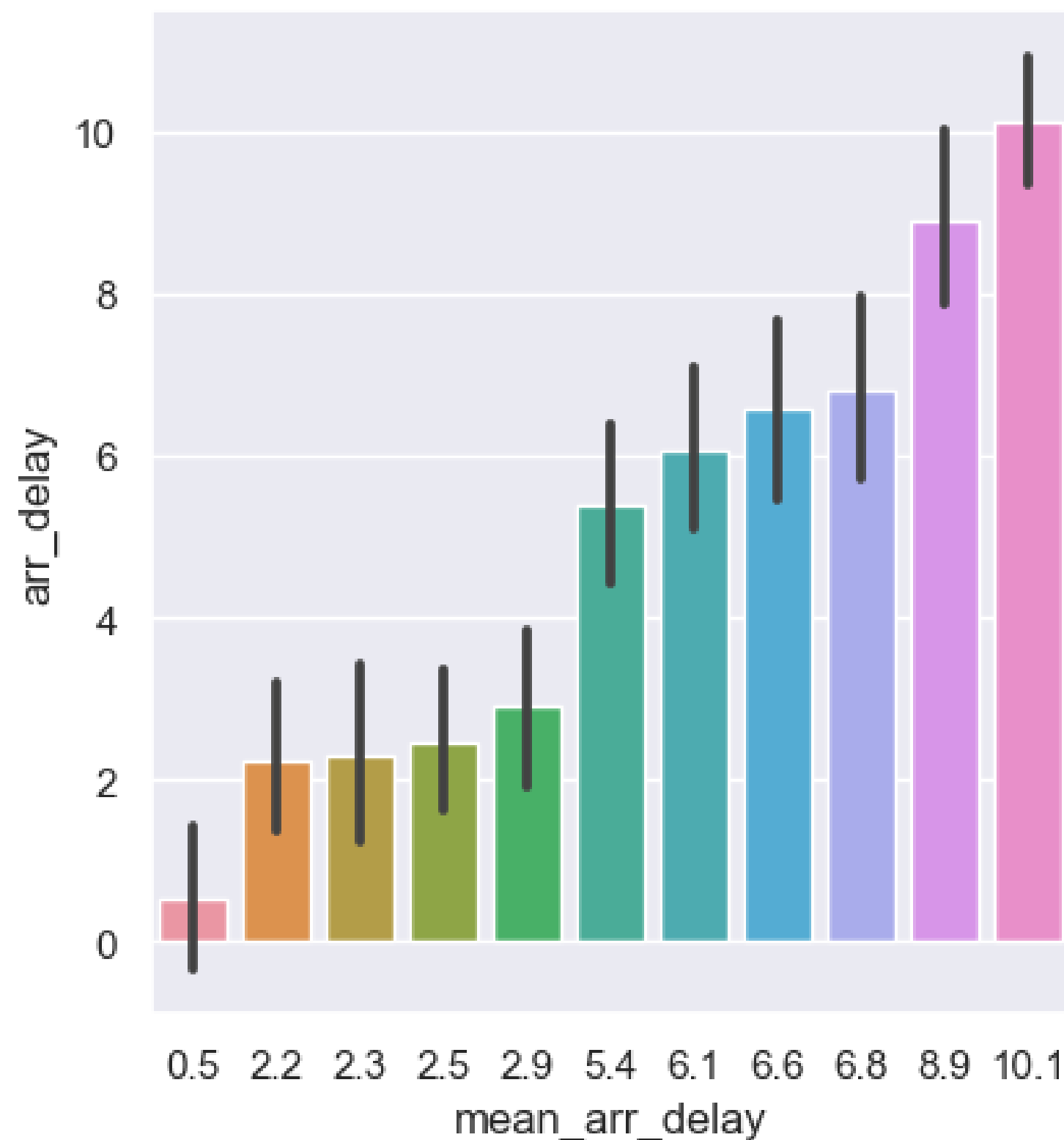


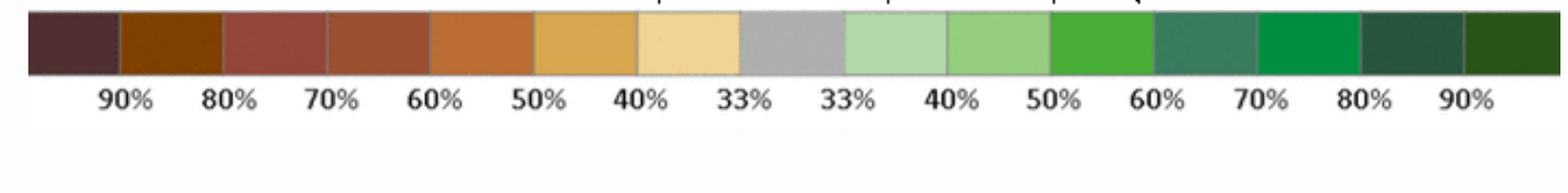
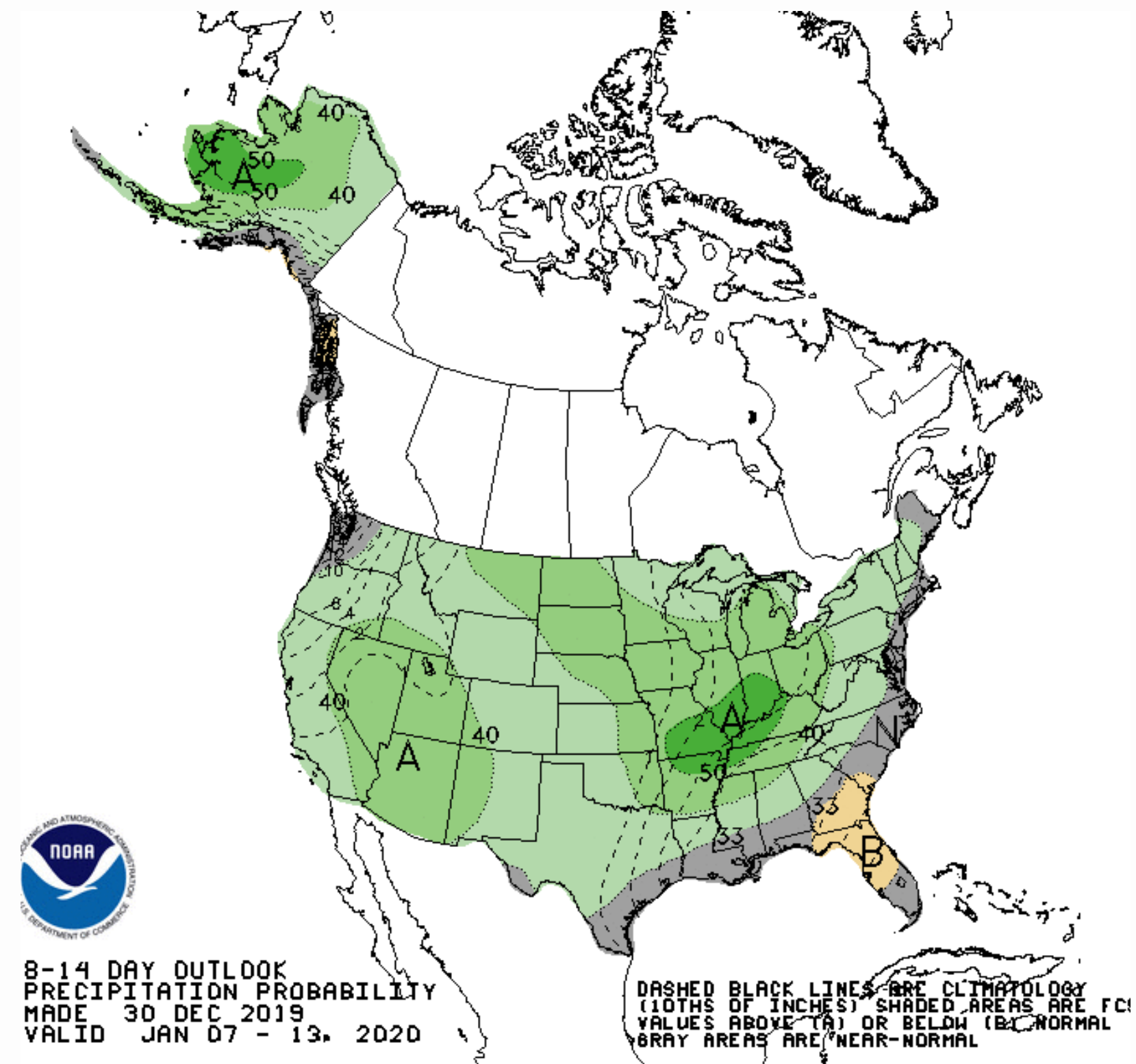
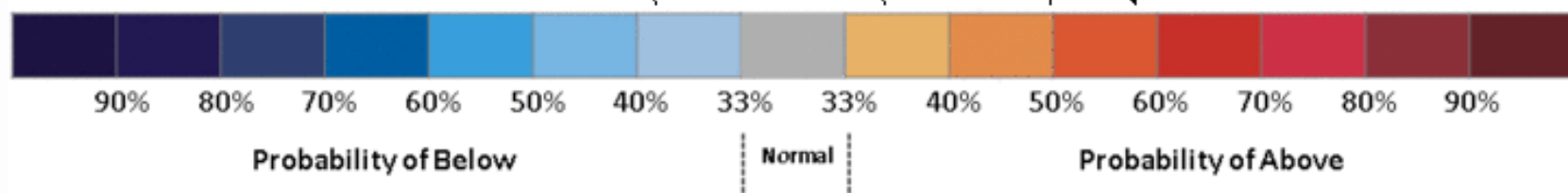
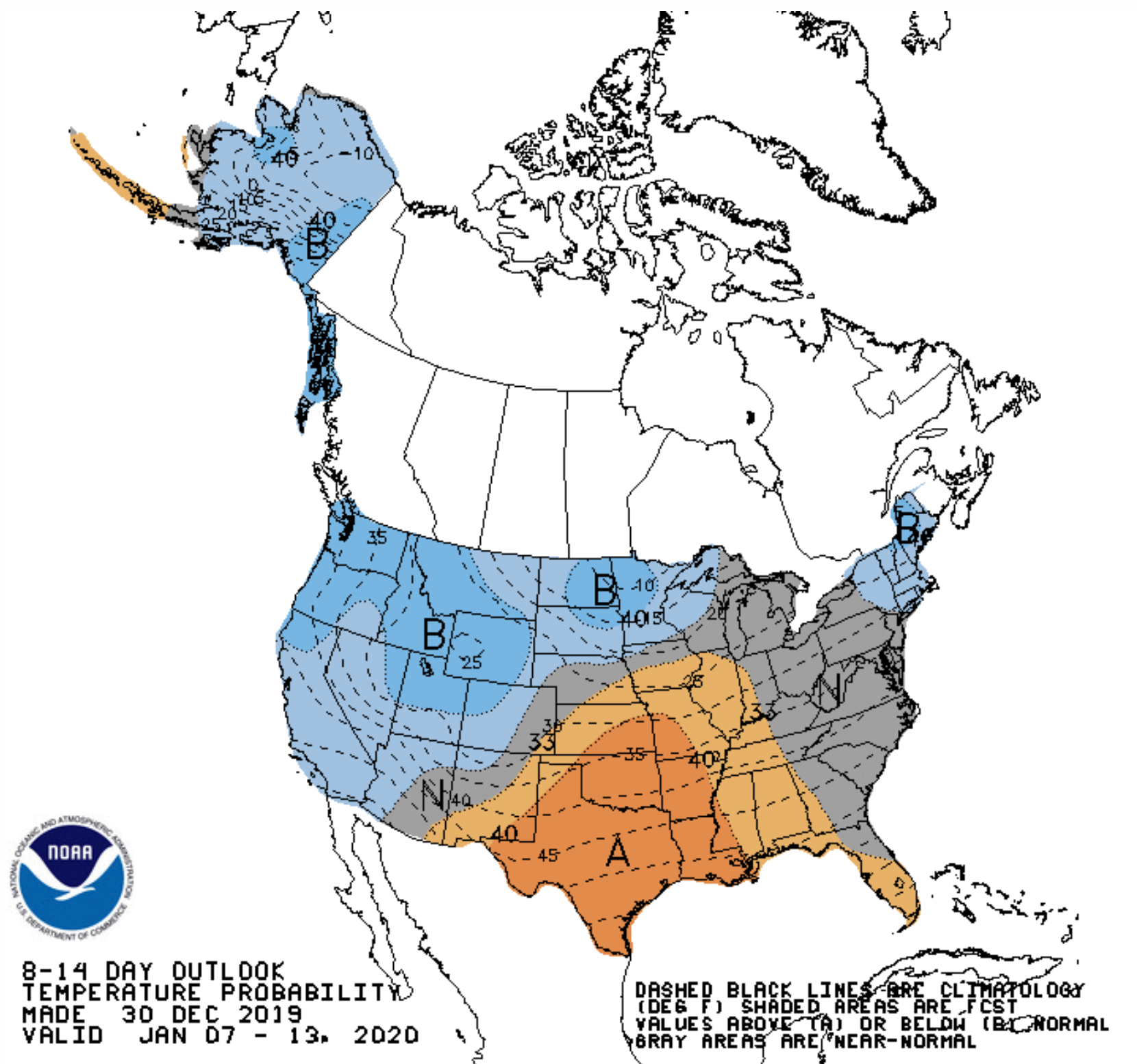
Arrival Delay by Season



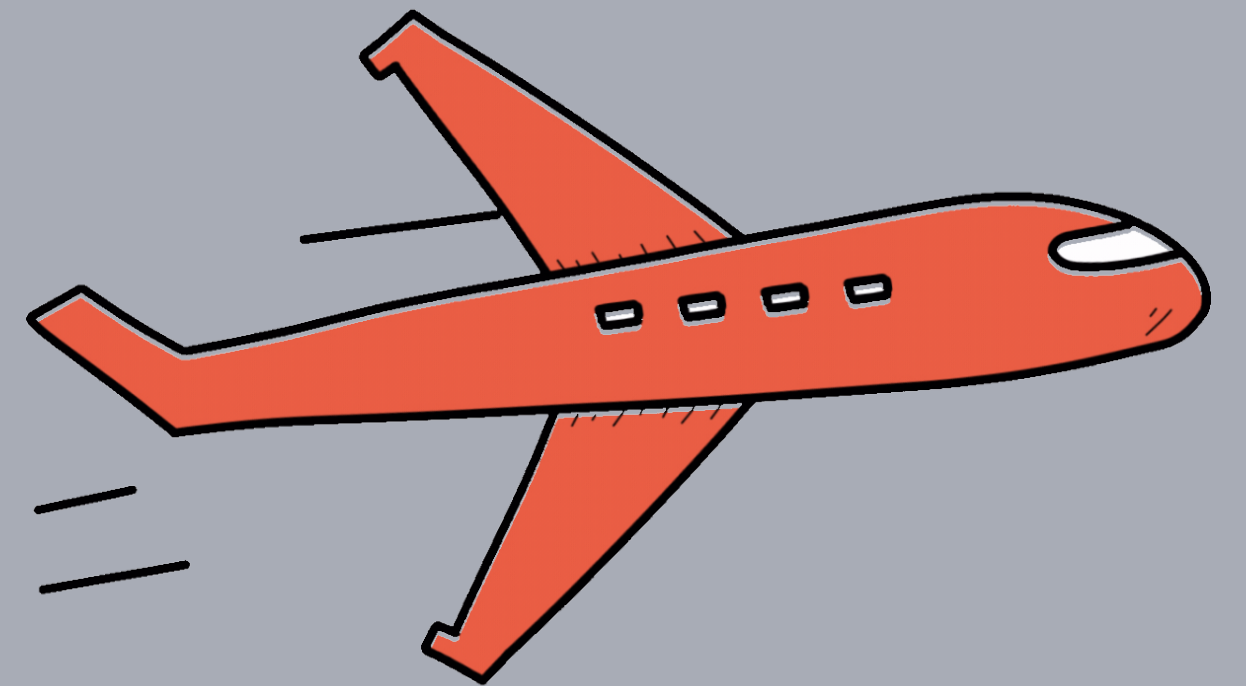


Arrival Delay by Mean Monthly Arrival Delay

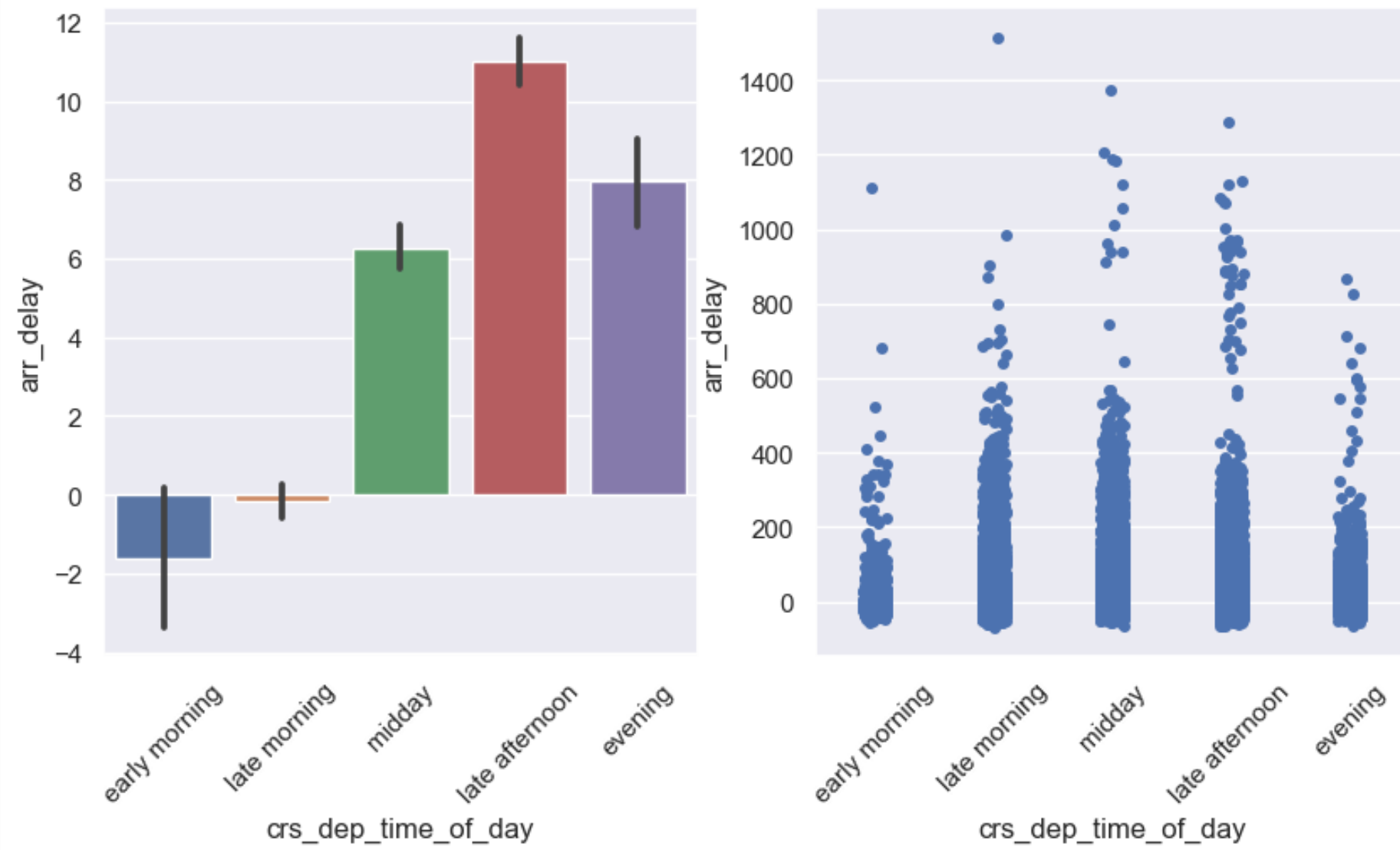




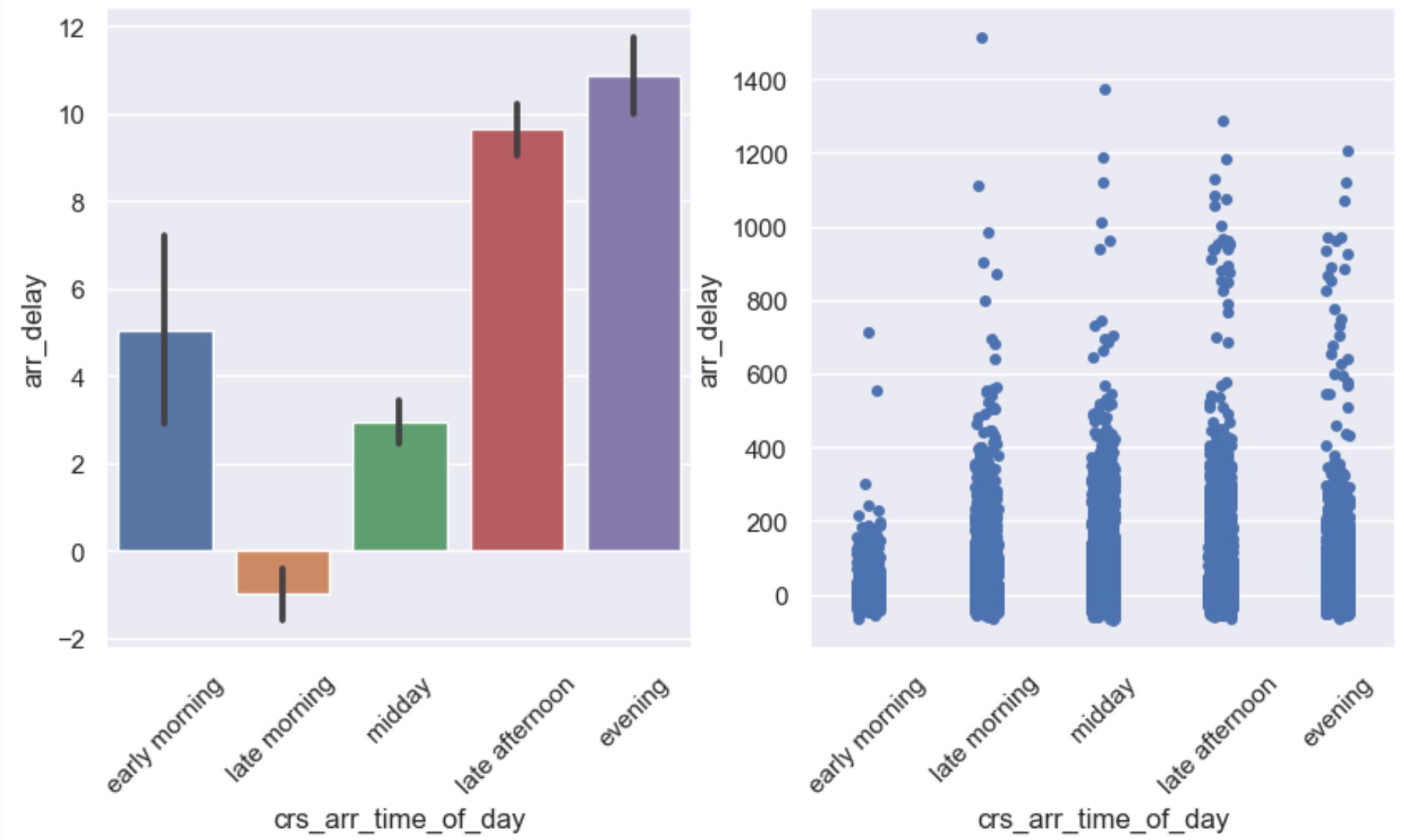
# Features based on Departure/Arrival Time



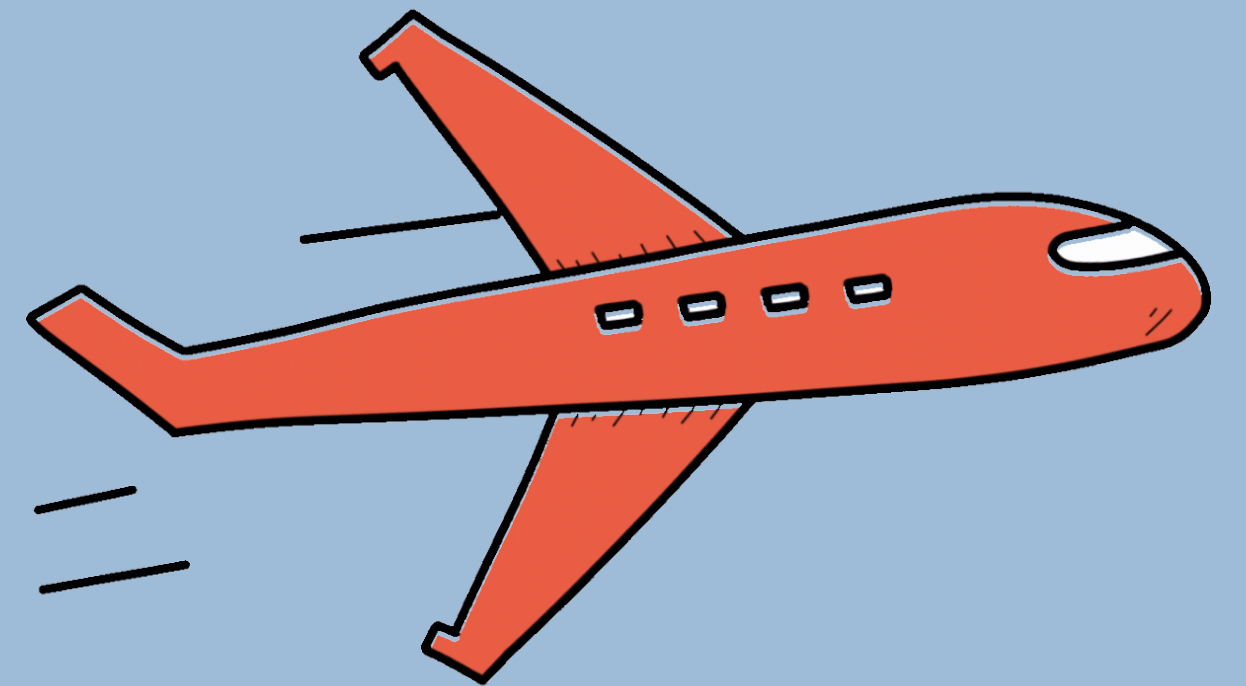
Arrival Delay by Departure Time of Day



Arrival Delay by Arrival Time of Day

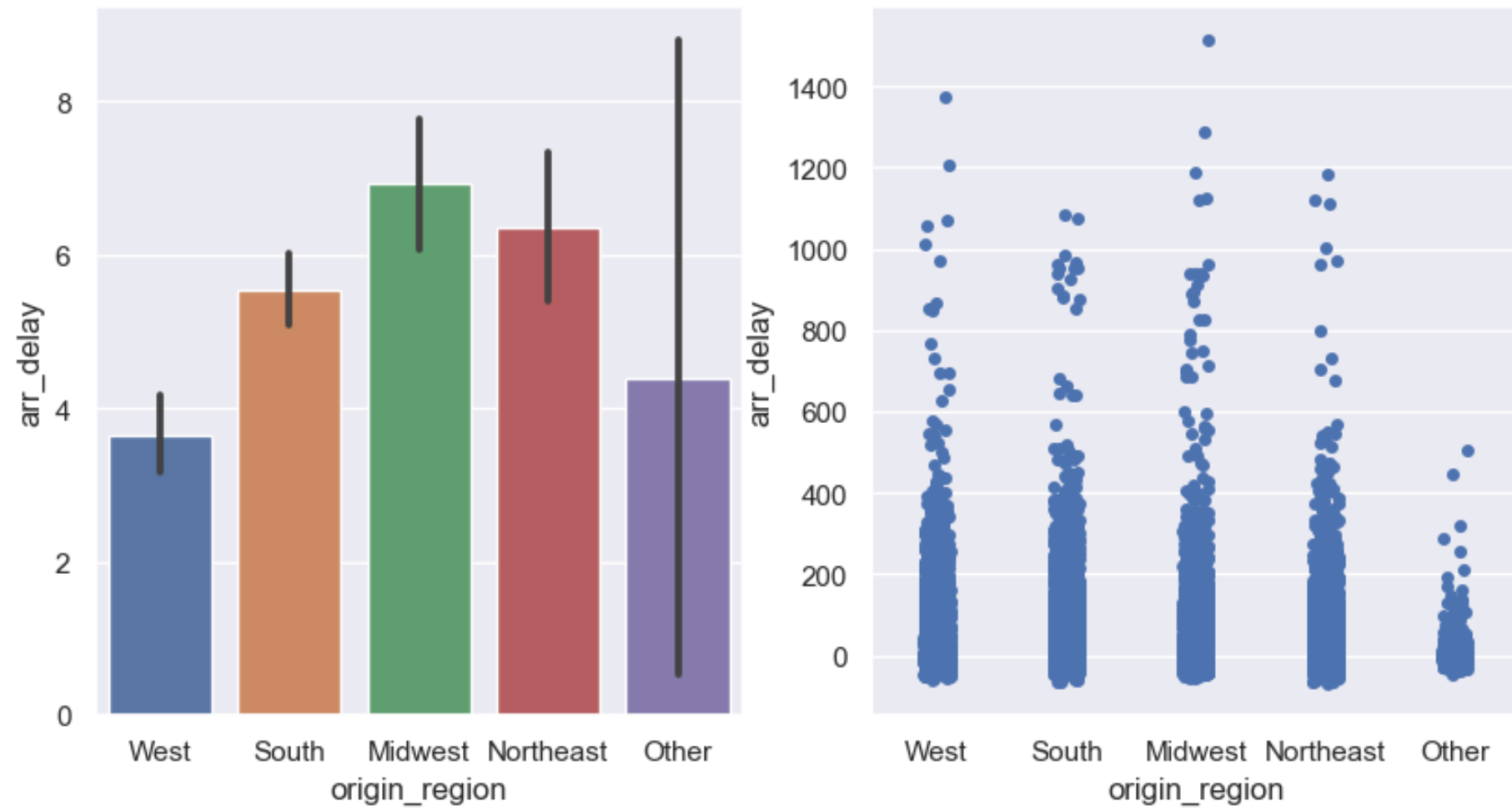


# Features based on Origin/Destination City

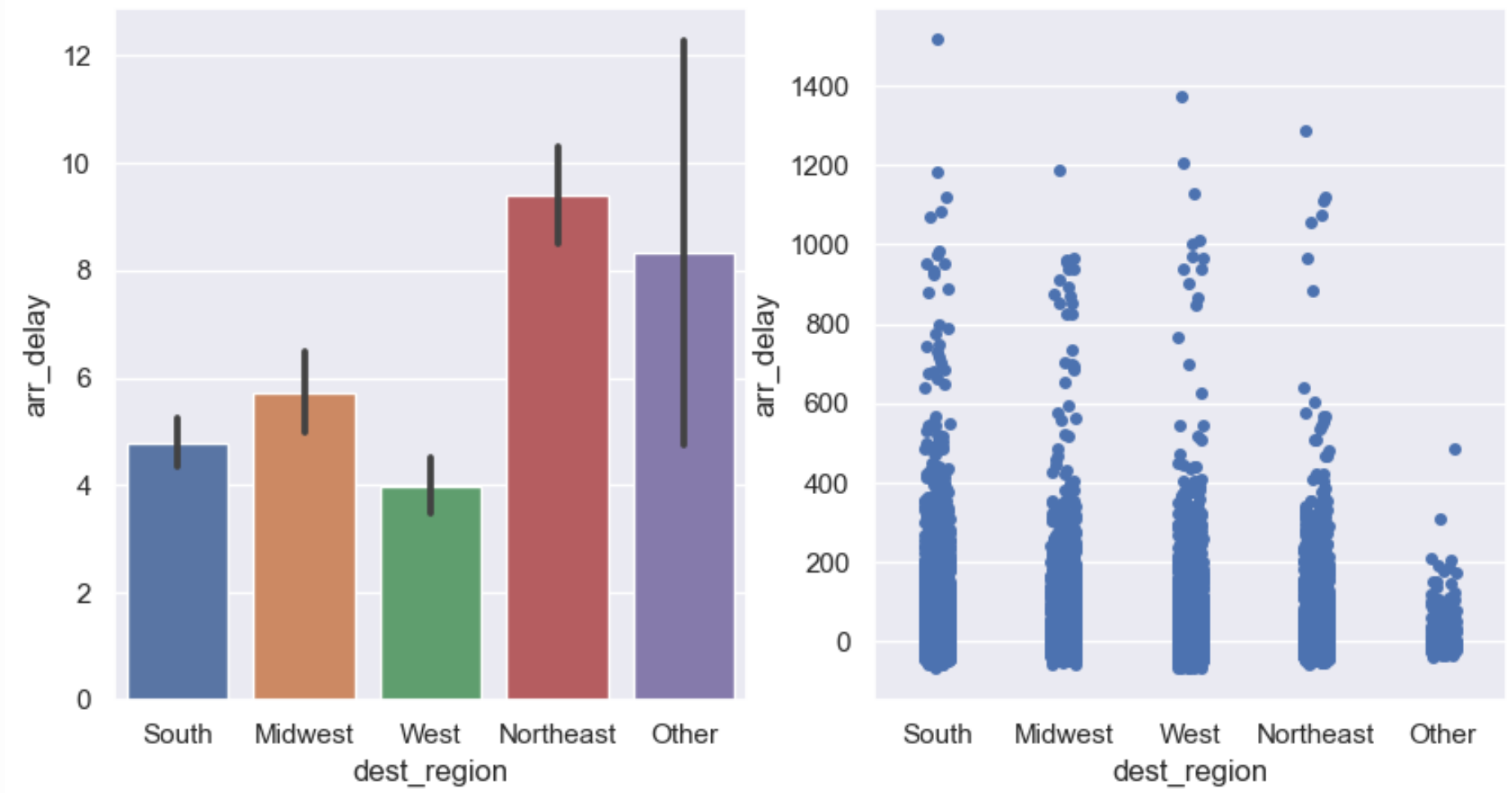




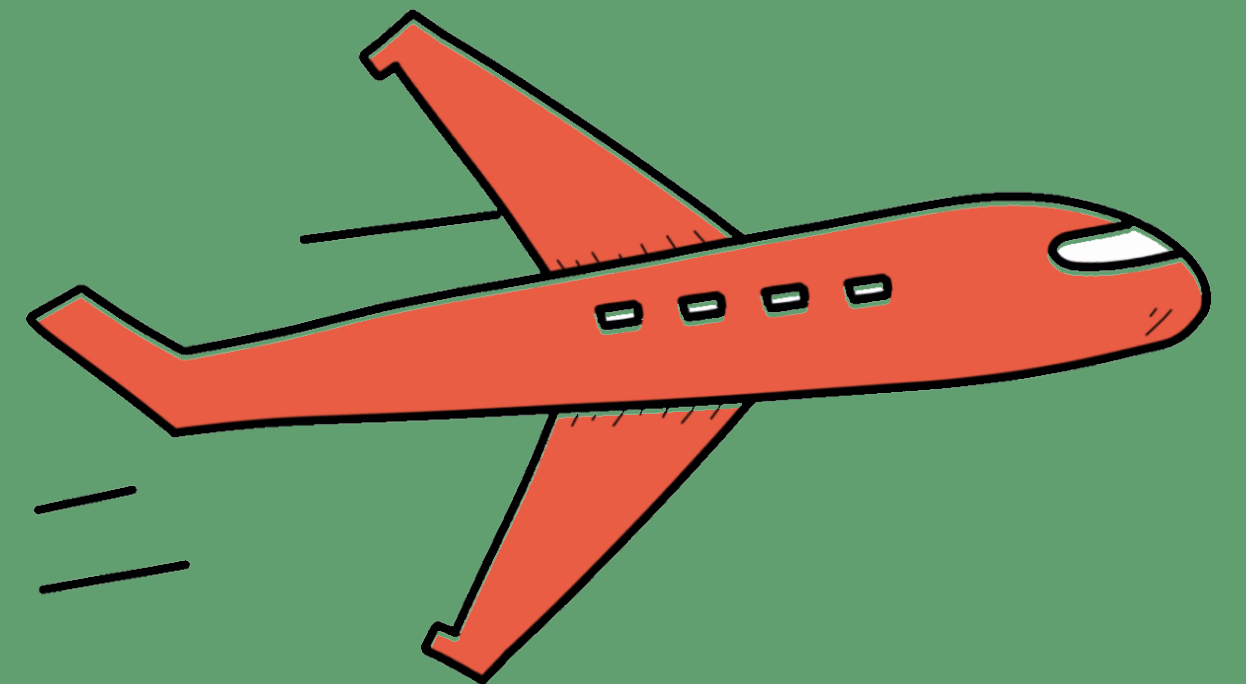
Arrival Delay by Flight Origin



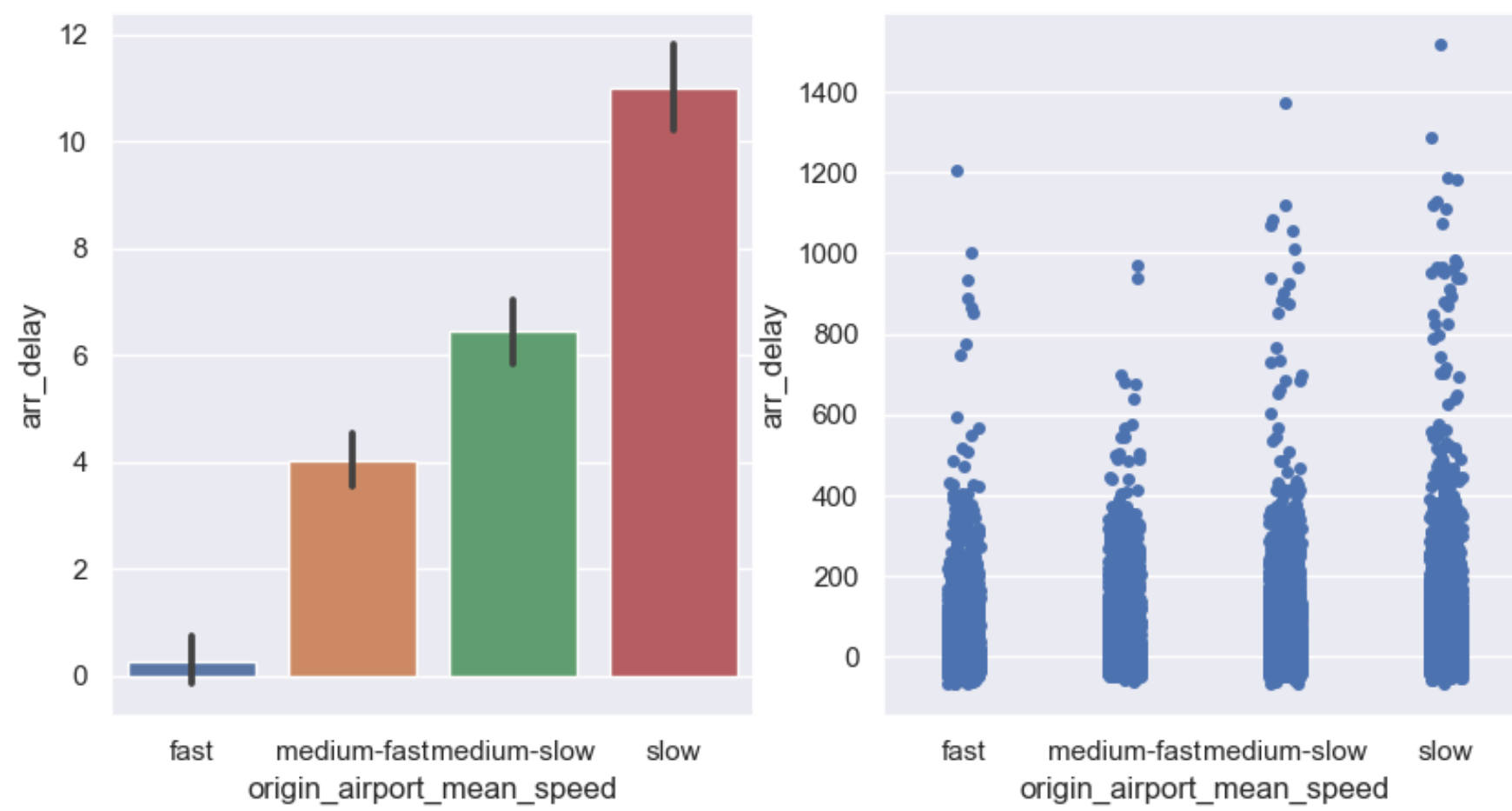
Arrival Delay by Destination Region



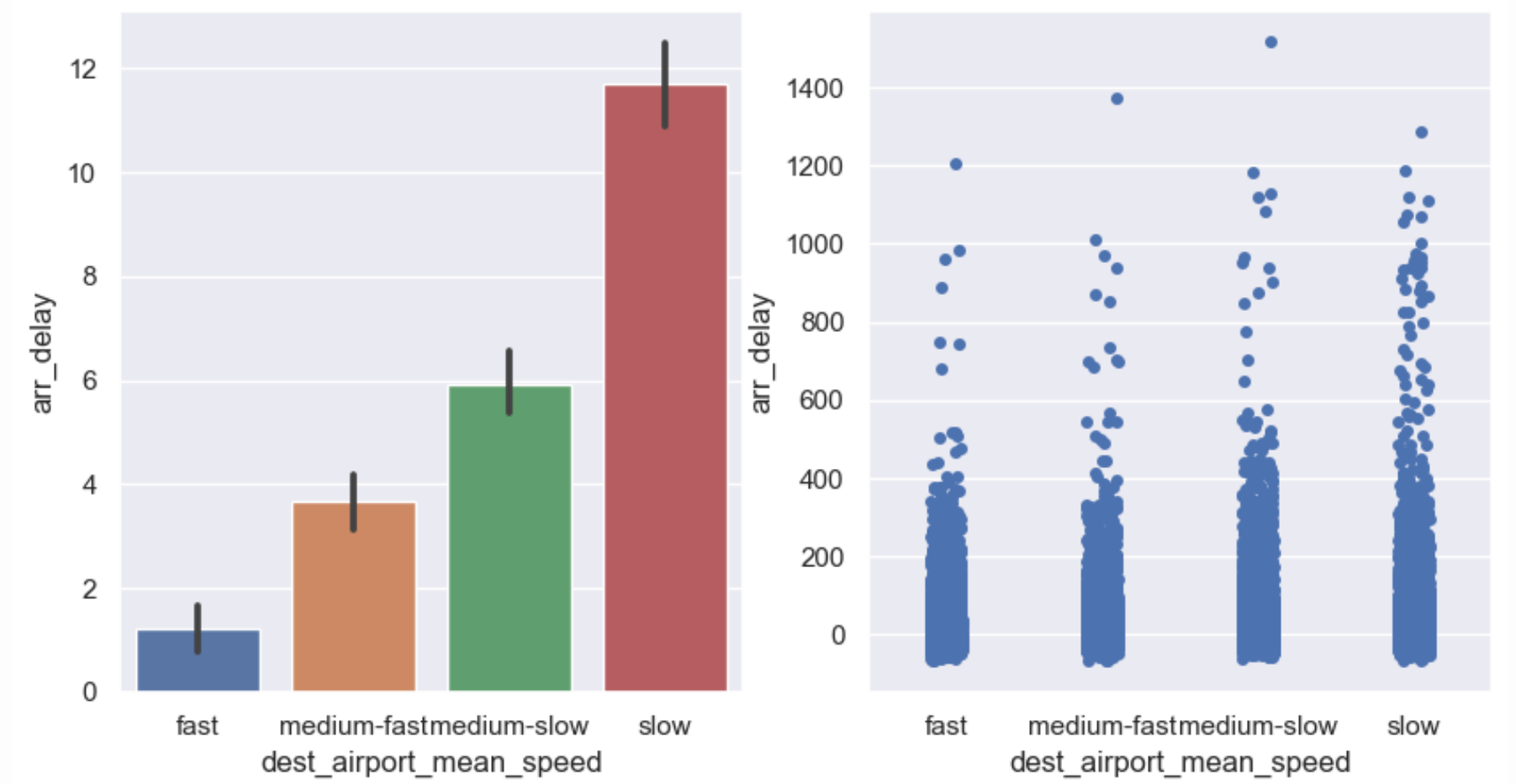
# Features based on Airport Id



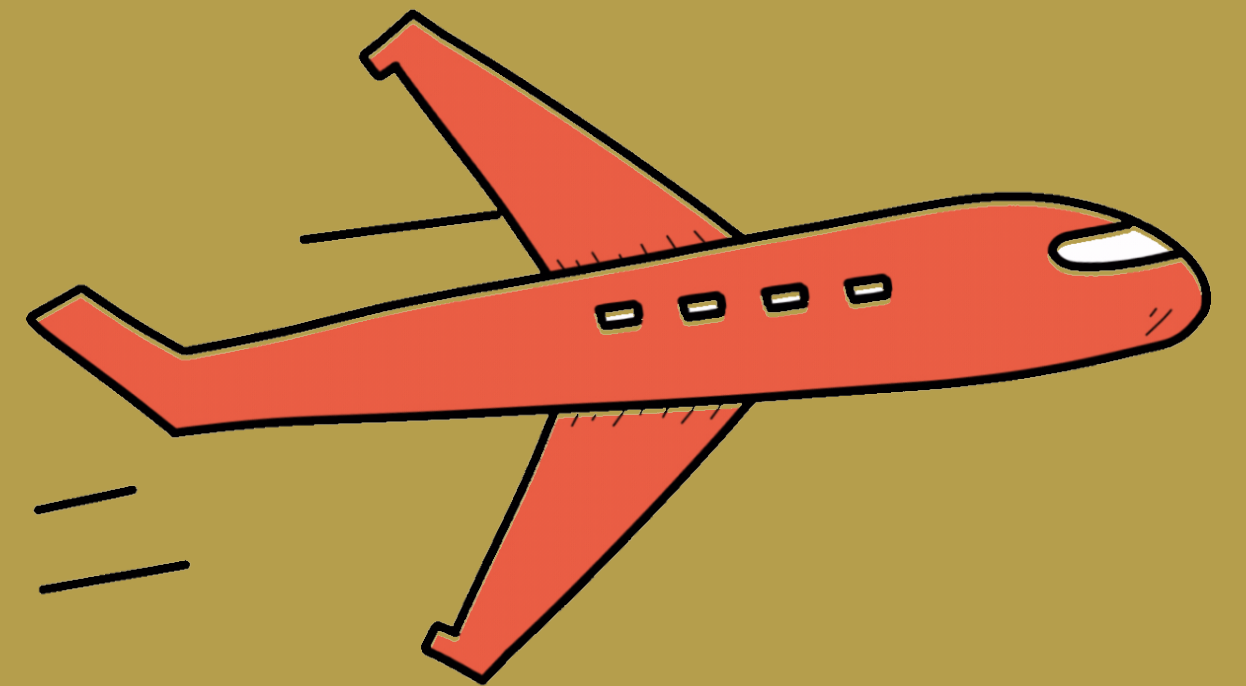
Arrival Delay by Origin Airport Speed (Mean)



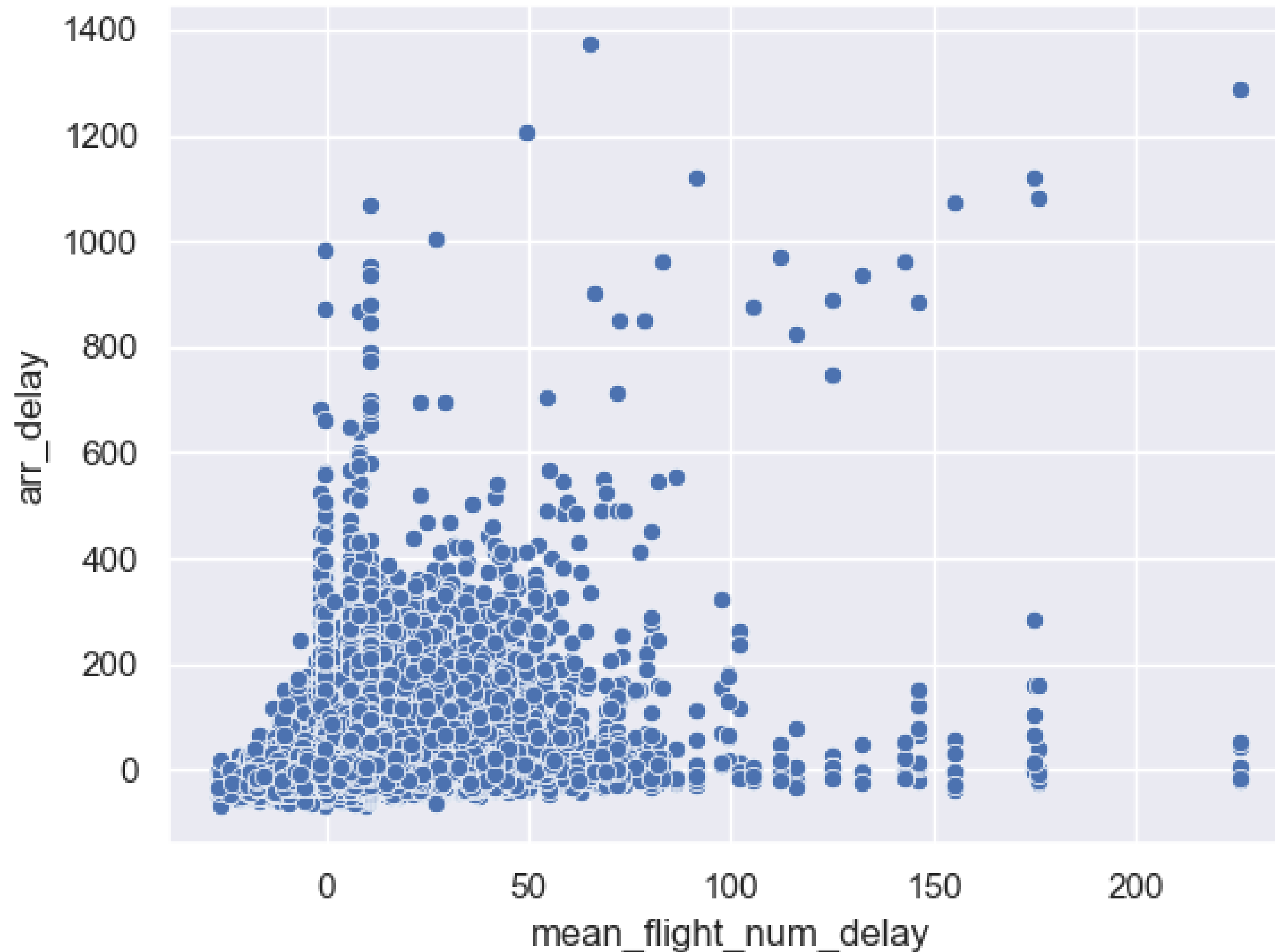
Arrival Delay by Destination Airport Speed (Mean)



# Features based on Flight Routes



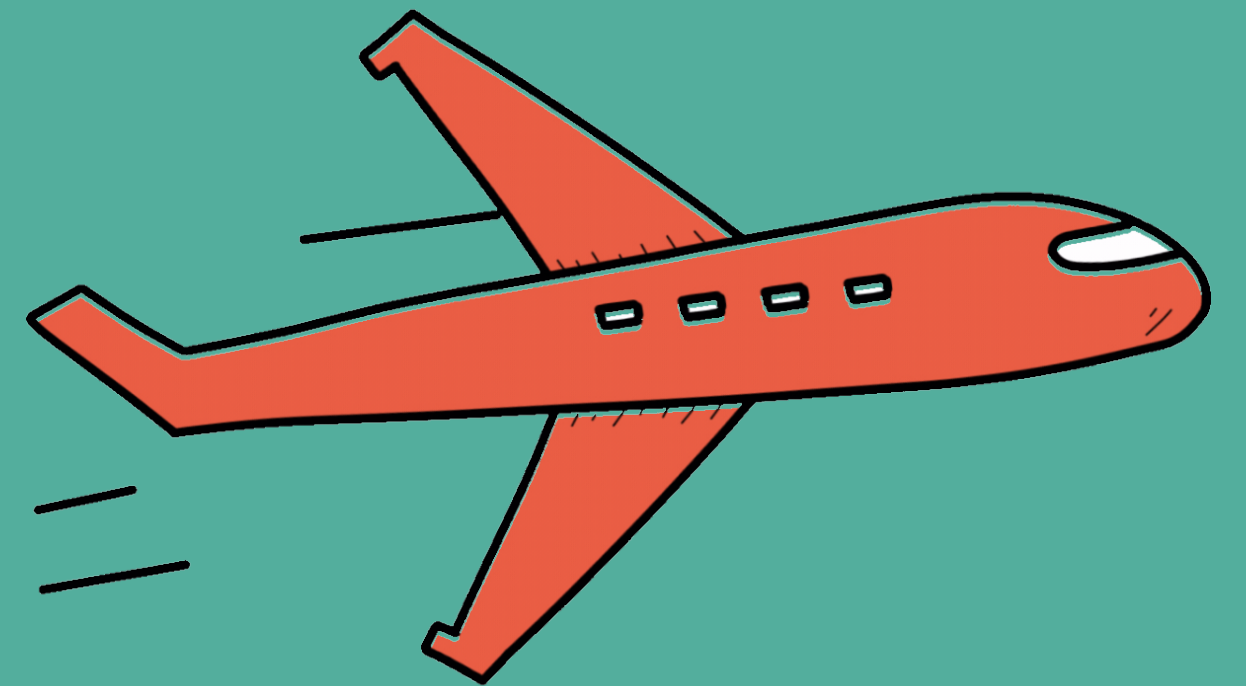
Arrival Delay by Mean Flight Number Delay



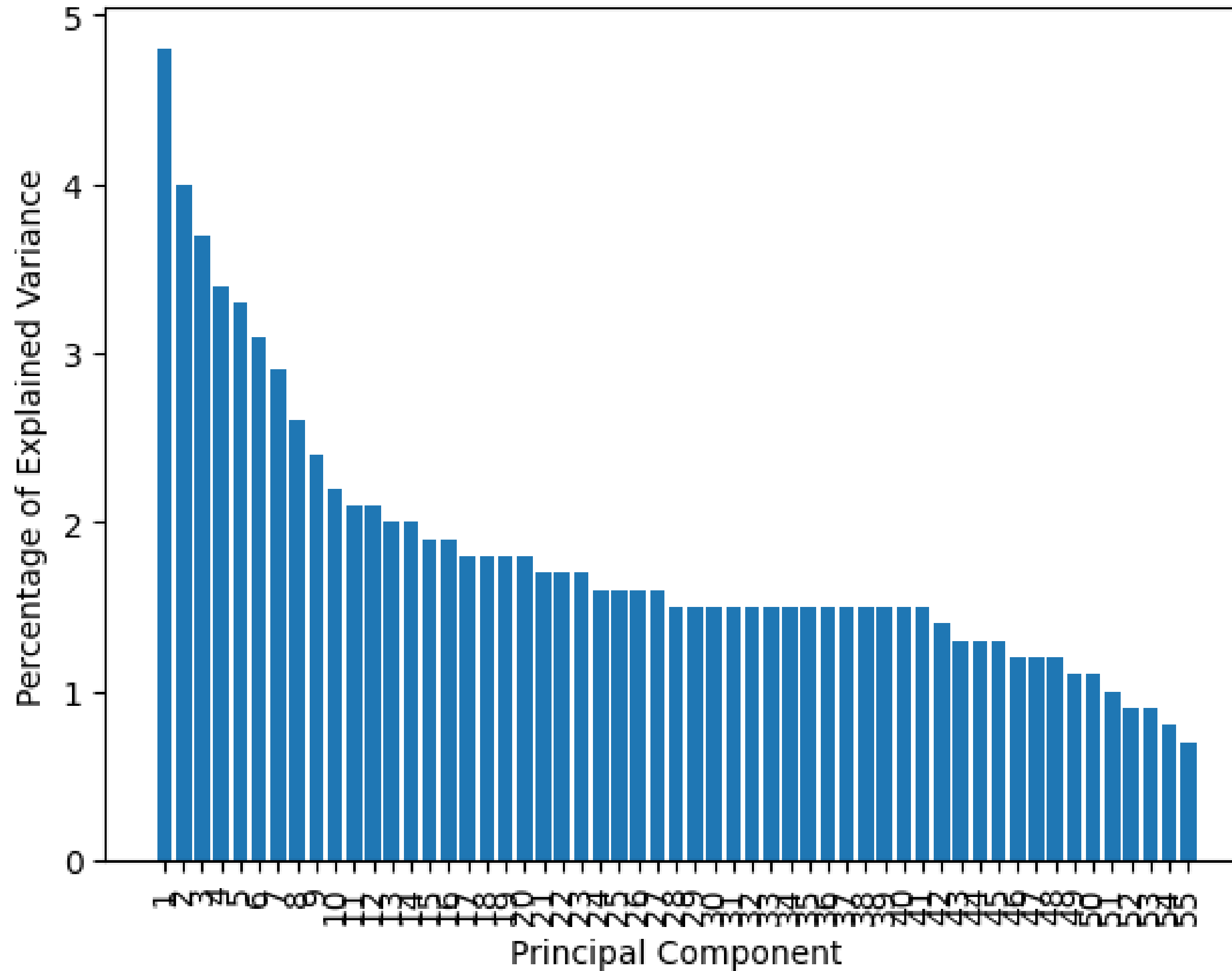
**Correlation  
Coefficient  
= 0.27**



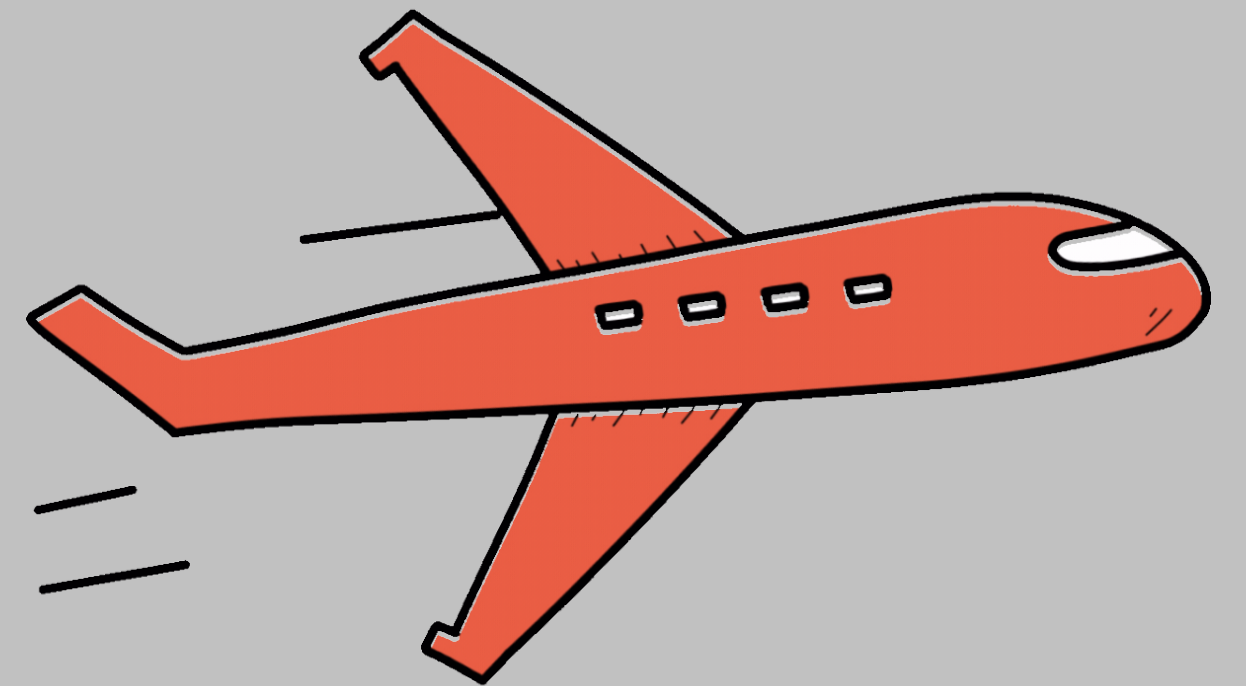
# Dimensionality Reduction



Scree Plot



# Training Models



|            | MAE_test  | MAE_train | RMSE_test | RMSE_train | R2_test  | R2_train | ADJR2_test | ADJR2_train |
|------------|-----------|-----------|-----------|------------|----------|----------|------------|-------------|
| linear     | 15.221647 | 15.392797 | 36.407460 | 34.151856  | 0.492709 | 0.485384 | 0.492709   | 0.485384    |
| ridge      | 15.221711 | 15.392862 | 36.407493 | 34.151856  | 0.492708 | 0.485384 | 0.492708   | 0.485384    |
| polynomial | 15.197644 | 15.355645 | 36.486101 | 34.140478  | 0.490516 | 0.485727 | 0.490516   | 0.485727    |
| lasso      | 15.370936 | 15.545292 | 36.493779 | 34.166493  | 0.490301 | 0.484943 | 0.490301   | 0.484943    |
| sgd        | 15.309494 | 15.750729 | 36.497936 | 34.245050  | 0.490185 | 0.482572 | 0.490185   | 0.482572    |
| r_forest   | 18.072724 | 18.122804 | 39.210367 | 35.901627  | 0.411593 | 0.431301 | 0.411593   | 0.431301    |
| g_boost    | 15.522921 | 15.364055 | 39.487800 | 33.617242  | 0.403237 | 0.501370 | 0.403237   | 0.501370    |
| voting_r   | 15.691872 | 15.270631 | 39.558562 | 33.401120  | 0.401096 | 0.507760 | 0.401096   | 0.507760    |
| xgb        | 15.983414 | 15.201328 | 43.680043 | 33.057213  | 0.269799 | 0.517845 | 0.269799   | 0.517845    |
| d_tree     | 16.588934 | 15.123663 | 45.439398 | 32.918182  | 0.209792 | 0.521892 | 0.209792   | 0.521892    |

# What's Next



## Check Tailnumber against the FAA's Registry

Examine how different types of airplanes impact delays

## Integrating the CPC's Hurricane Season Outlook

Examine the impact of extreme weather conditions

## Fine Tuning our Hyper-Parameters



# Lessons Learned



Don't automatically drop outliers - only if they clearly don't make sense

Multiple dimensions presenting the same information in different ways weakens model performance

Split into train / test set before feature engineering

Thank you!

