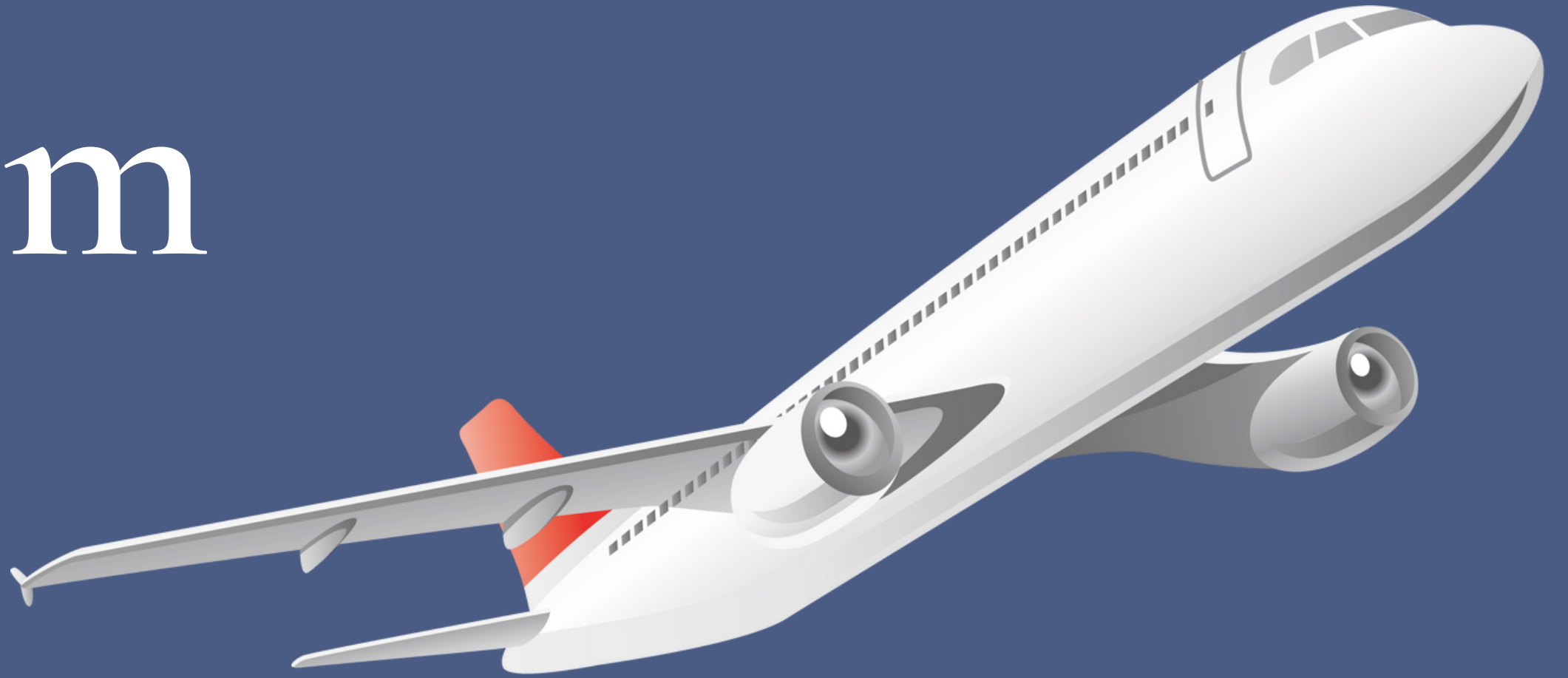


Mid-Term Project



SHILPI & SHUBHANGI



DELTA 2106	B10	11:05am	On Time
DELTA 4547	D01	11:15am	On Time
DELTA 780	B12	1:30pm	Boarding
DELTA 4649	C03	11:05am	On Time
DELTA 5296	E83	3:00pm	Boarding
DELTA 6729	E83	2:00pm	On Time
DELTA 7383	D09	11:00am	Boarding
DELTA 7383	E70	11:10am	On Time
DELTA 7383	B7	11:09am	On Time

Objective

PREDICTING DELAYS

WHY ???

- Challenging, lot of learning
- Can check accuracy
- Supervised ML Algorithms

Steps



Approach

RESEARCH

- Research what factors affect flight delay
- Deciding how to sample dataset for modeling

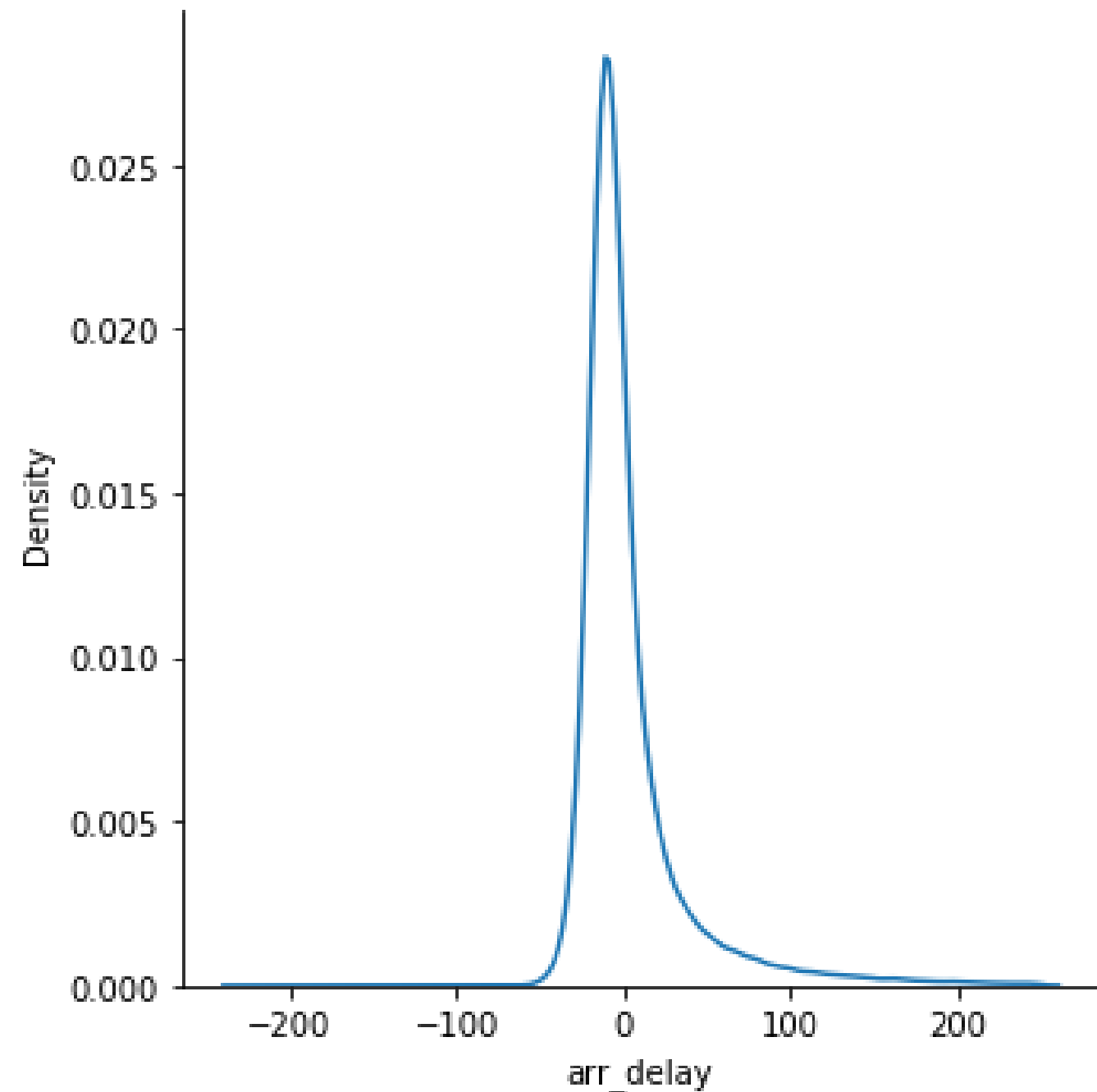
TOOLS & TECHNIQUES

- Jupyter Lab/Notebook
- Google Collab
- Pycaret
- Google Drive/Slack

CODE QUALITY

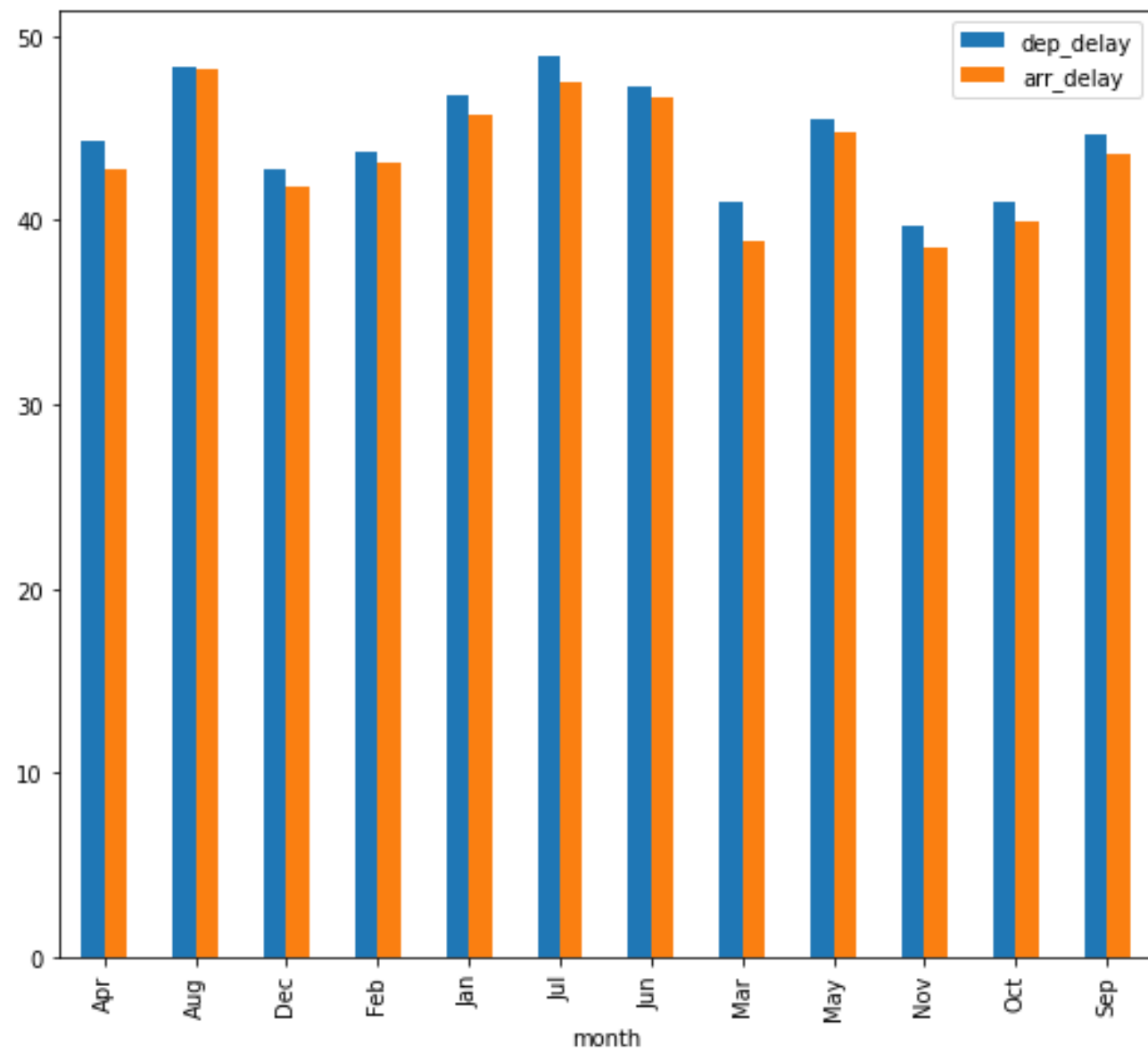
- Organize Codes
- Add comments & docstring
- Relevant names of variable
- Organized Dataset
- Storey-Telling

Exploratory Data Analysis

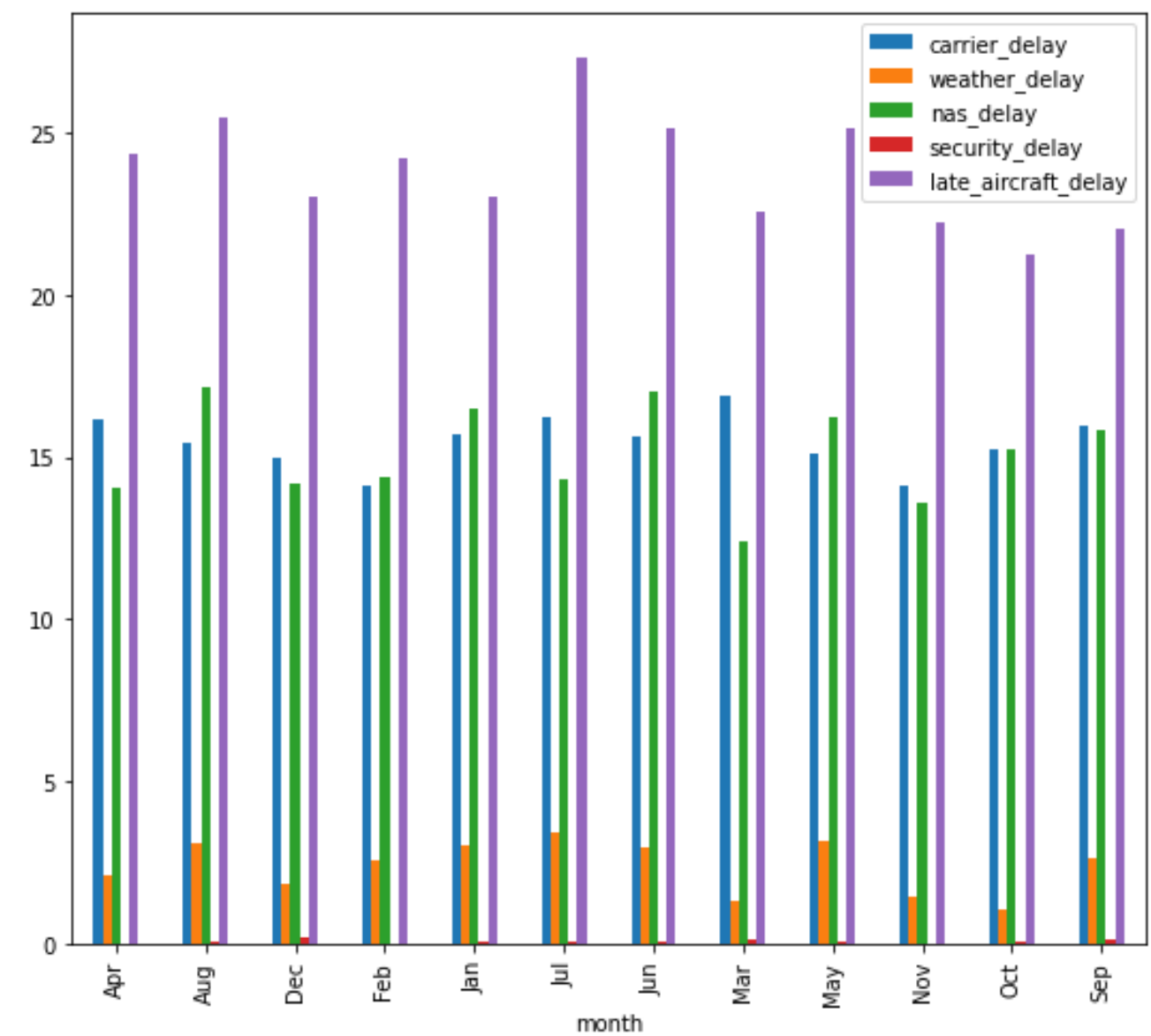


TEST HYPOTHESIS FOR NORMAL DISTRIBUTION

- Sampled Dataset Randomly of sample size 150K
- Plotted graph to check distribution
- Shepiro test on sample data < 5000

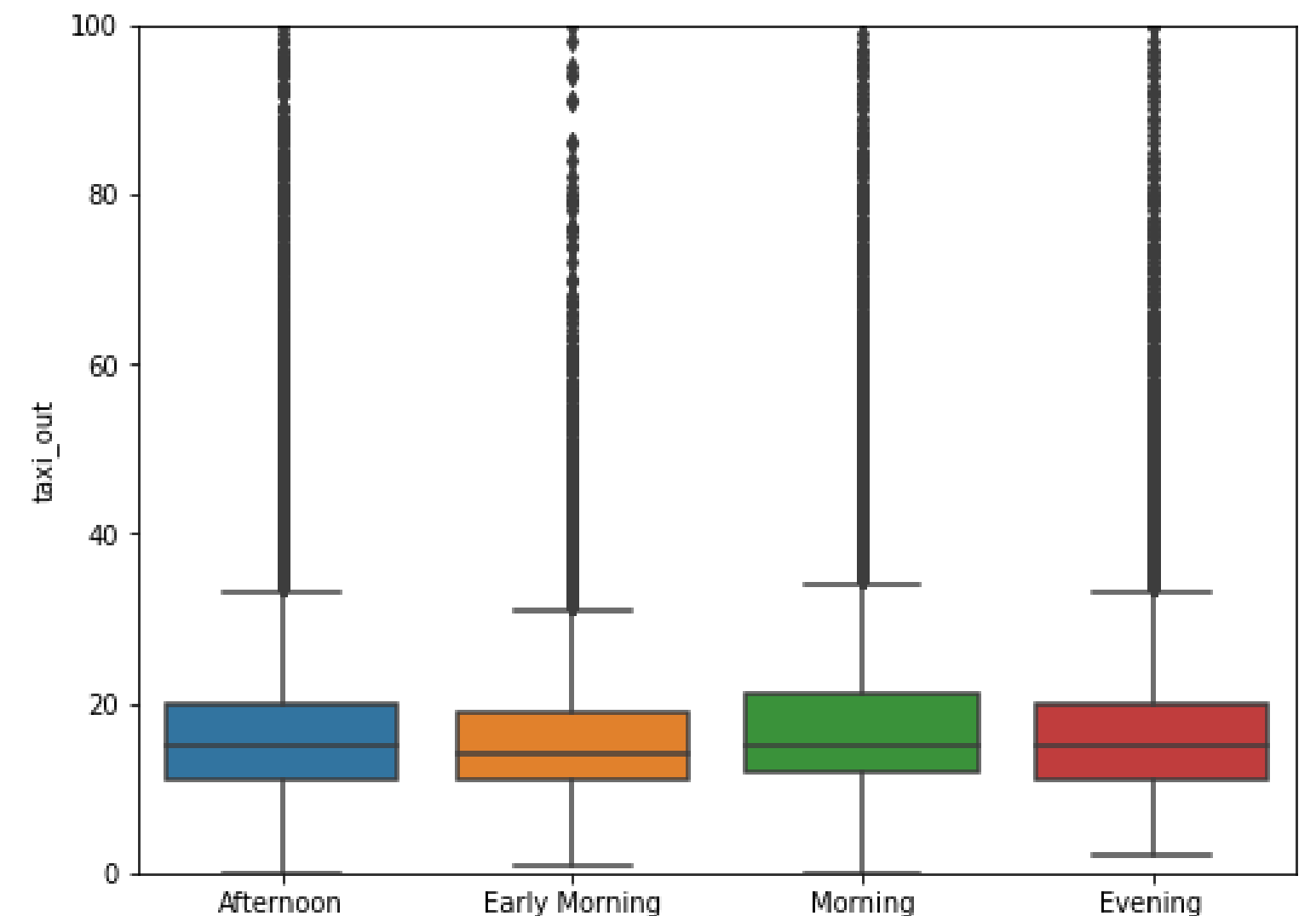
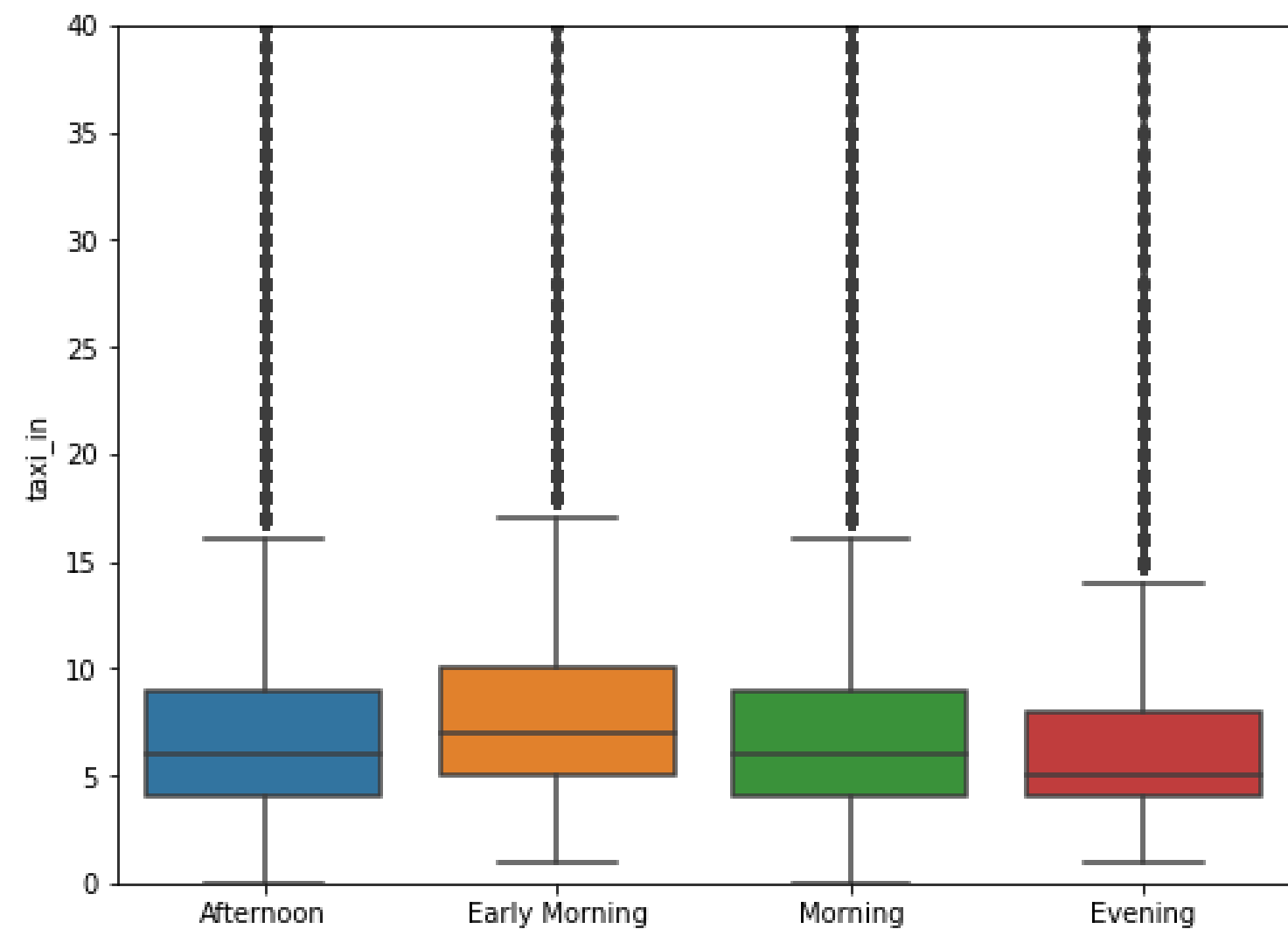


Monthly distribution of arrival and departure delays

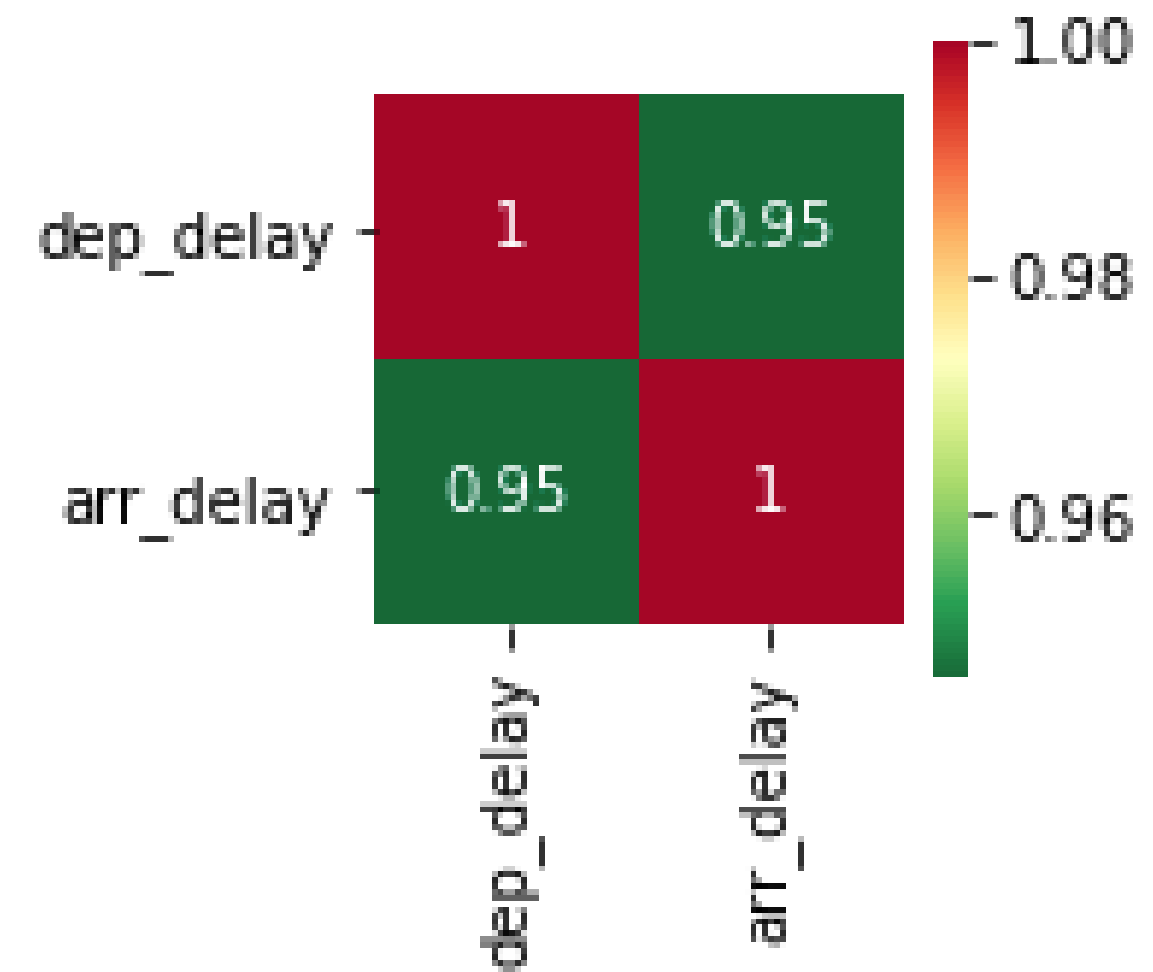
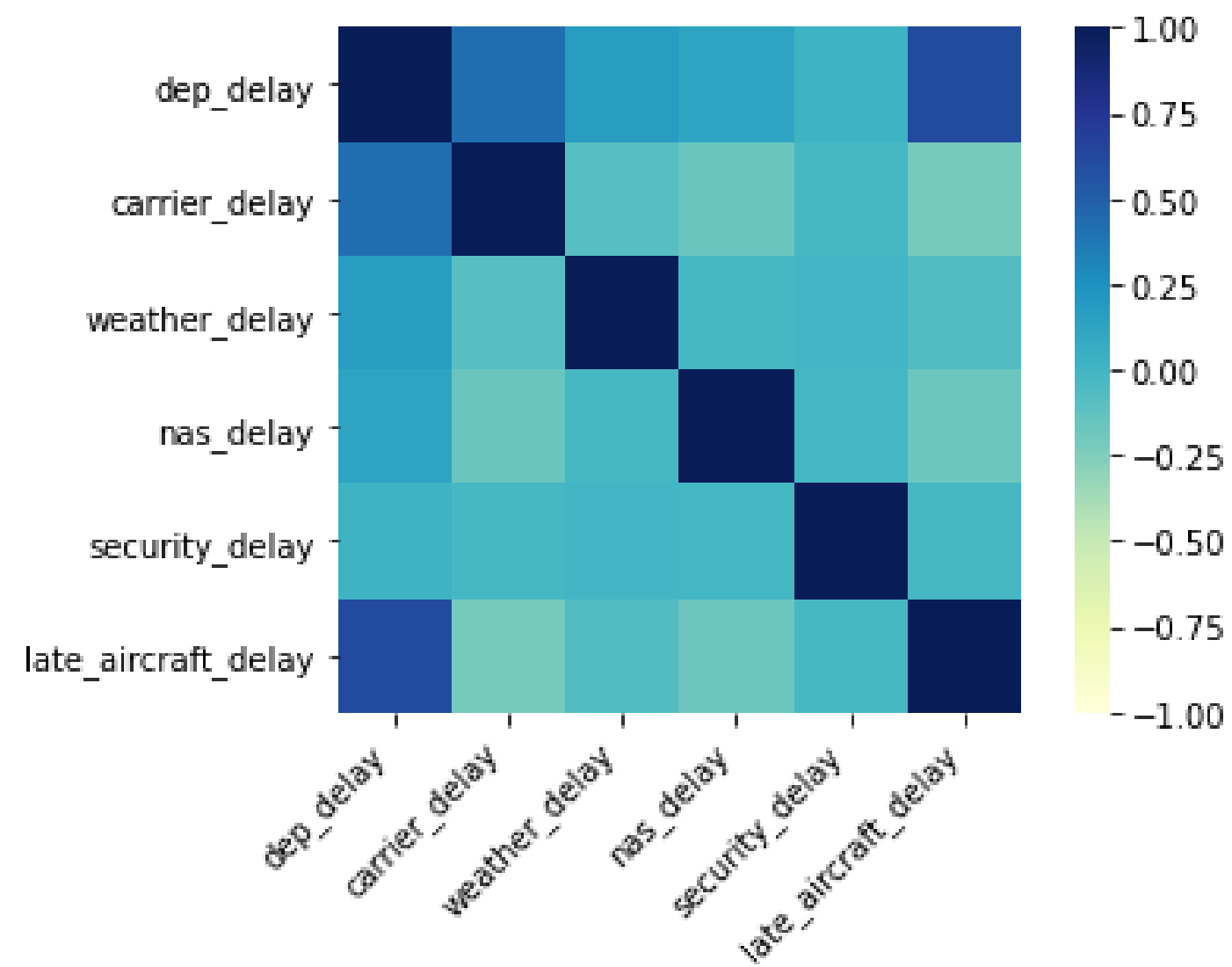


Monthly distribution of reasons of arrival delays

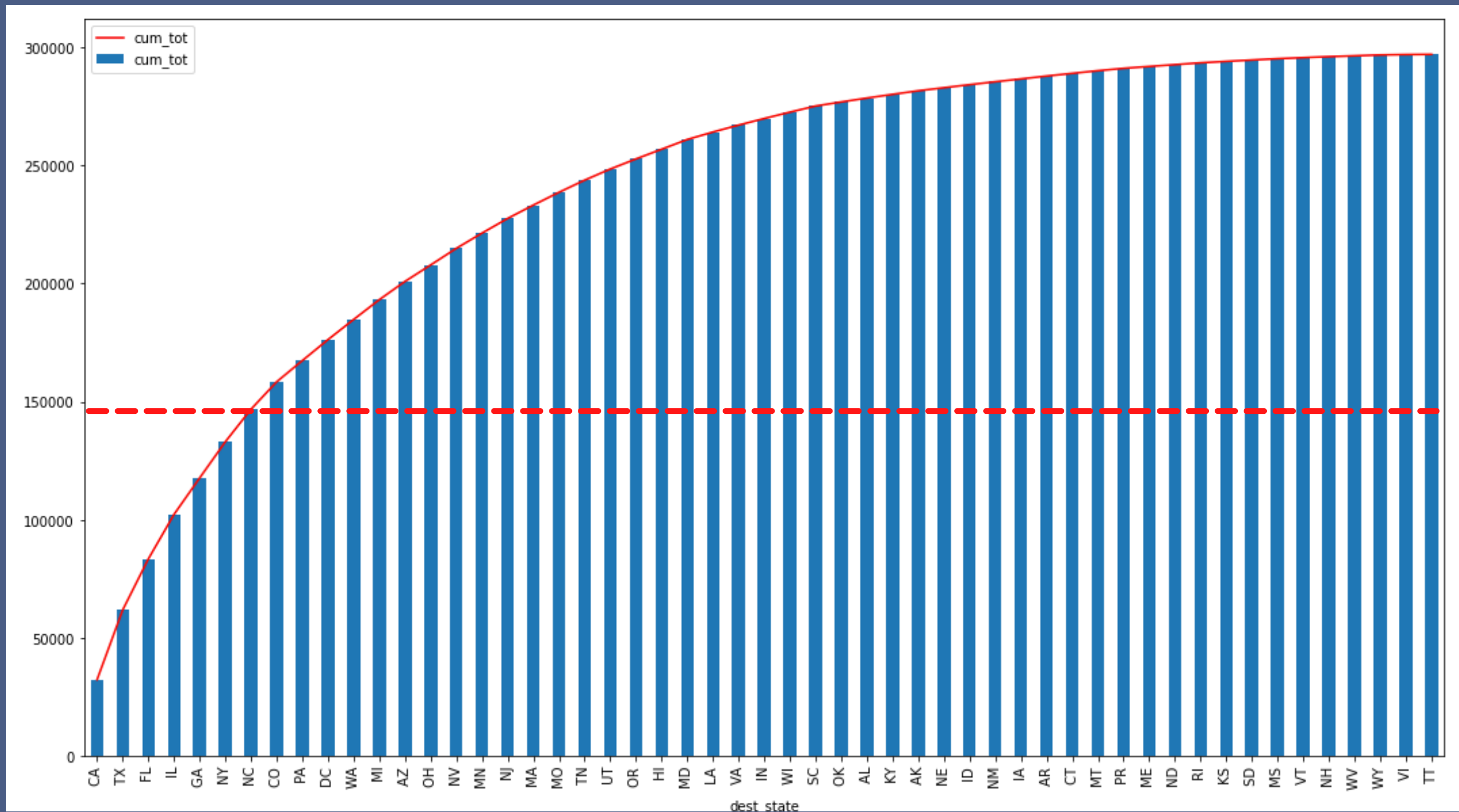
Taxi in and out during the day



Relationship between arrival and departure

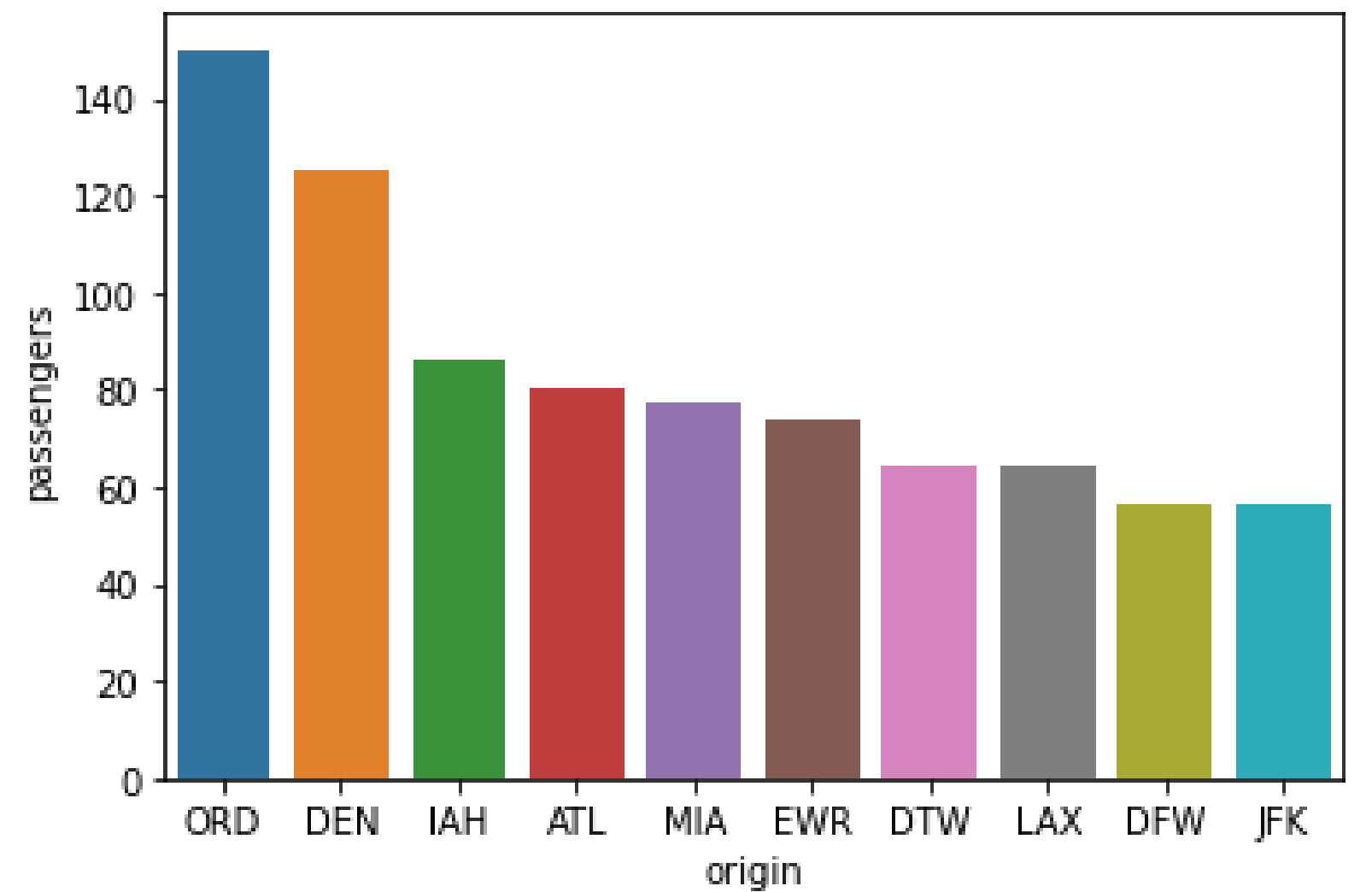
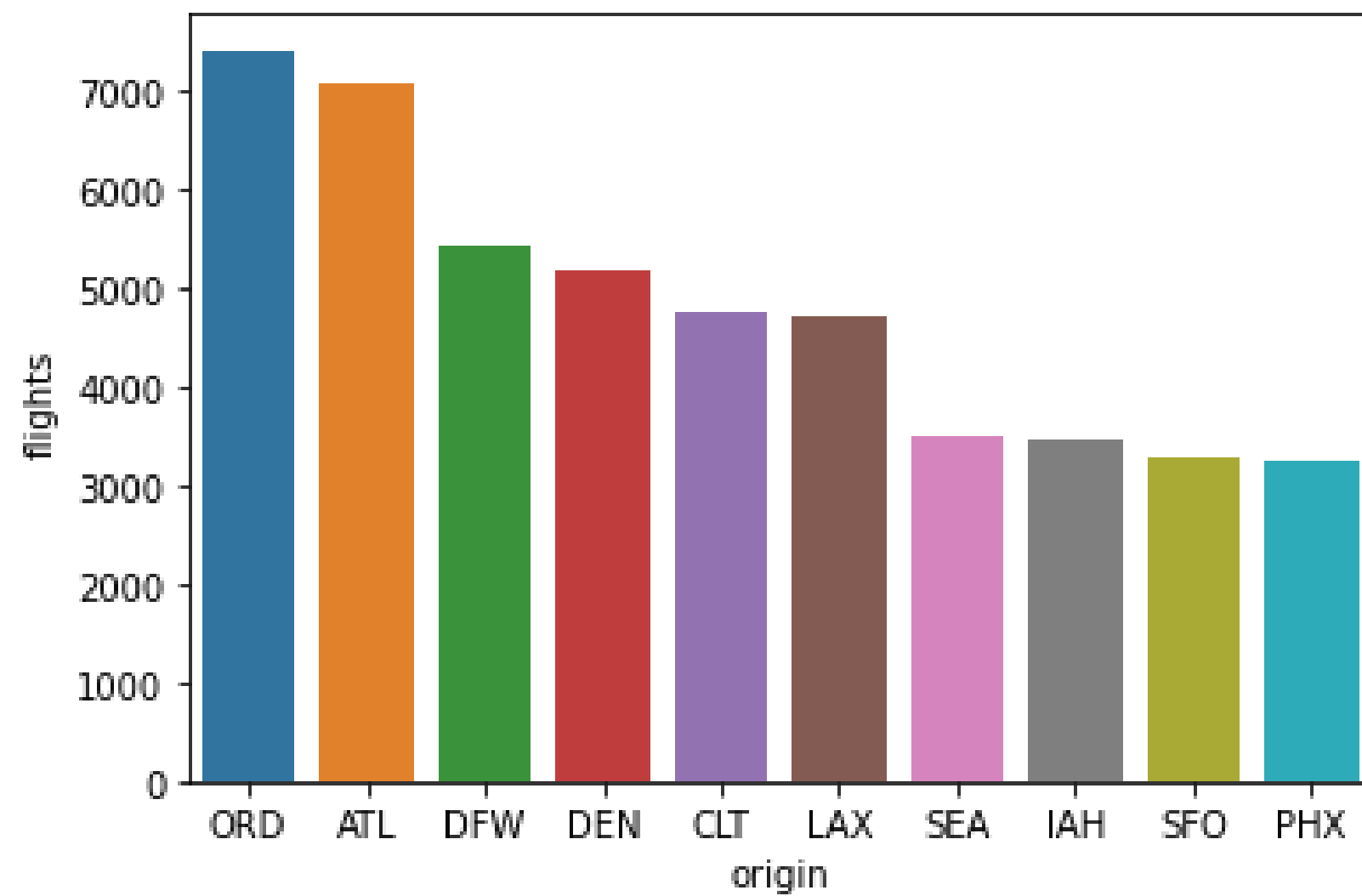


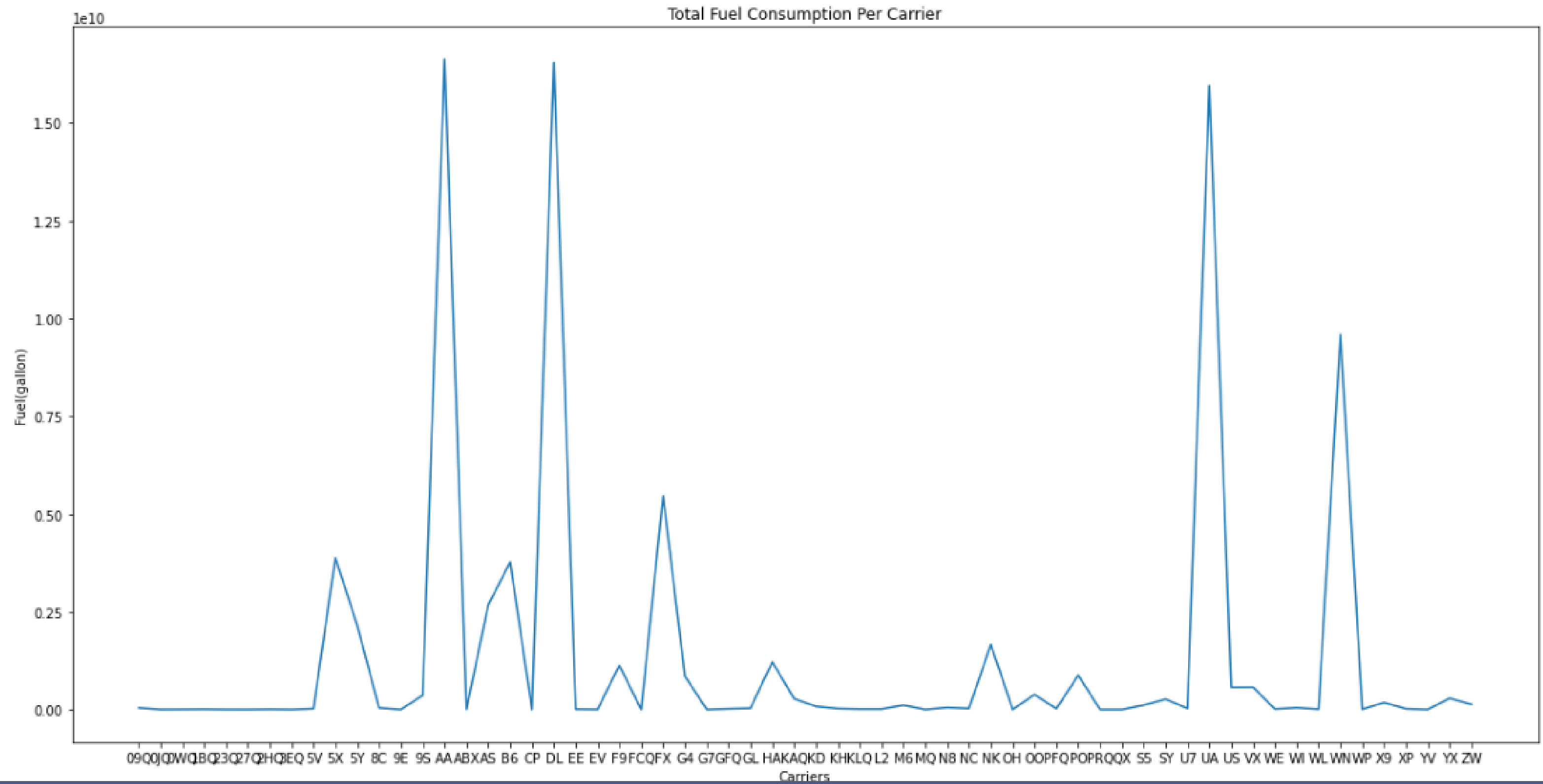
Air traffic on the basis of states



**States covering
50% air traffic**

- California
- Texas
- Florida
- Illinois
- Georgia
- New York
- North Carolina





ML Algorithms

Logistic Regression

Random Forest

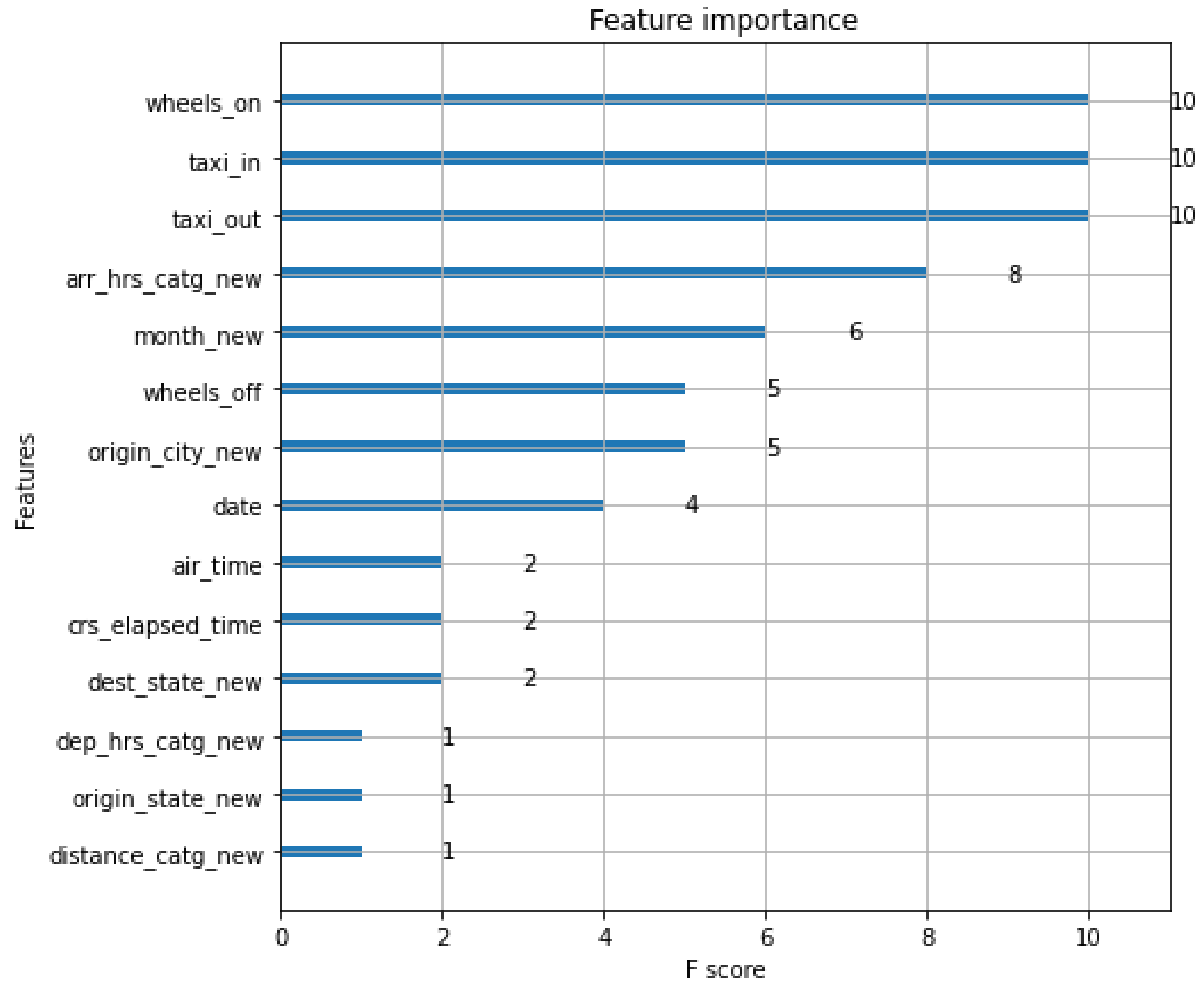
XG Boost



Model Evaluation

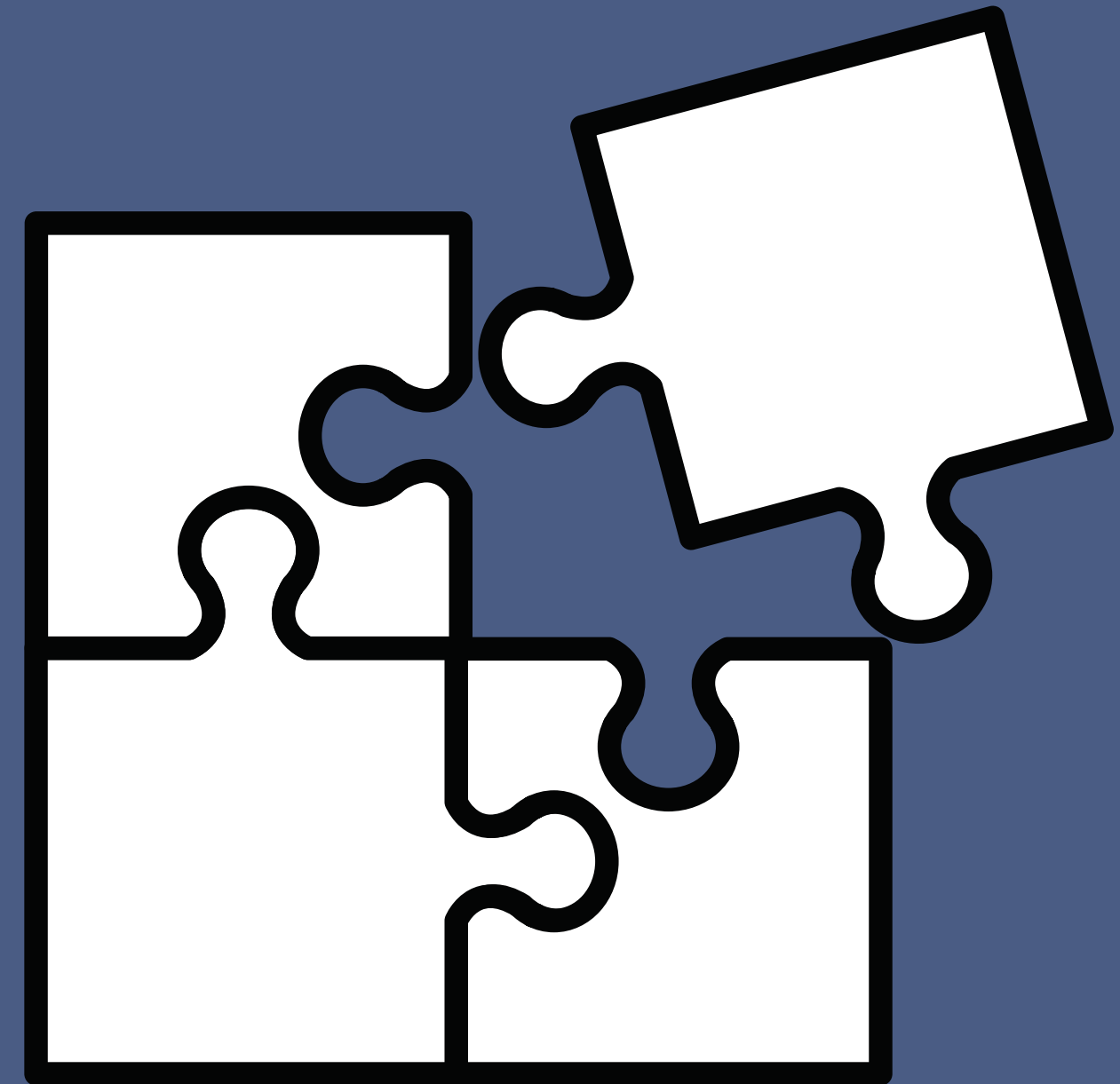
Random Forest: Accuracy 5%

XG Boost: MSE 33.81, MSE 29



Challenges

- Technical issues
- Sampling dataset
- Many new concepts to learn
- Tools - Google Colab, Pycaret
- No definite answer
- Limited Timeline



THANKYOU