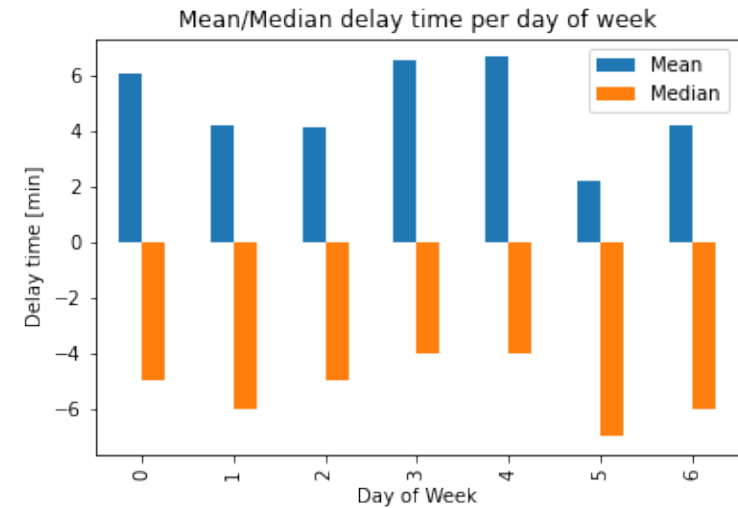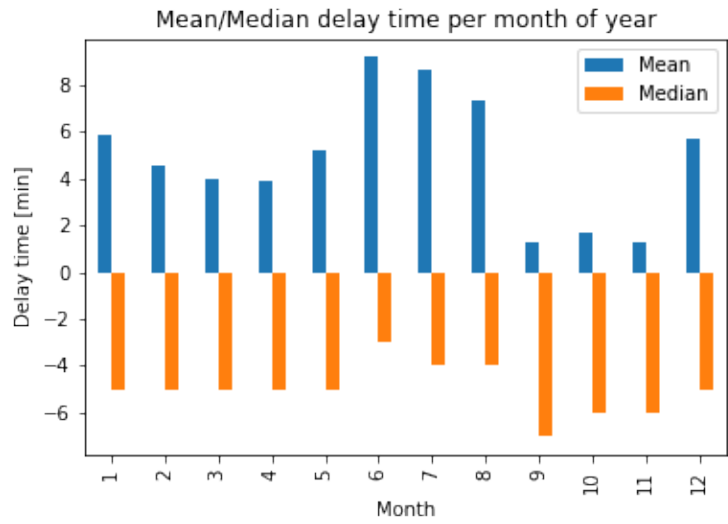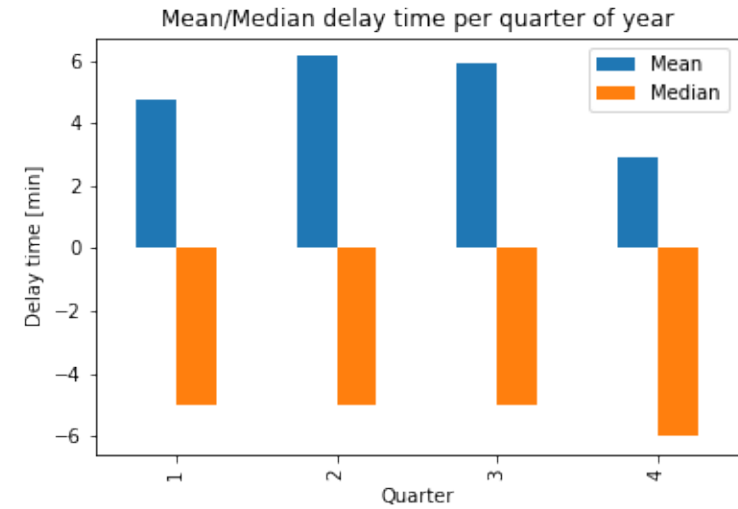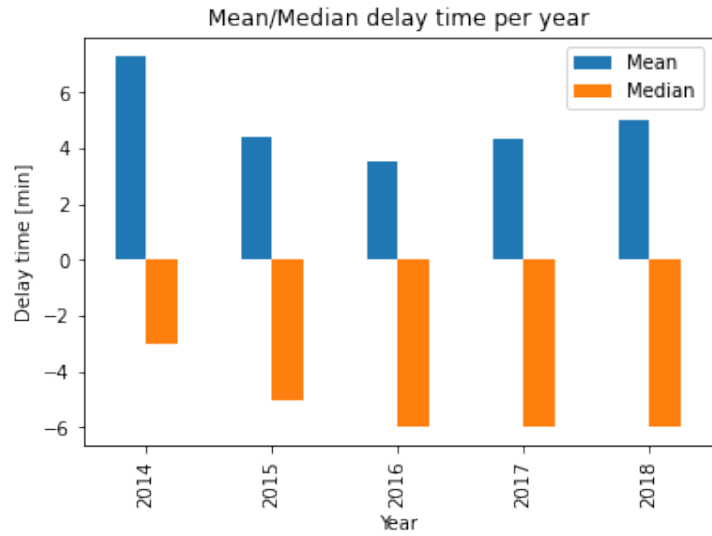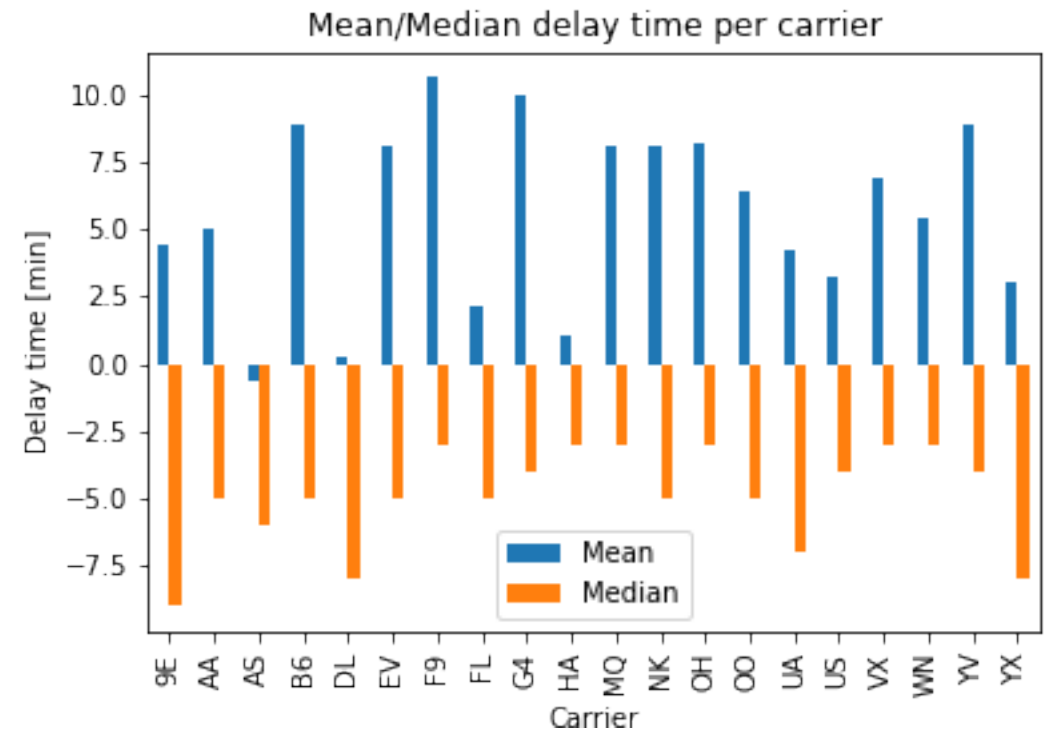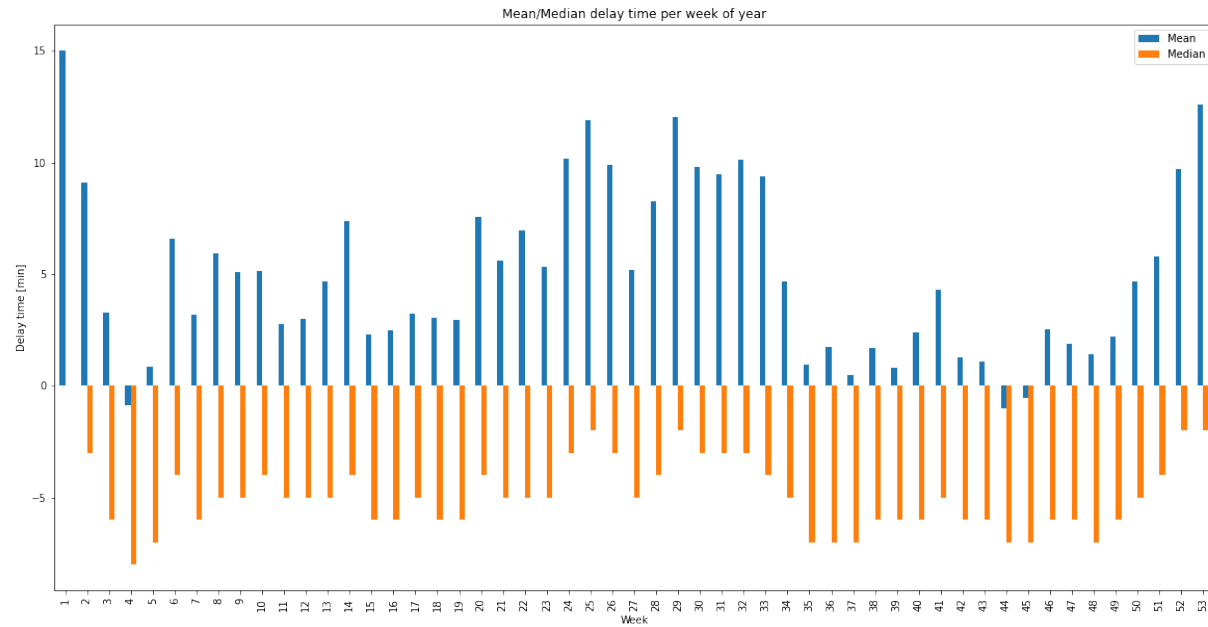# Case study

Vojislav Staletovic

# Mean/Median/Skewness/Kurtosis of delay time

- Mean = 4.94 min

- Median = -5 min

- Skewness = 7.8

- Kurtosis = 132.56

- Comparing mean and median, mean is more sensitive to outliers, so that's the reason why mean has higher value.

- Skewness is a measure of the asymmetry of a distribution, and if it is zero, there is no asymmetry. Skewness is 0 for normal distribution. We have positive skew, and because of it, median has lower value than mean.

- Kurtosis is a statistical measure that defines how heavily the tails of a distribution differ from the tails of a normal distribution. Kurtosis is 0 for normal distribution when it is used Fisher's definition. Extreme positive kurtosis (our case) indicates a distribution where more of the numbers are located in the tails of the distribution instead of around the mean. That means there are more chances of outliers. When kurtosis is positive, it has a leptokurtic distribution.

- Theoretically, I would fit beta distribution for delay time, because it can fit my skewness and kurtosis. It is right skewed.
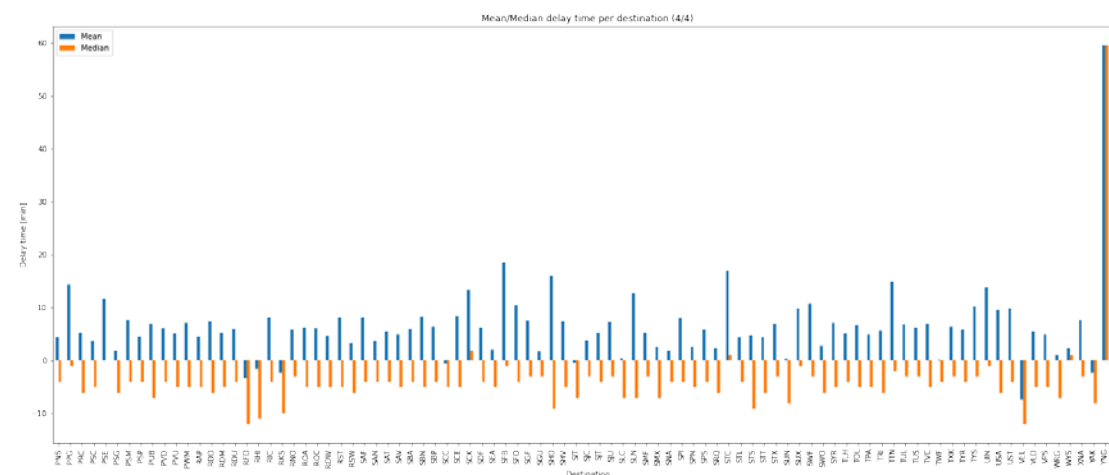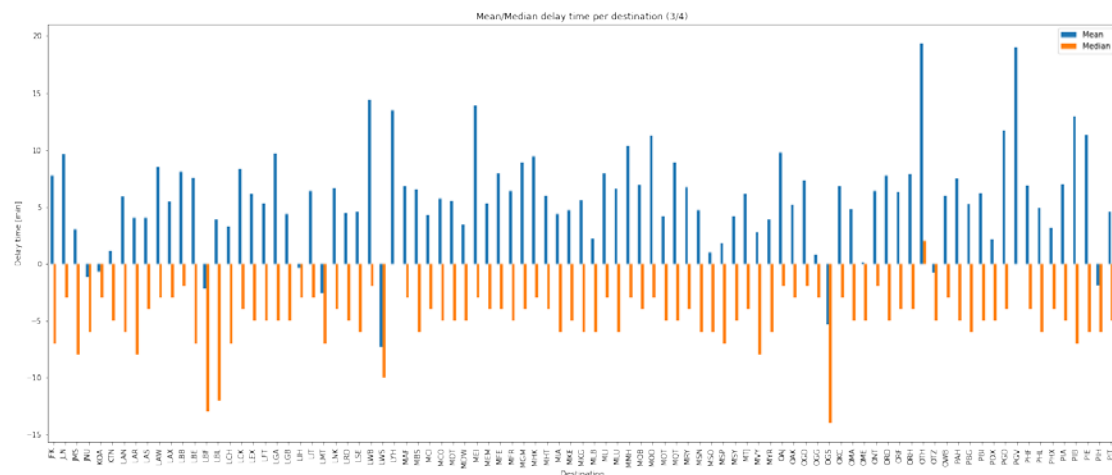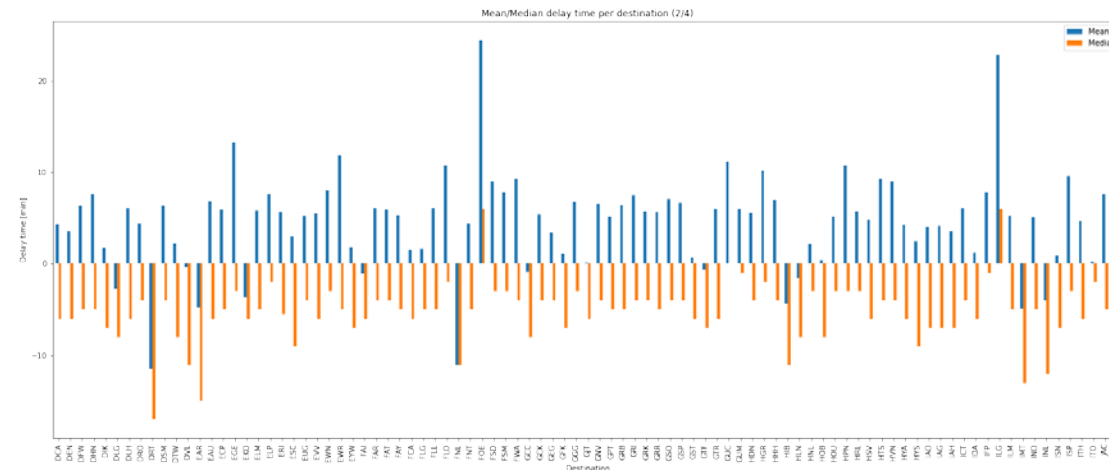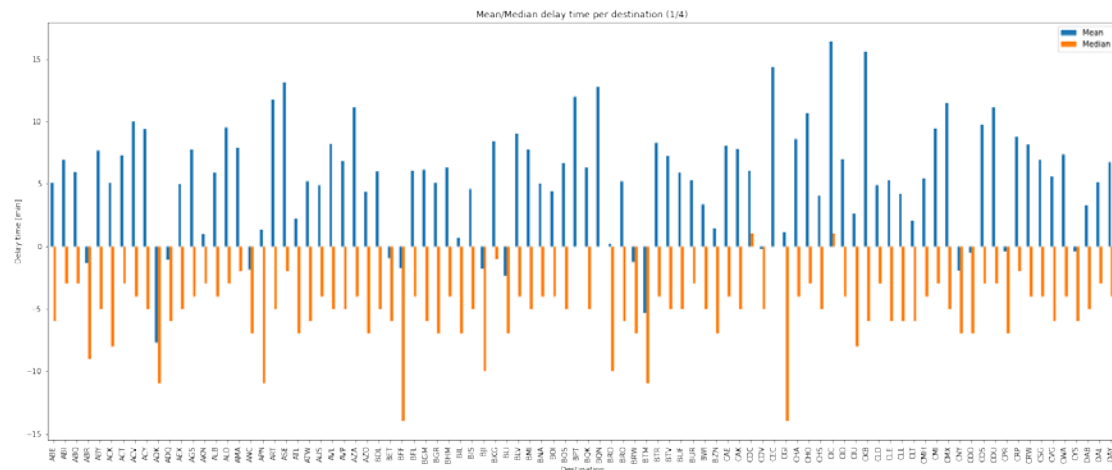
# Mean/Median

# Mean/Median

# Mean/Median



Mean/Median delay time per destination (1/4)

Mean/Median delay time per destination (2/4)

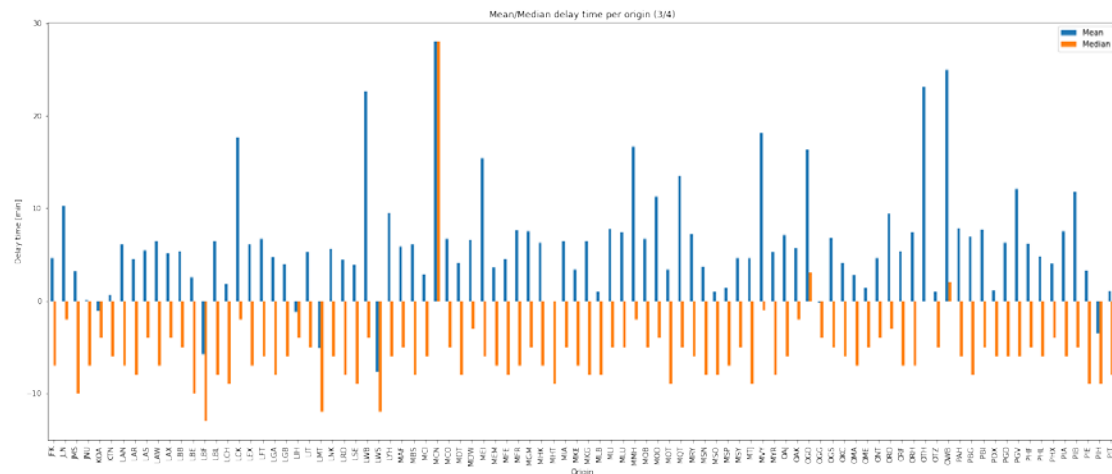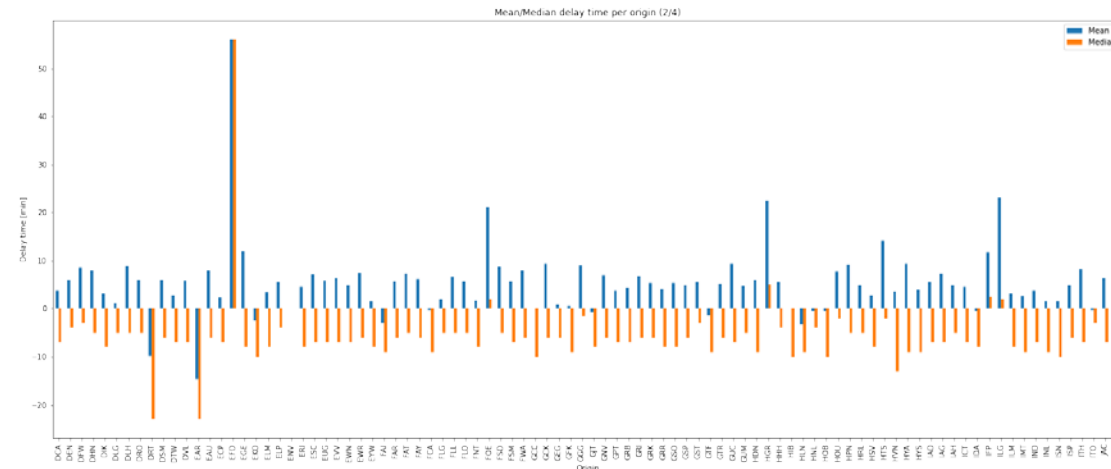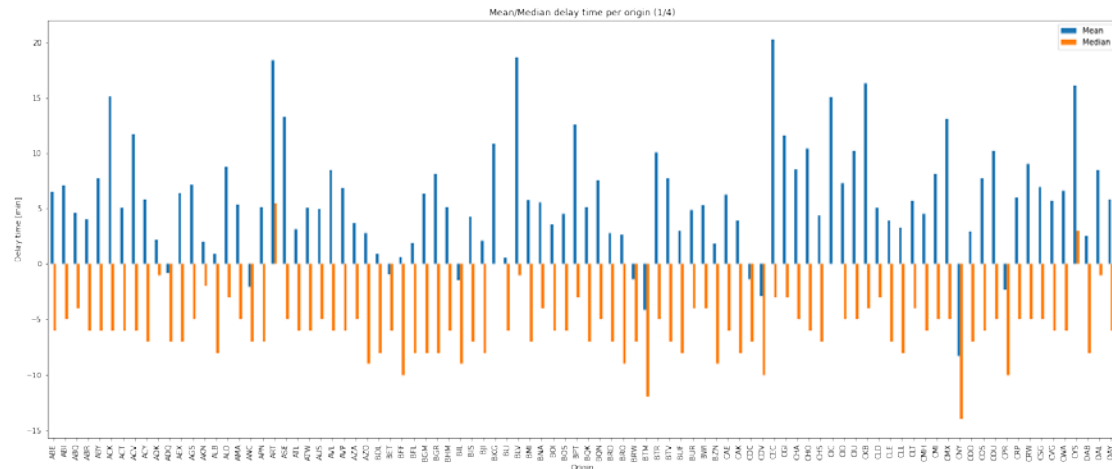Mean/Median delay time per destination (3/4)

Mean/Median delay time per destination (4/4)

# Mean/Median

- Year: minimum delay time was in 2016, and maximum delay time was in 2014
- Quarter of year: minimum delay time is in the fourth quarter of year, and maximum delay time is in the second quarter of year
- Month of year: minimum delay time is in September, and maximum delay time is in June
- Week of year: minimum delay time is in the 44. week, and maximum delay time is in the 1. week of year
- Day of week: minimum delay time is on Saturday, and maximum delay time is on Friday
- Carrier: minimum delay time is for carrier AS, and maximum delay time is for carrier F9
- Origin: minimum delay time is for origin EAR, and maximum delay time is for origin YNG
- Destination: minimum delay time is for destination DRT, and maximum delay time is for destination YNG

# Gradient Boosting Regressor

- Gradient Boosting method is used for the prediction of delay time, and also for the importance of features

- All features are used except Cancellation code because most of the data are missing values

- Month and Day of Week are used as features instead of Date

- Metrics for model evaluation:
  - $R^2 = 0.993$
  - Mean Squared Error = 39.46
  - Mean Absolute Error = 4.19

# Gradient Boosting Regressor

- Order of features, given by measuring the impact on delay time (Gini importance)
    1. DEP_DELAY - departure delay in minutes
    2. NAS_DELAY - delay due to National Aviation System in minutes
    3. TAXI_OUT - taxi out time in minutes
    4. LATE_AIRCRAFT_DELAY - previous flights caused delays (in minutes)
    5. CARRIER_DELAY - delay occurred due to carrier in minutes
    6. TAXI_IN - taxi in time in minutes
    7. WEATHER_DELAY - delay due to extreme weather conditions in minutes
    8. ACTUAL_ELAPSED_TIME - actual elapsed time in minutes
    9. CRS_ELAPSED_TIME - planned elapsed time in minutes
    10. OP_CARRIER - carrier
    11. CRS_DEP_TIME - planned departure time
- Other features don't have any impact on delay time