



FLIGHT DELAY PREDICTION

MACHINE LEARNING

PILAR ESCRIG
VALERIA VALEVASHNIKOVA
ROCÍO ROMO
IVÁN SIMÓN

INTRO

- Explore and analyse USA flights over a 12 months period
- Predict whether a flight will be delayed/cancelled based on different metrics.

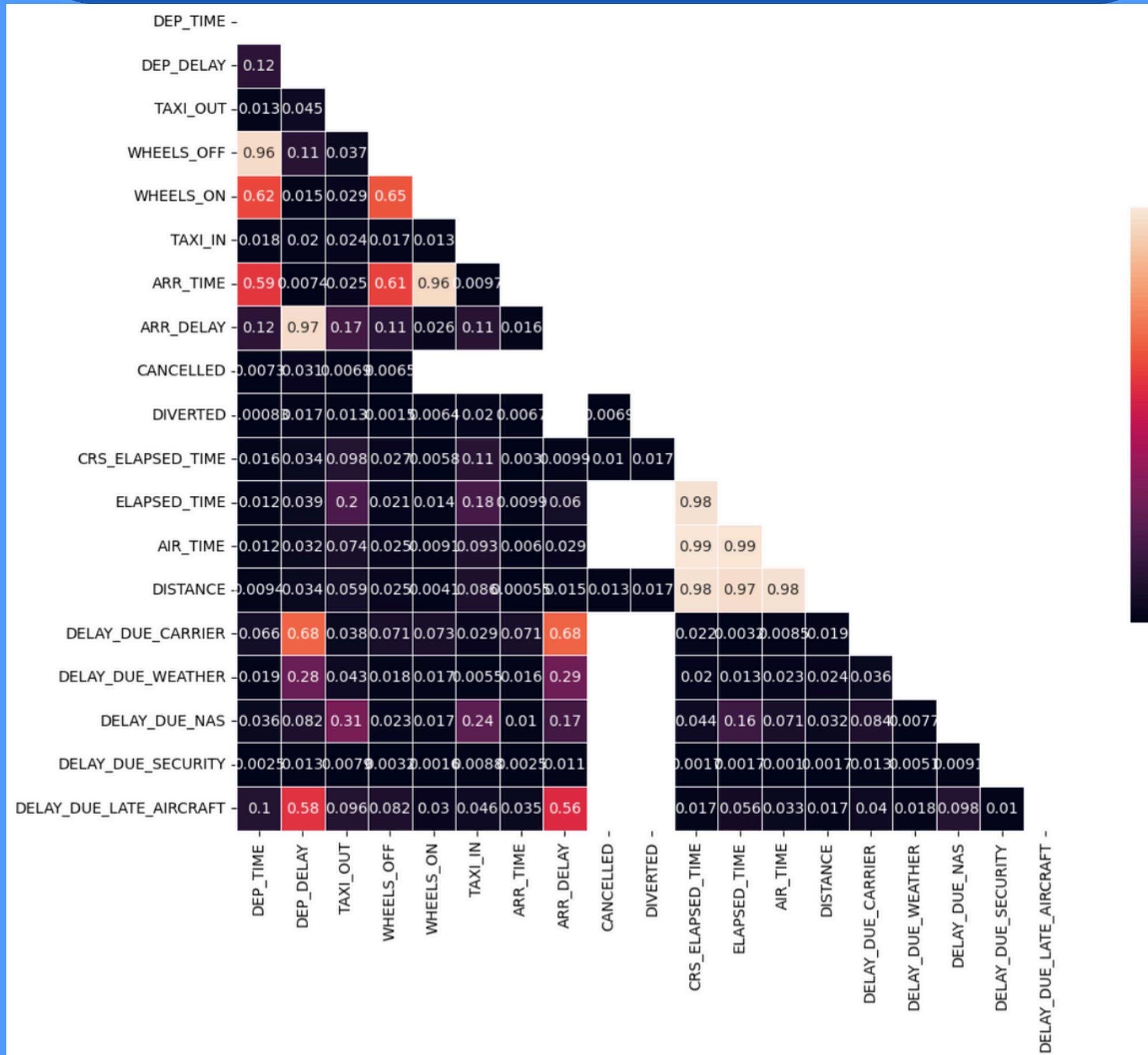
DATASET SELECTION



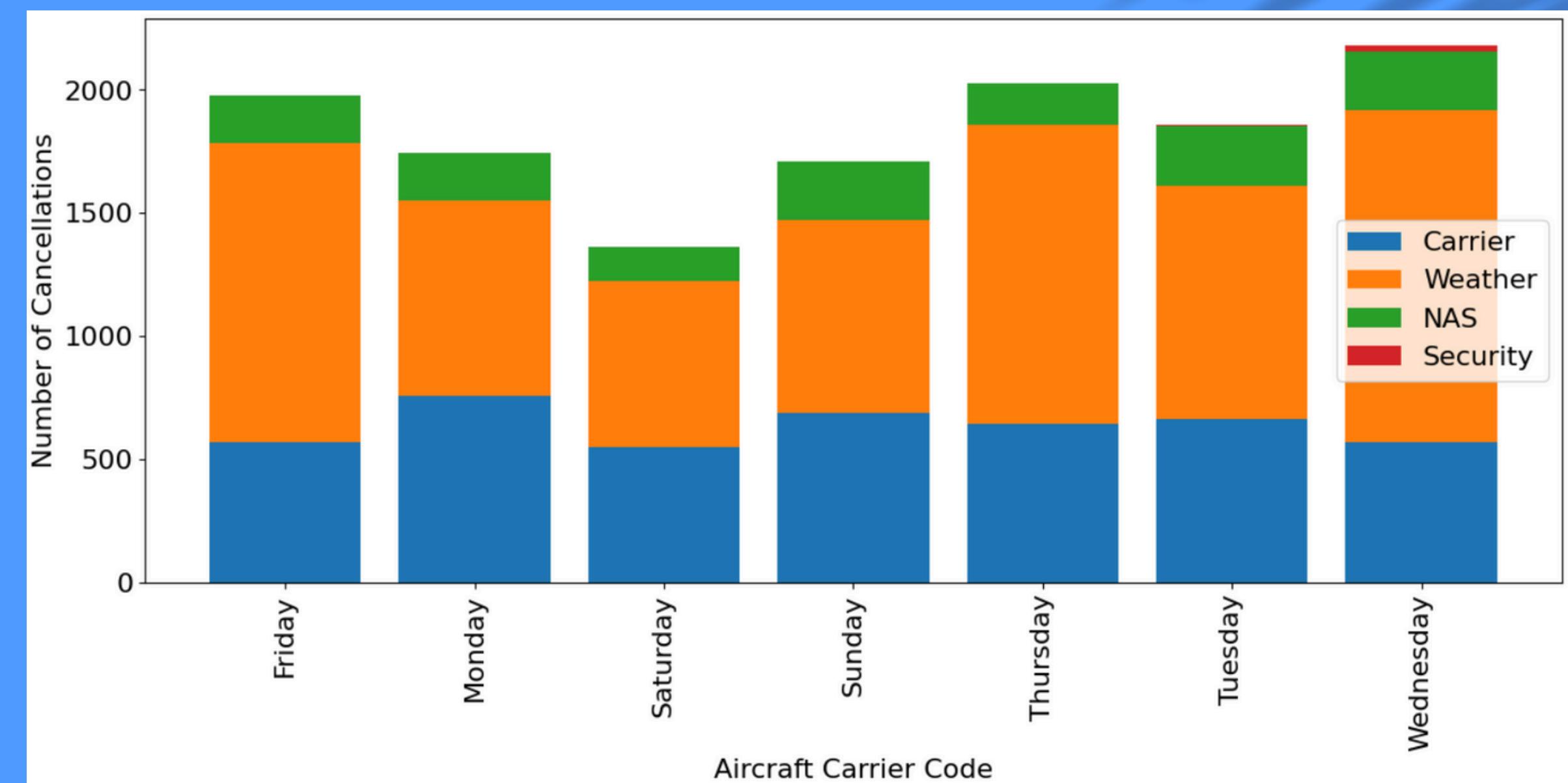
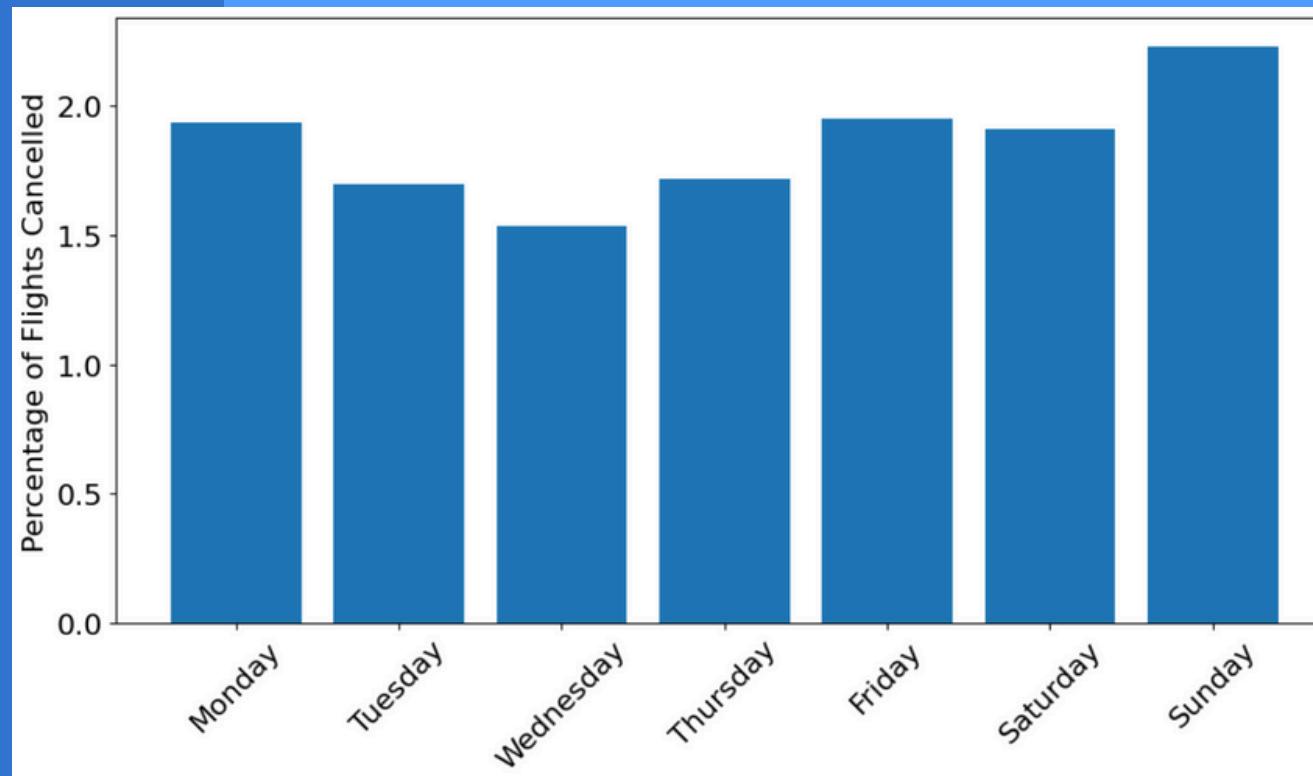
- Kaggle
- Flight Delay and Cancellation Dataset (2019-2023)

SEPTEMBER 2022 - AUGUST 2023

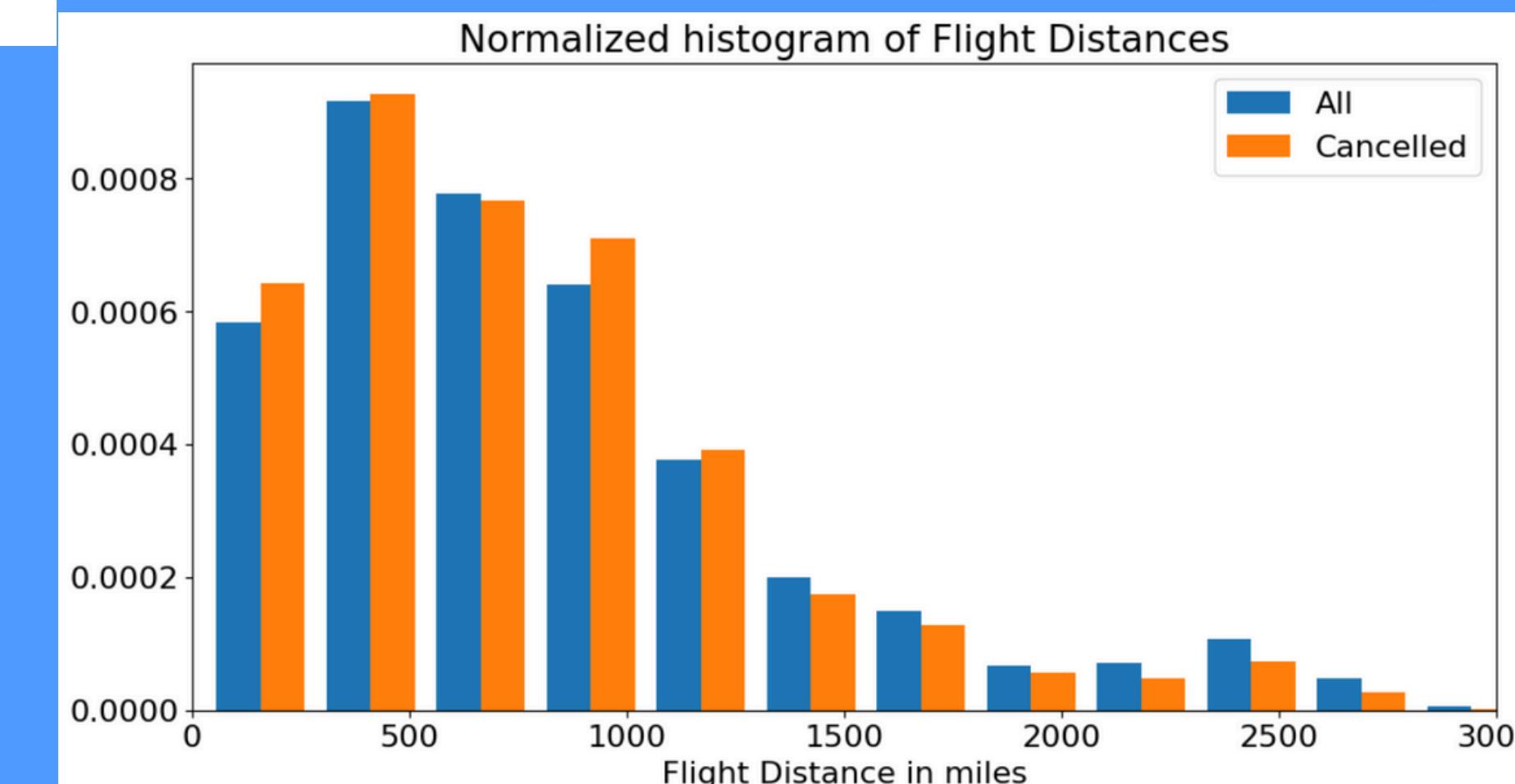
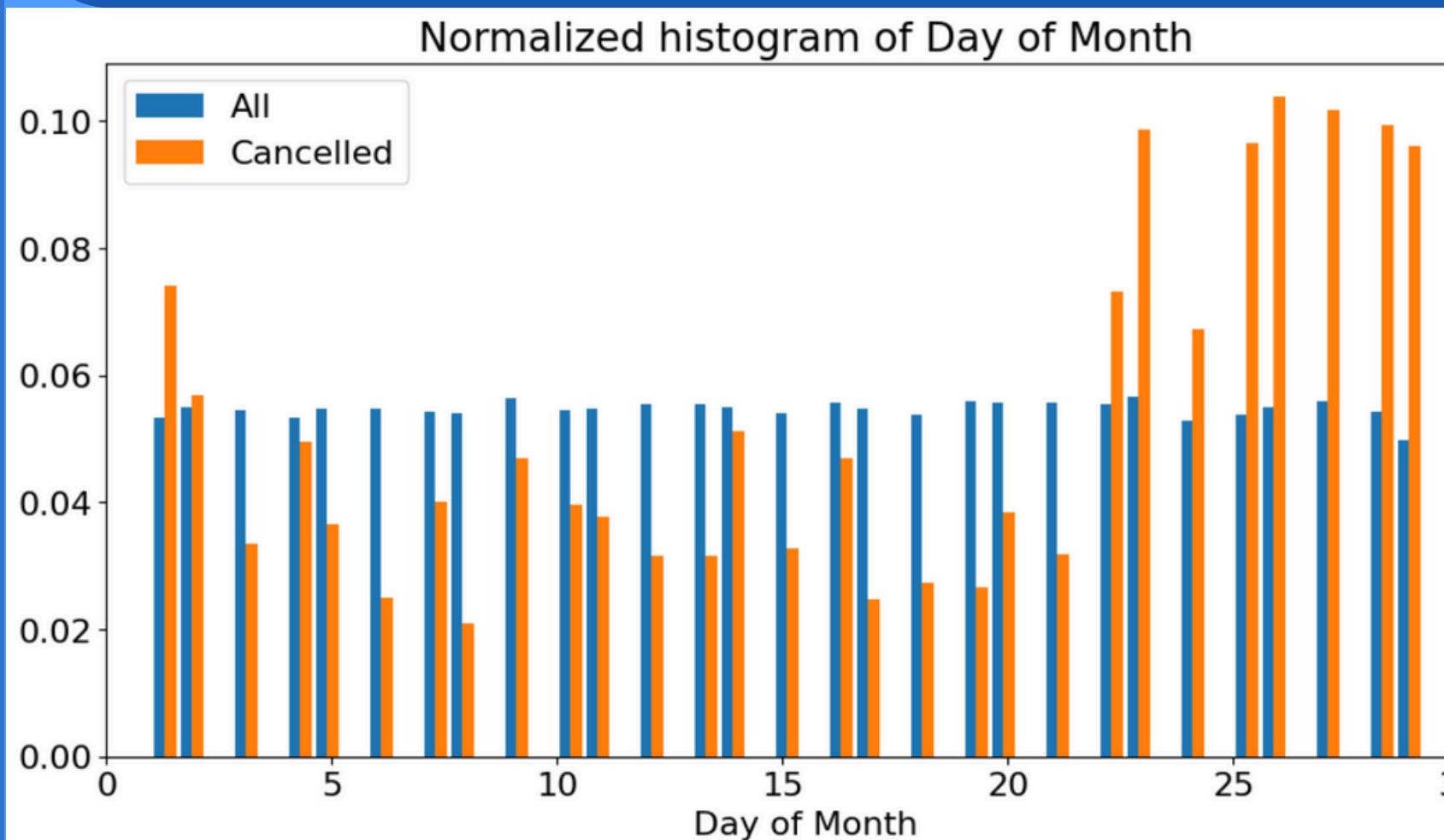
FEATURE ENGINEERING



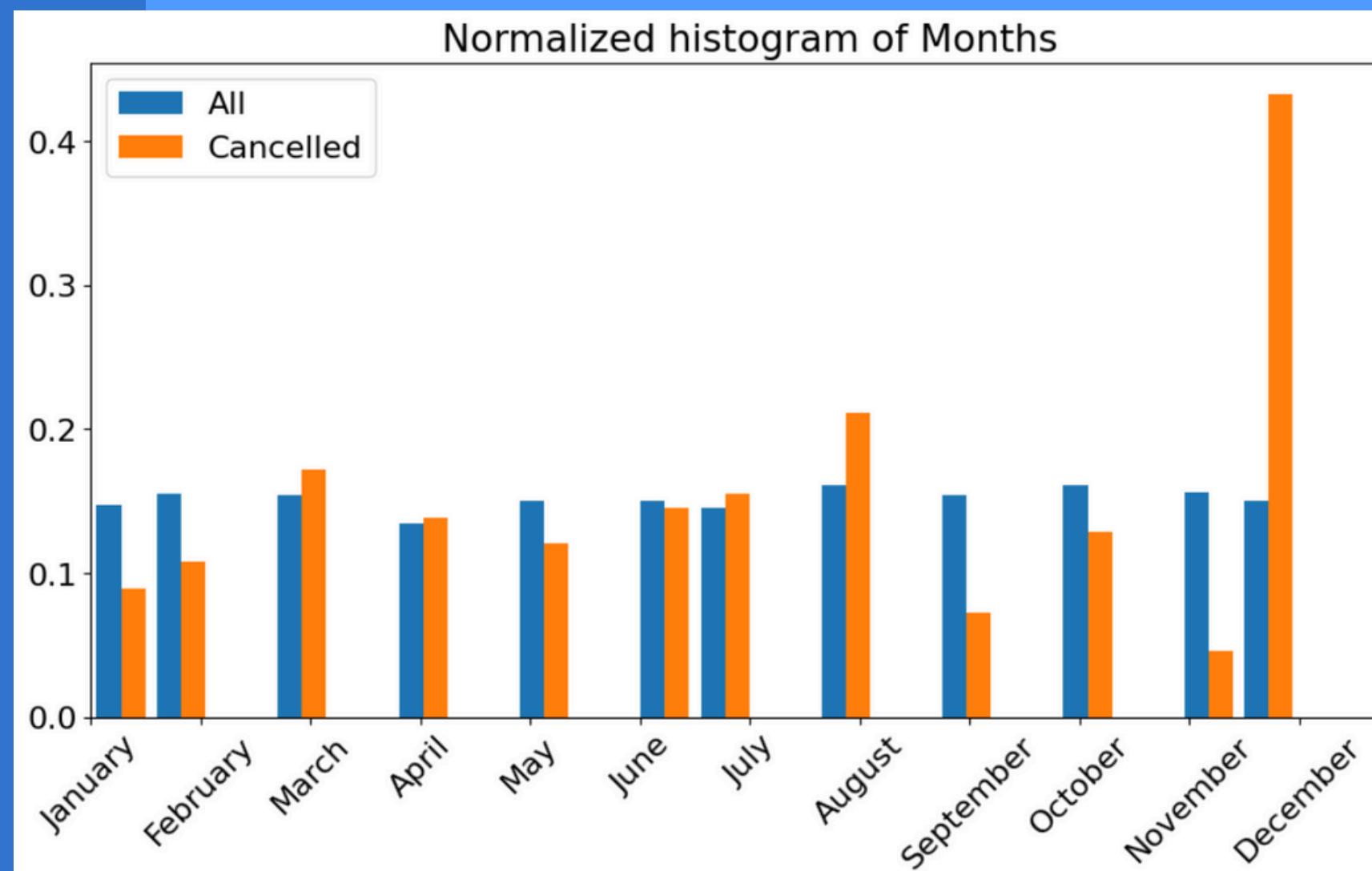
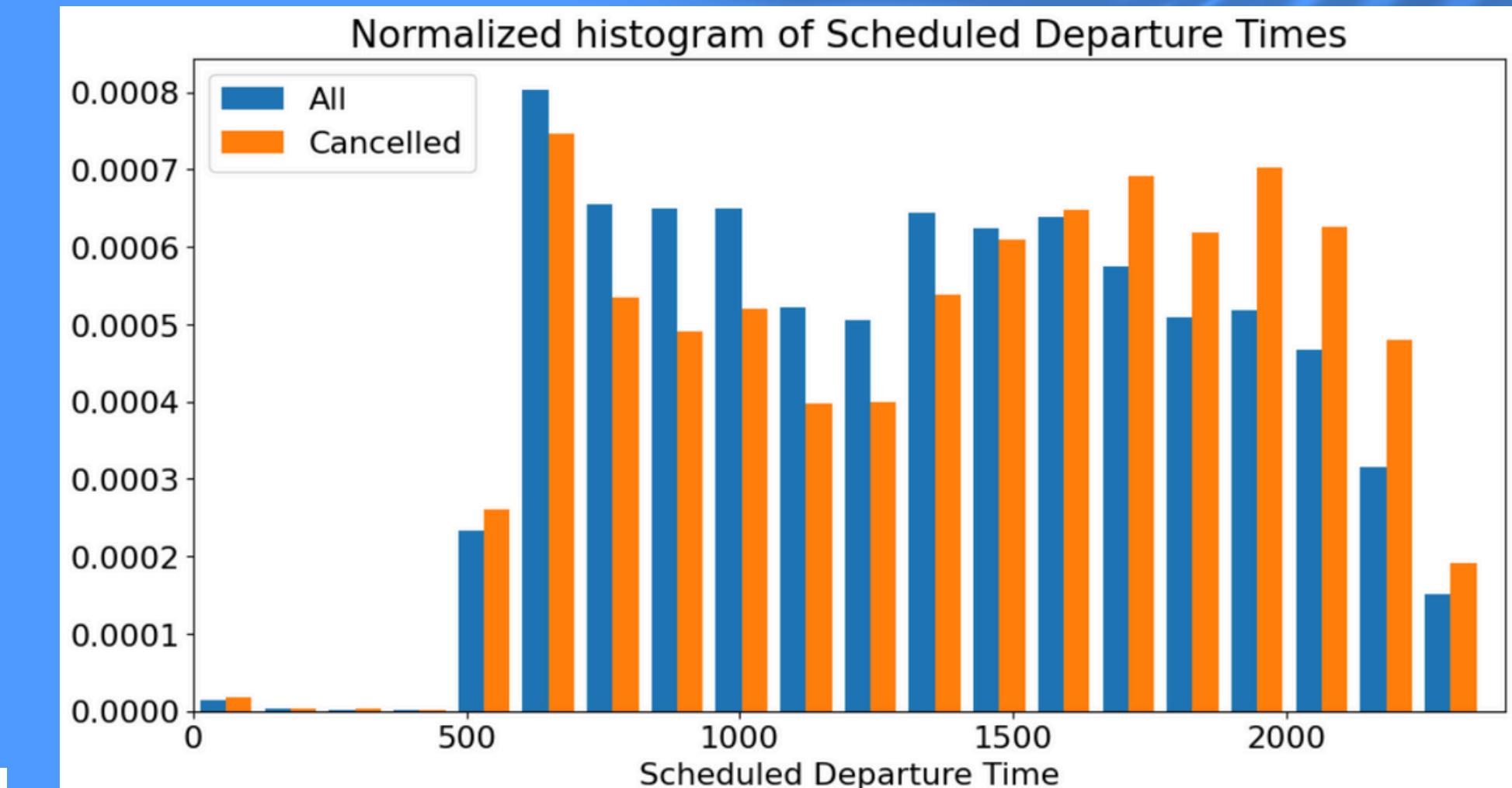
CANCELLATIONS DAY OF WEEK



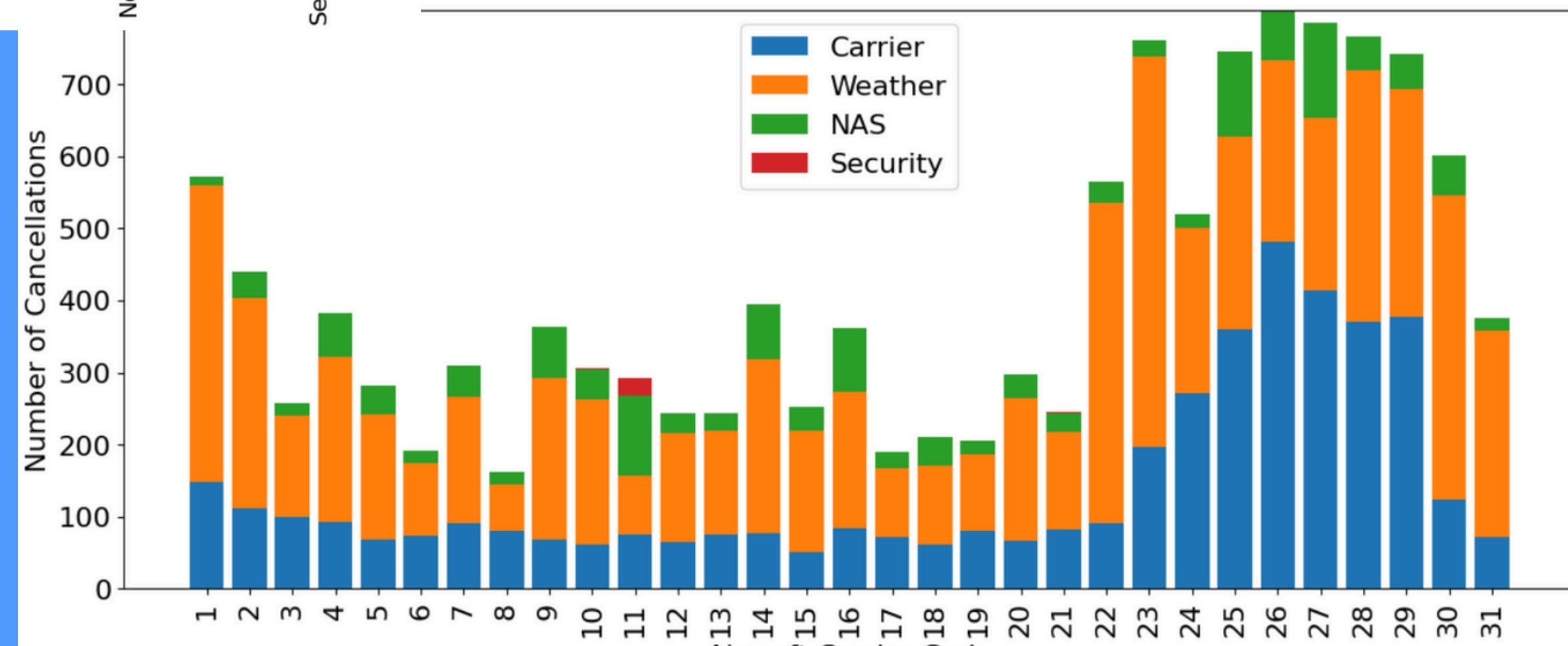
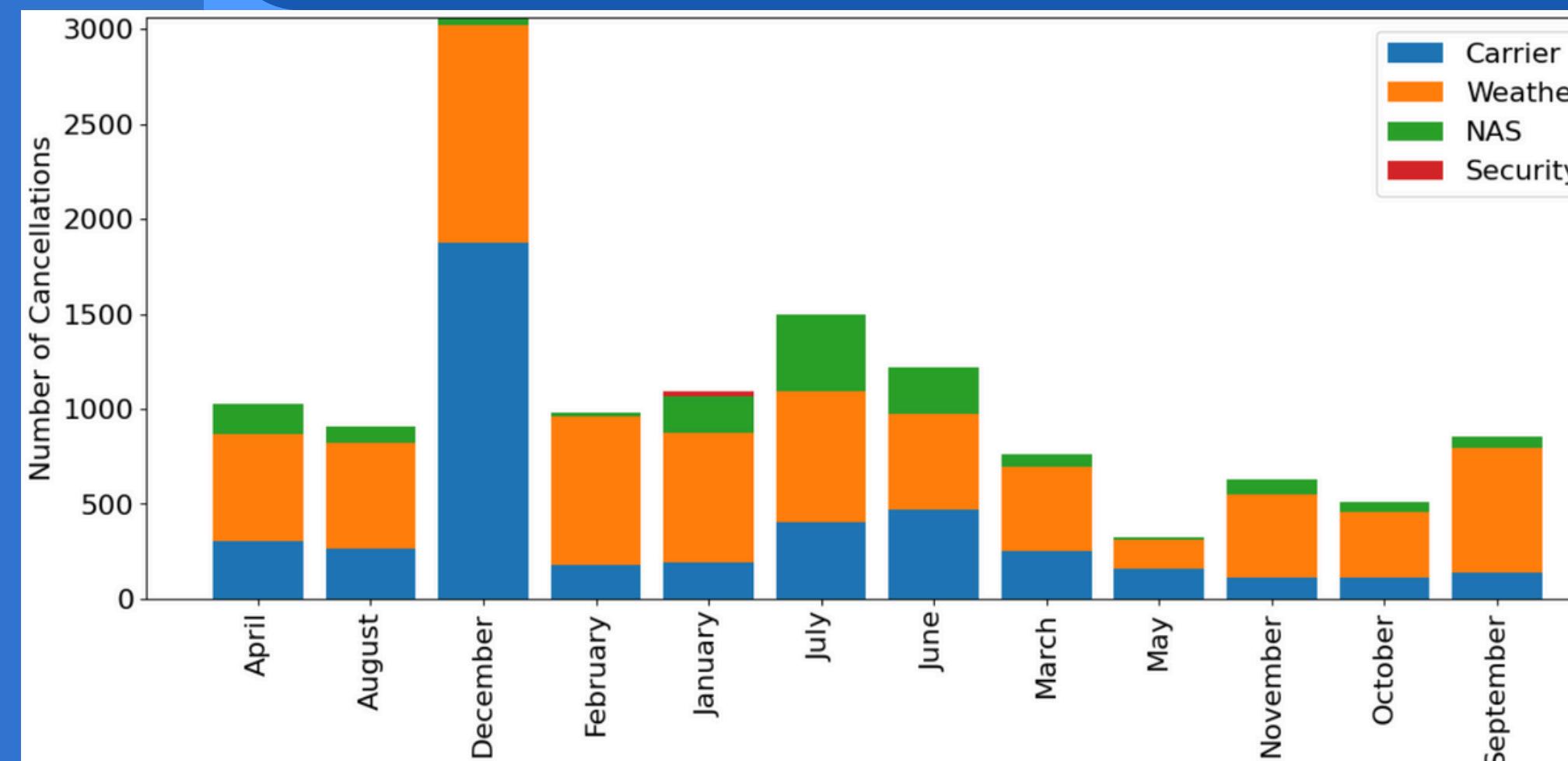
CANCELLATIONS DAYS OF MONTH AND DISTANCE



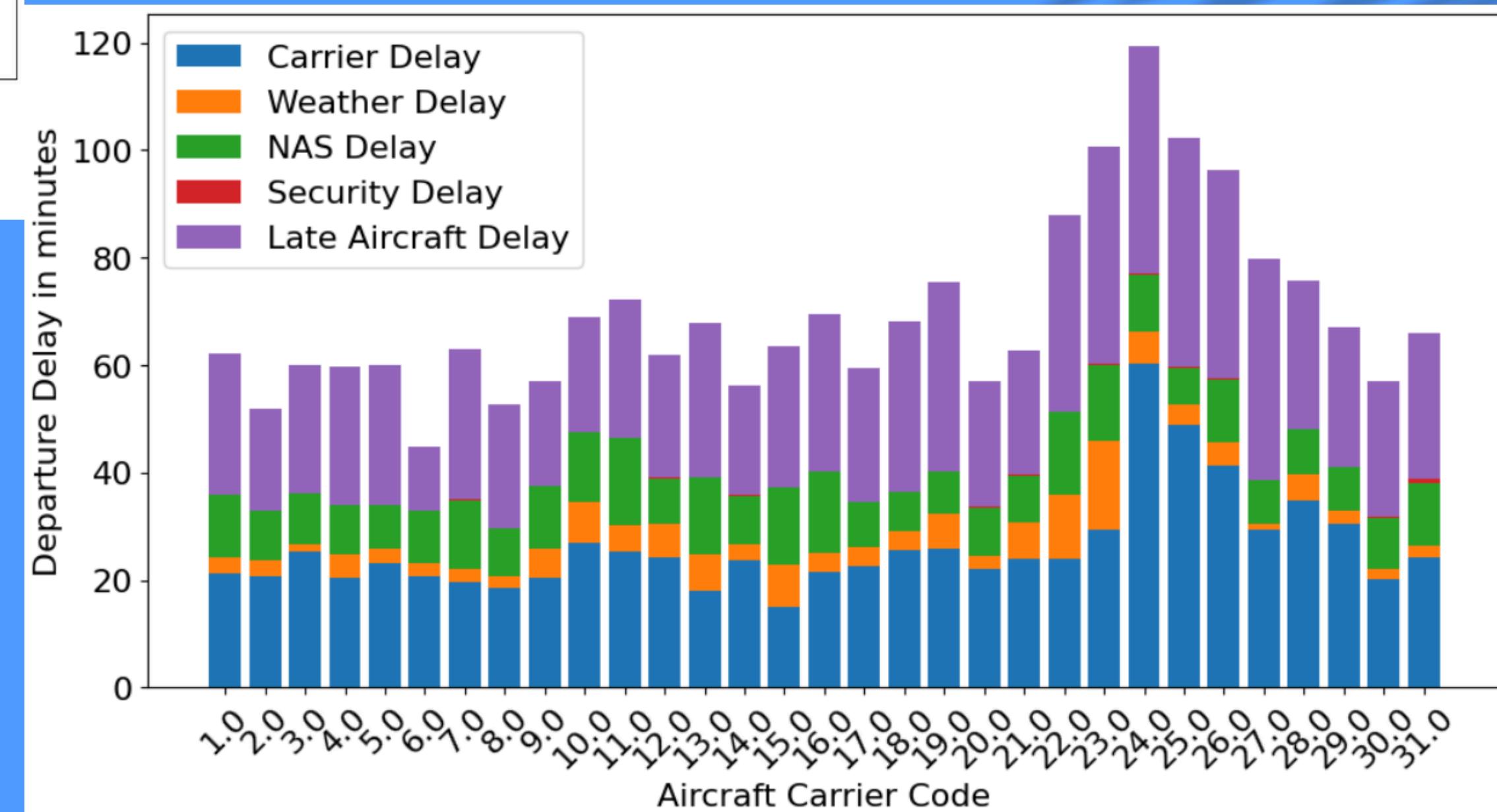
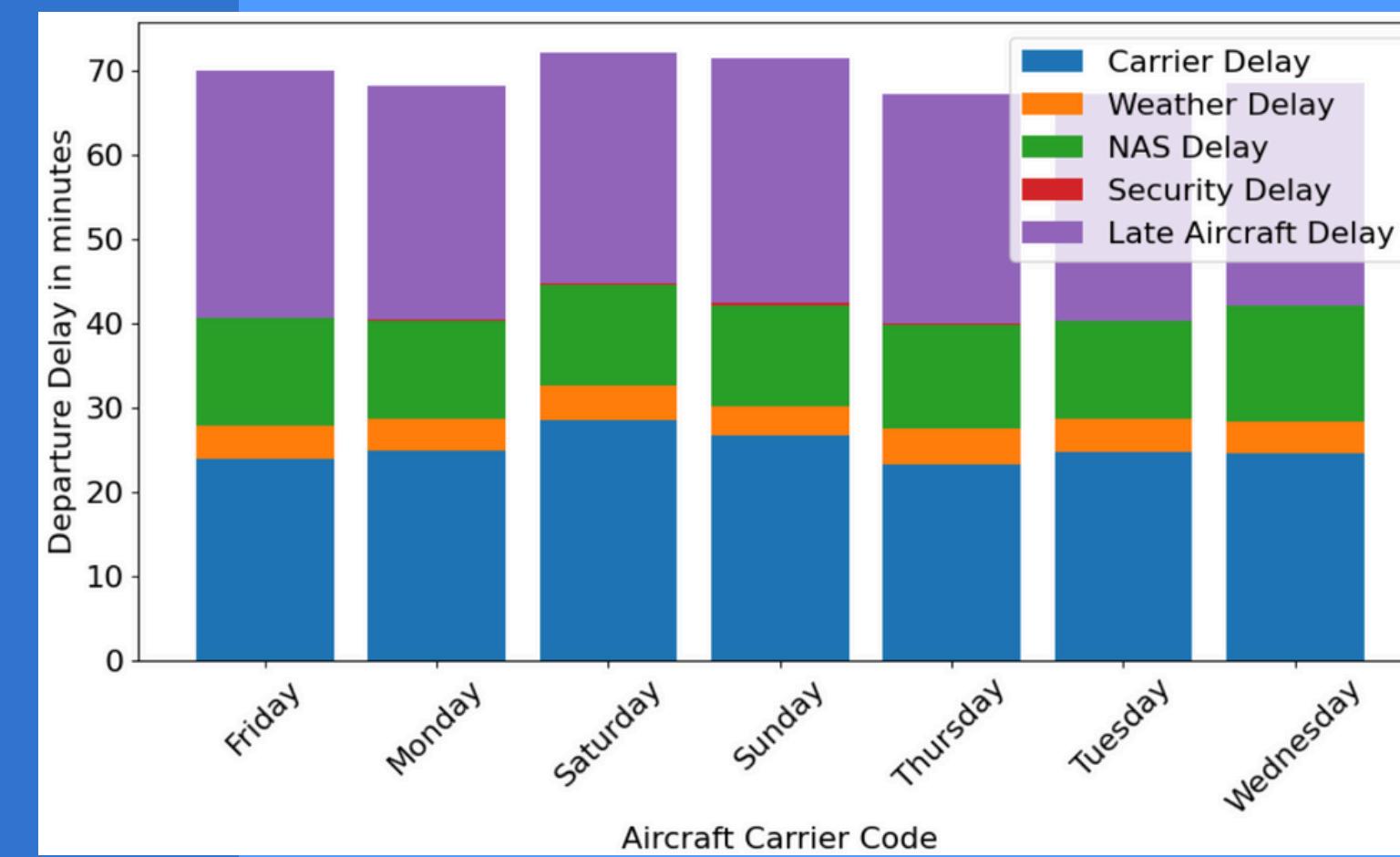
CANCELLATIONS



CANCELLATIONS DAYS OF DECEMBER

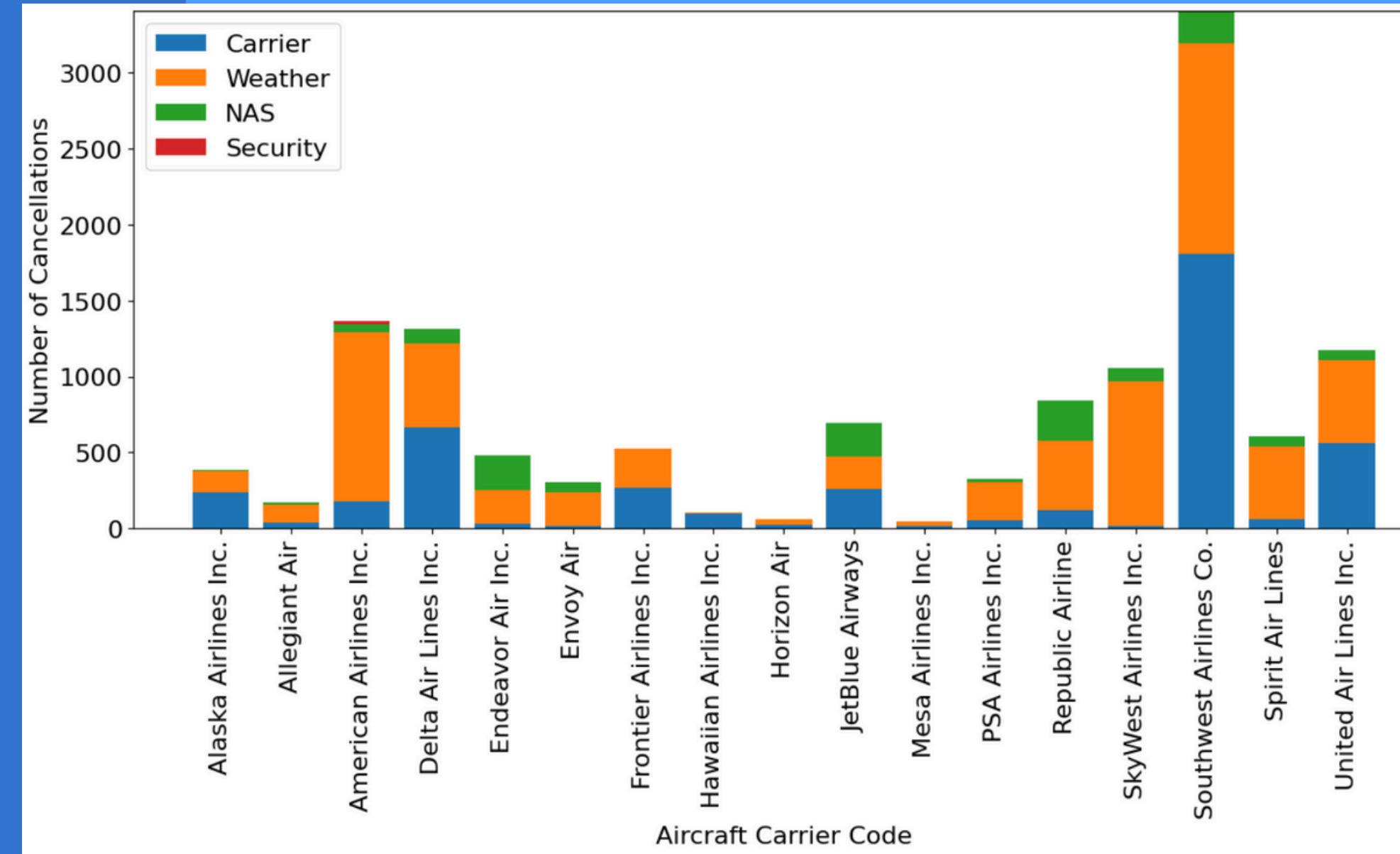


DELAY

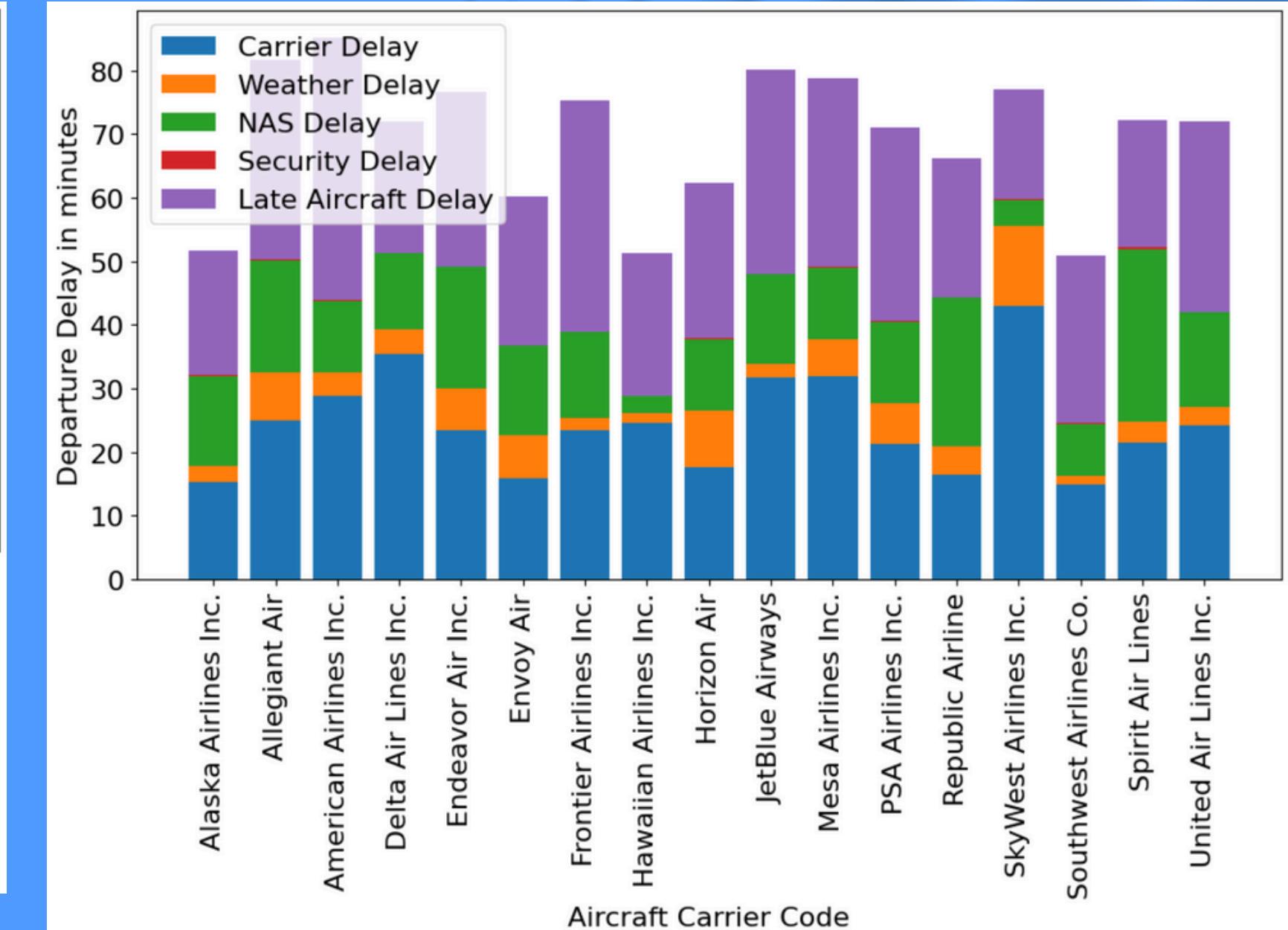


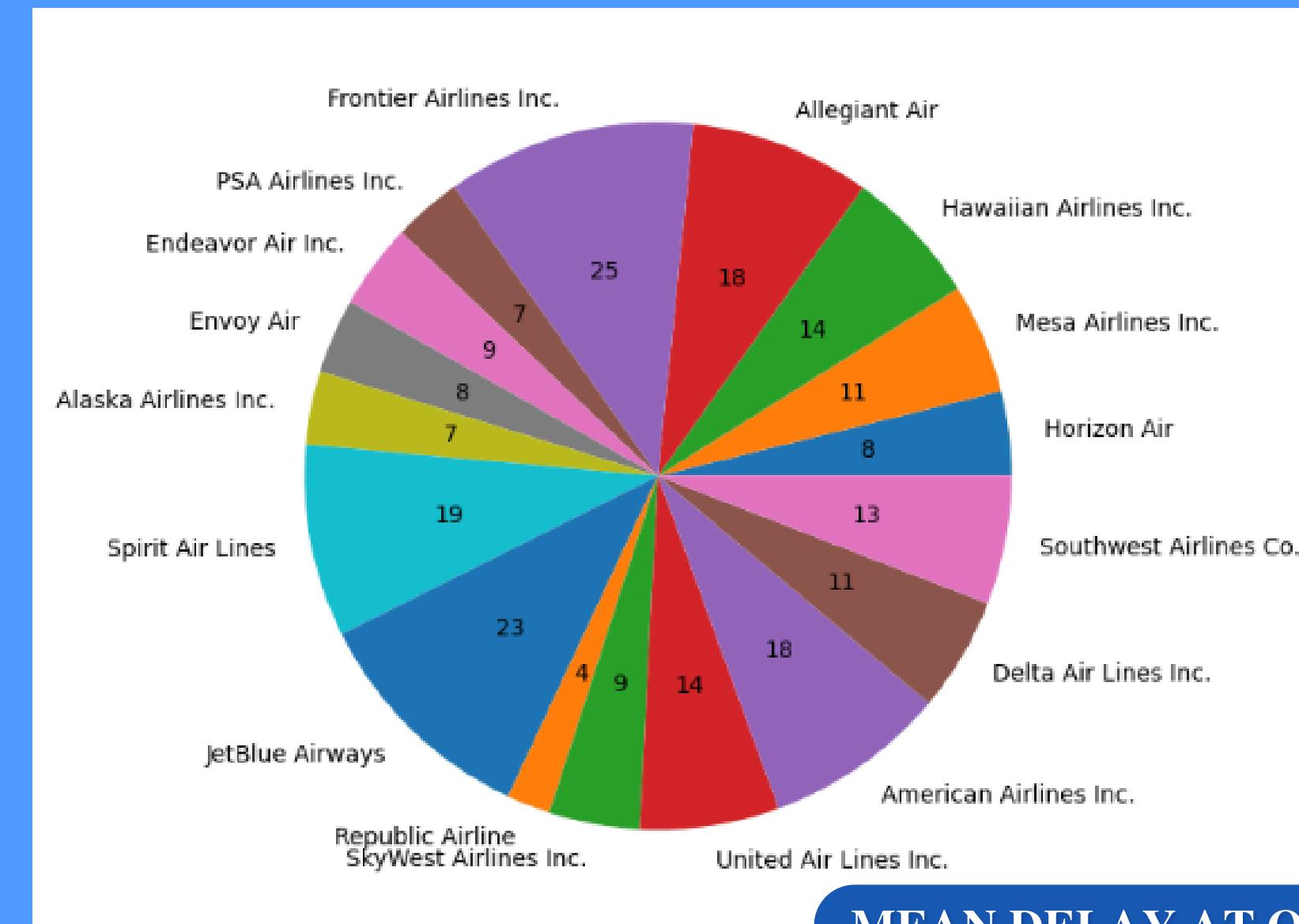
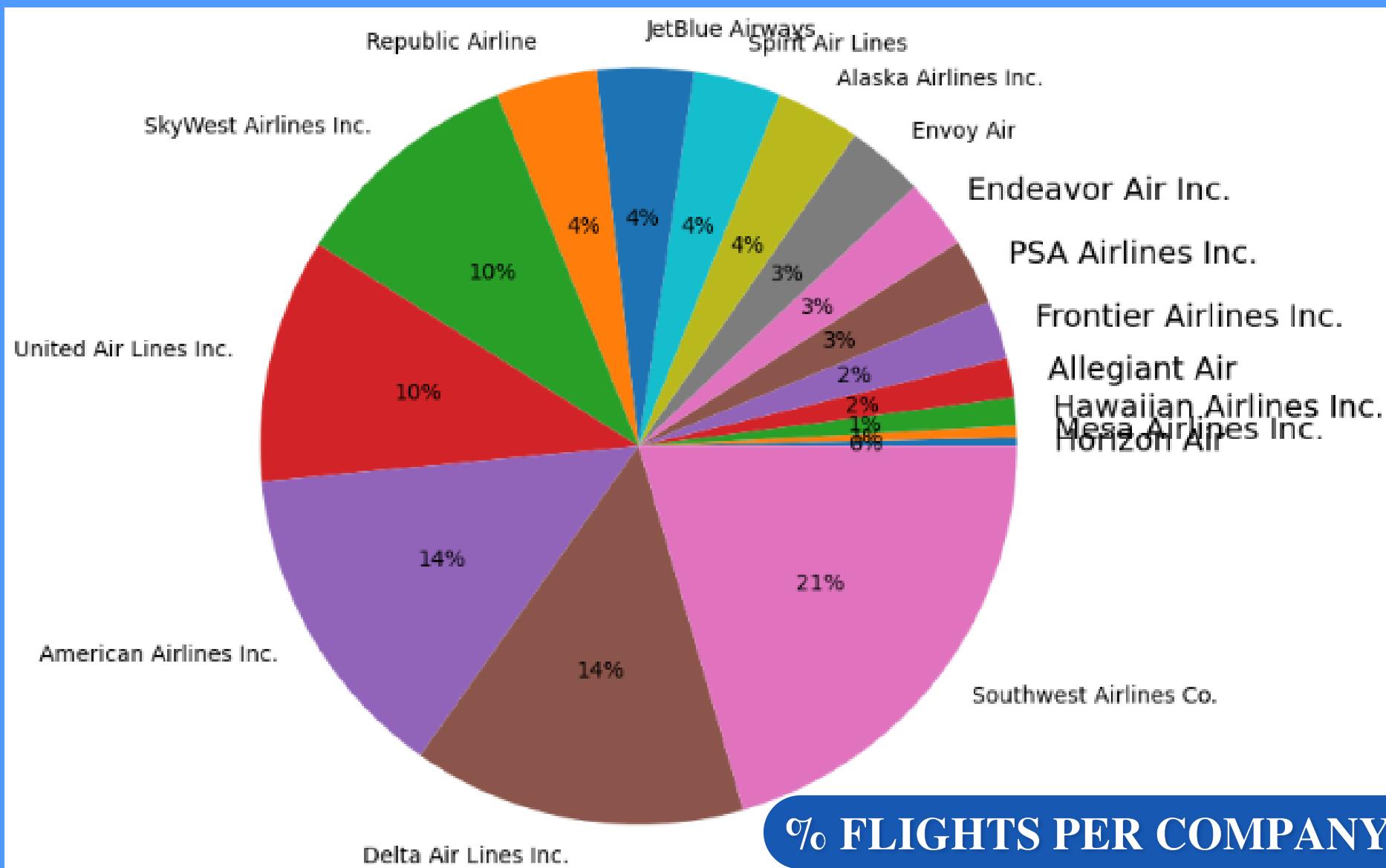
AIRLINE

CANCELLATION

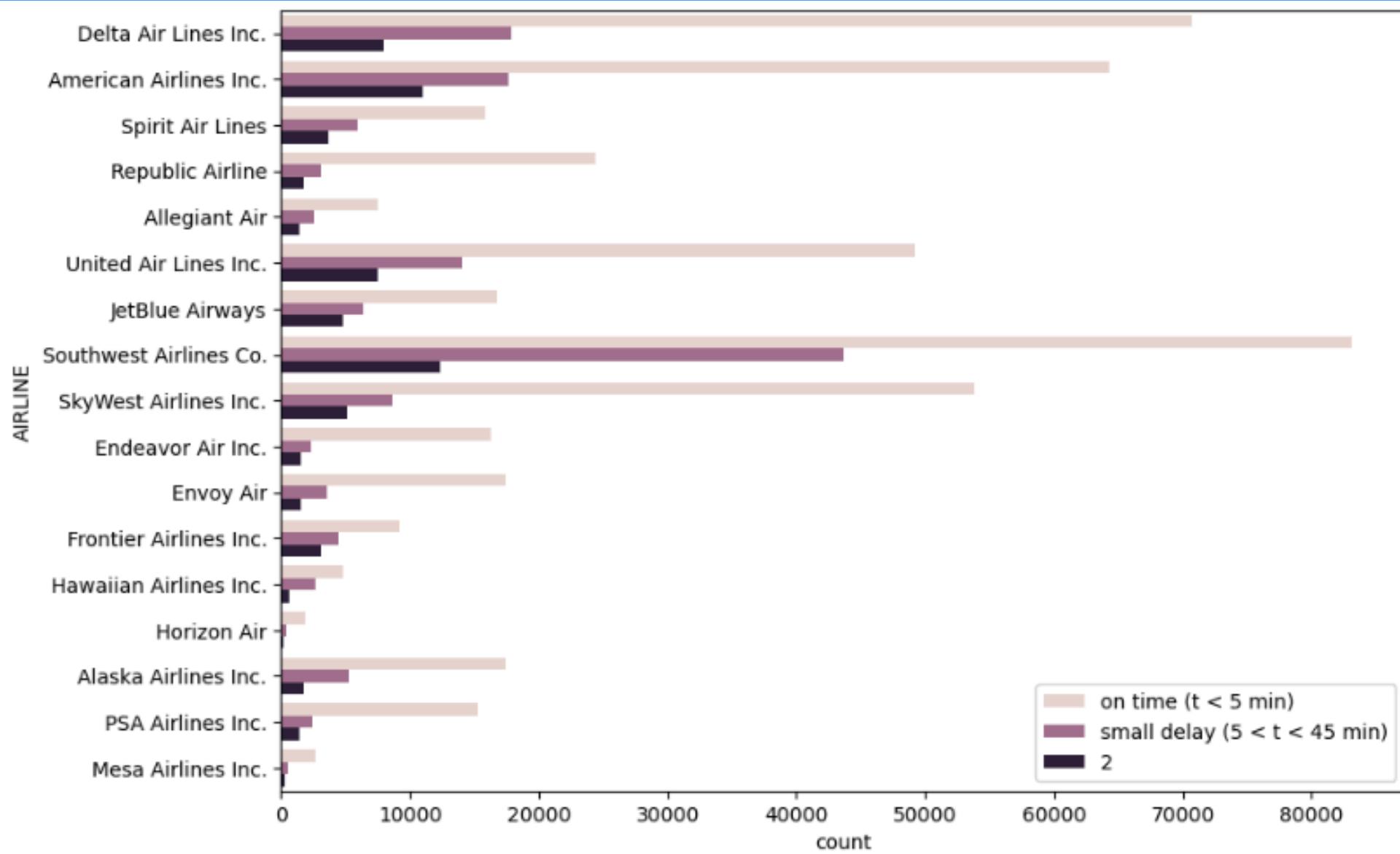


DELAY

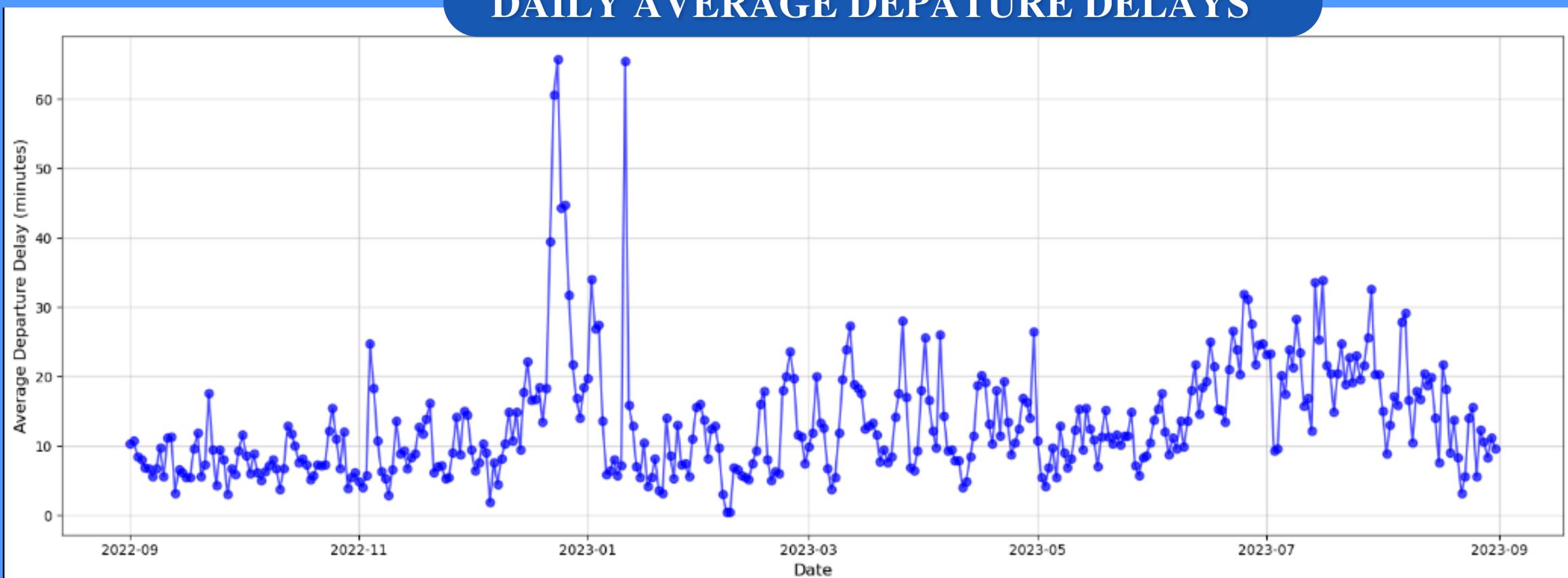




DELAY MAGNITUDE

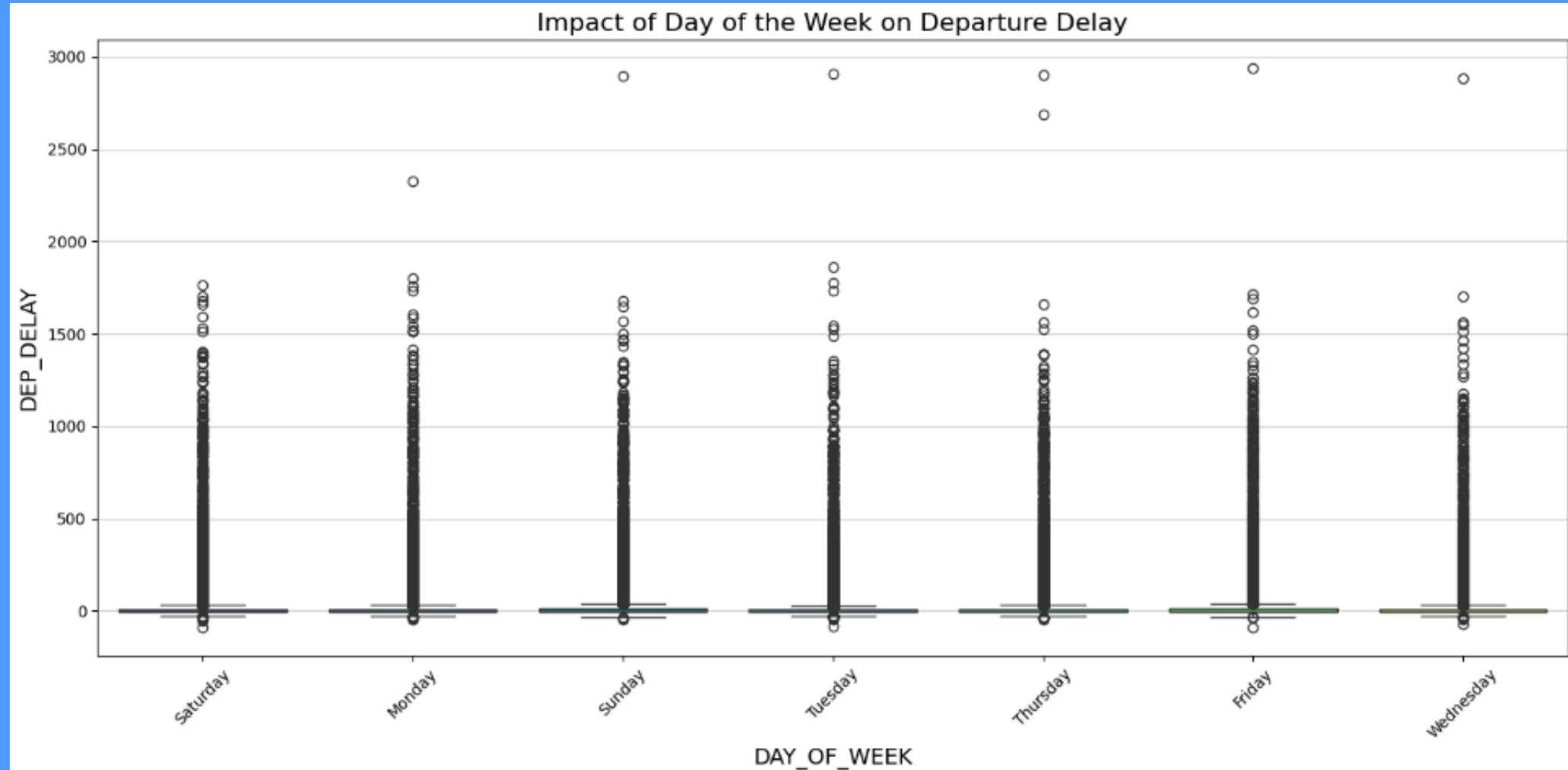


COUNTING THE NUMBER OF DELAYS
AS A FUNCTION OF TIME

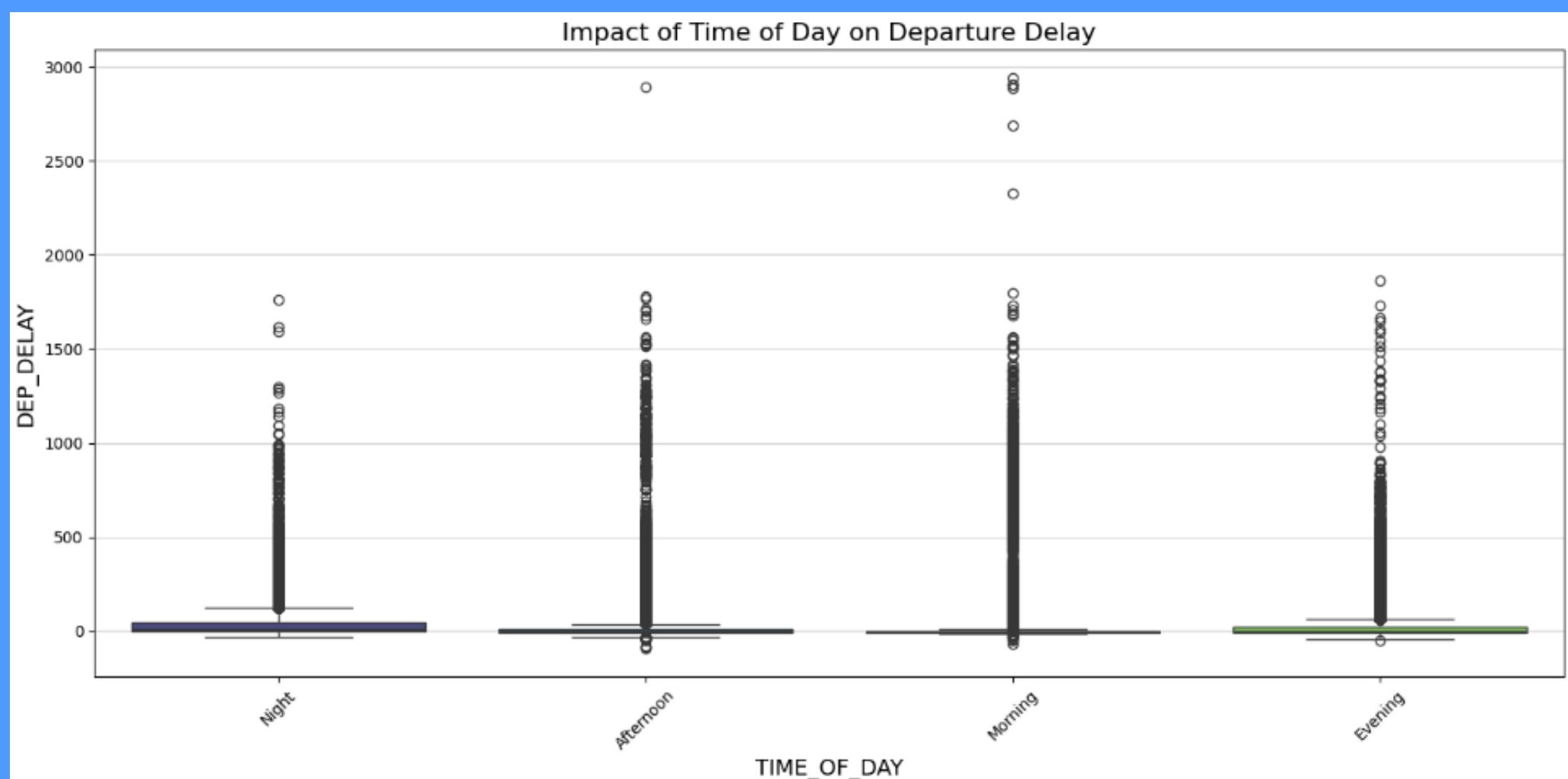


DELAY INFORMATION

The day with the fewest delays is Monday, they increase throughout the week until they reach a peak on Saturdays and begin to decrease on Sundays, with the lowest peak of the week being Mondays.



Most delays happen in the morning, and as the day goes on, they decrease.



FEATURE ENGINEERING

KNN NO NORMALIZED

```
knn.score(X_test_norm, y_test)  
✓ 1m 0.2s  
0.9082972931929555
```

KNN NORMALIZED

```
knn = KNeighborsRegressor(n_neighbors=10)  
knn.fit(X_train_reduced, y_train)  
  
knn.score(X_test_reduced, y_test)  
✓ 50.1s  
0.9269747148646493
```



MODEL BUILDING



MODEL 1

```
bagging_reg = BaggingRegressor(DecisionTreeRegressor(max_depth=20),  
                                n_estimators=100,  
                                random_state=42)
```

```
bagging_reg.fit(X_train_norm, y_train)
```

```
▶ BaggingRegressor
    (i) (?)
```



```
▶   estimator:
        DecisionTreeRegressor
            ▶   DecisionTreeRegressor
                (?)
```

```
pred = bagging_reg.predict(X_test_norm)
bag_acc = bagging_reg.score(X_test_norm, y_test)

print(f"MAE: { mean_absolute_error(pred, y_test): .2f}")
print(f"RMSE: { mean_squared_error(pred, y_test): .2f}")
print(f"R2 score: { bagging_reg.score(X_test_norm, y_test): .2f}")
```

MAE: 7.21
RMSE: 102.97
R2 score: 0.97

MODEL 2

```
forest = RandomForestRegressor(n_estimators=100,  
| | | | | | | random_state=42)
```

```
forest.fit(X_train_norm, y_train)
```

RandomForestRegressor

```
pred = forest.predict(X_test_norm)
forest_acc = forest.score(X_test_norm, y_test)
```

```
print(f"MAE: { mean_absolute_error(pred, y_test): .2f}")
print(f"RMSE: { mean_squared_error(pred, y_test): .2f}")
print(f"R2 score: { forest.score(X_test_norm, y_test): .2f}")
```

MAE: 6.94
RMSE: 97.44
R2 score: 0.97

Bagging Accuracy: 0.9710040023054798

Random Forest Accuracy: 0.972561693125157

HYPERPARAMETER TUNING & MODEL OPTIMIZATION

GRID SEARCH



Mejor puntuación: 0.9634065124014897

MAE: 7.63

RMSE: 118.67

R2 score: 0.97

RANDOM SEARCH



Mejor puntuación: 0.9638021297043807

MAE: 7.55

RMSE: 117.13

R2 score: 0.97



KEY FINDINGS

- Cancellations and delays in December
- Cancellations and delays because of the airline
- Delays in the morning and fridays



REAL WORLD APPLICATIONS

- Airlines
- Airports
- Frequent travelers

CHALLENGES & LEARNINGS

Study in more detail the reasons for cancellation and delay

Help airlines detect the biggest problems they have and try to solve them

Be prepared to prevent cancellations and delays in special periods



PILAR ESCRIG
VALERIA VALEVASHNIKOVA
ROCÍO ROMO
IVÁN SIMÓN

THANK
YOU

