# Capstone Project
## NYC Taxi Trip Time Prediction

**Sandesh Gaikwad**

# Content

- Problem Statement

- Exploratory Data Analysis

- Feature Engineering

- Modeling and comparing results

- Conclusions

# Problem Statement

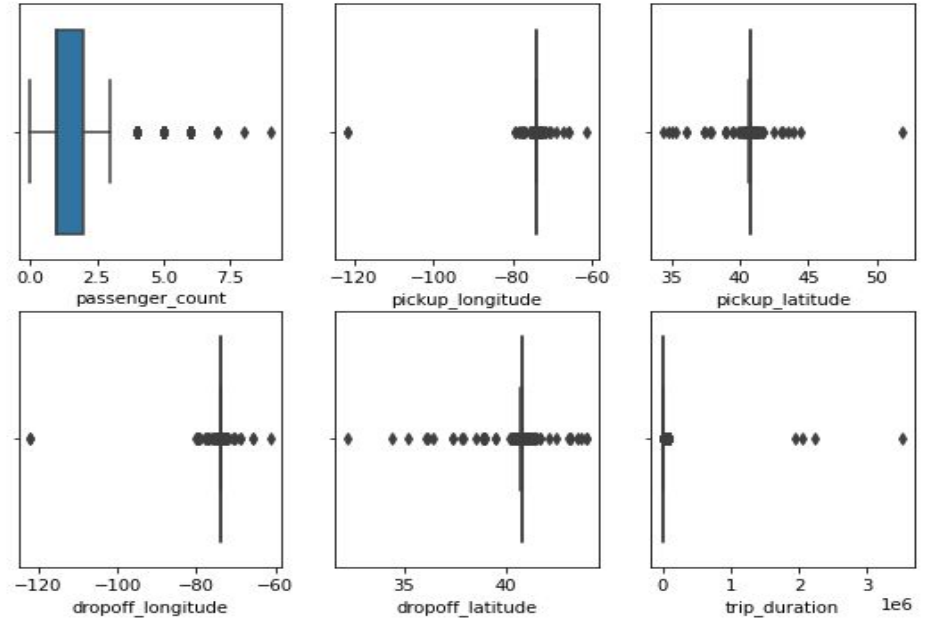Problem : Build a model that predicts the total ride duration of taxi trips in New York City.

Data Description - The dataset is based on the 2016 NYC Yellow Cab trip record data made available in Big Query on Google Cloud Platform. The data was originally published by the NYC Taxi and Limousine Commission (TLC). The data was sampled and cleaned for the purposes of this project. Based on individual trip attributes, we need predict the duration of each trip in the test set.

# Data Summary

- Id : A unique identifier for each trip

- Passenger count : Number of passengers in the vehicle (driver entered value)

- Location : Four columns with longitude and latitude of pickup and drop Locations

- Time: Two columns with pickup and  dropoff date time

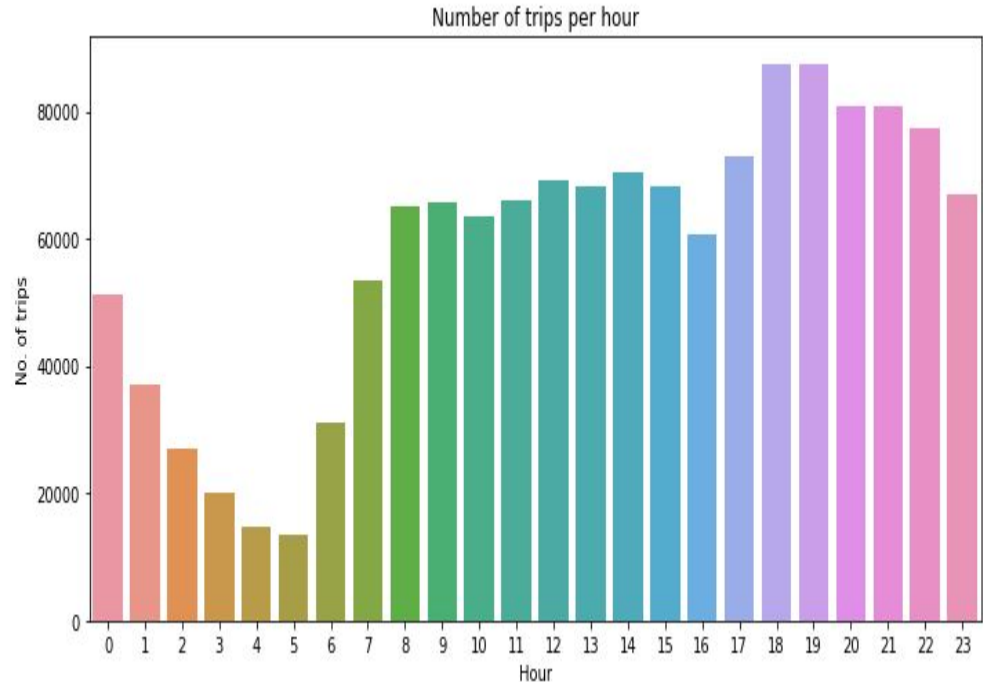- Trip duration : duration of the trip in seconds

# Outliers detection (Box Plots)

- Every Parameter has extreme values
- We can remove most extreme value and later remove remaining anomalies
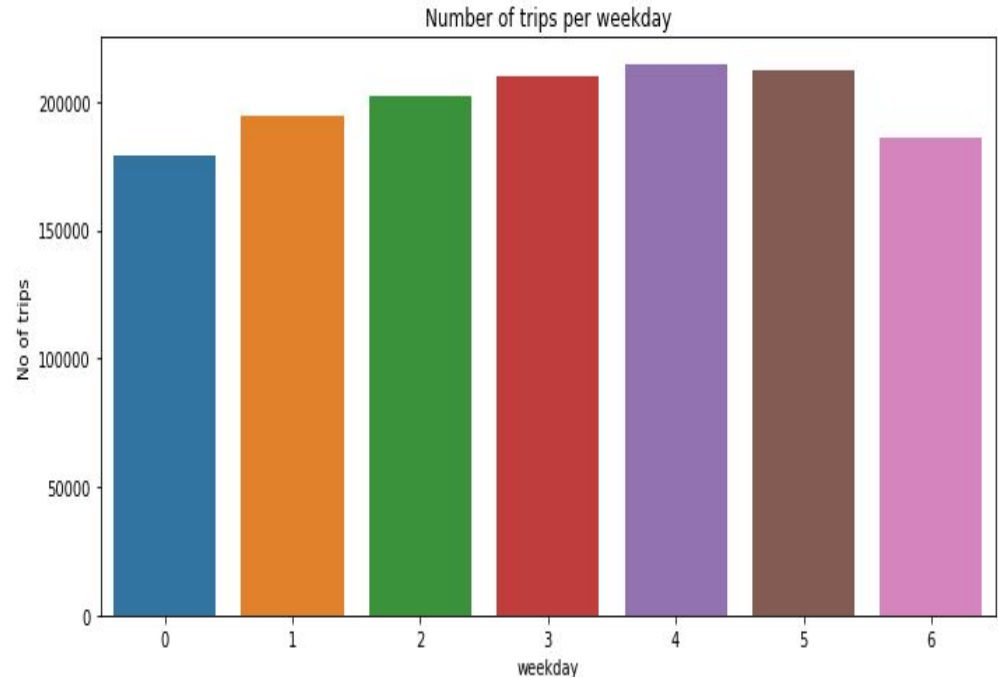
# Avg trip per hour

- Number of trips start increasing at 6 early morning
- Then morning to late afternoon trips remain almost same
- Evening has highest number of trip per hour
- Late night trip are least



Number of trips per hour

# Avg trip per day

- There is no significant change in trips per day

- Sunday has least tris per day while friday and saturday highest
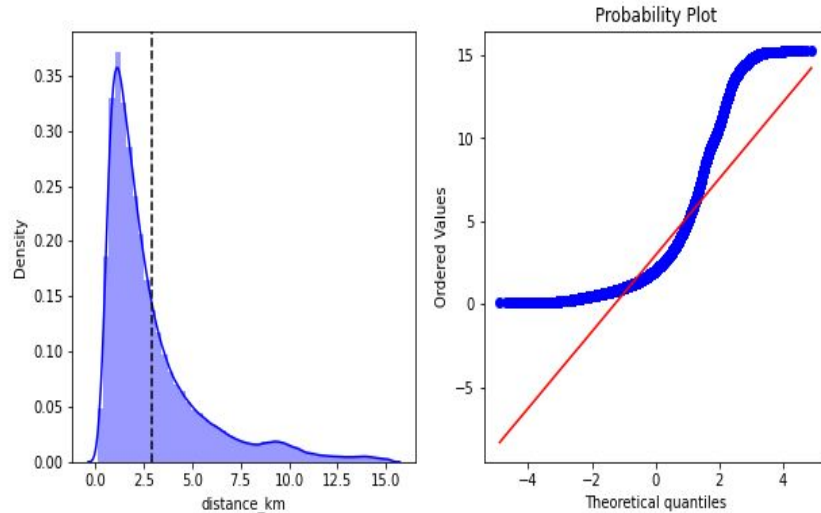
# Correlation matrix

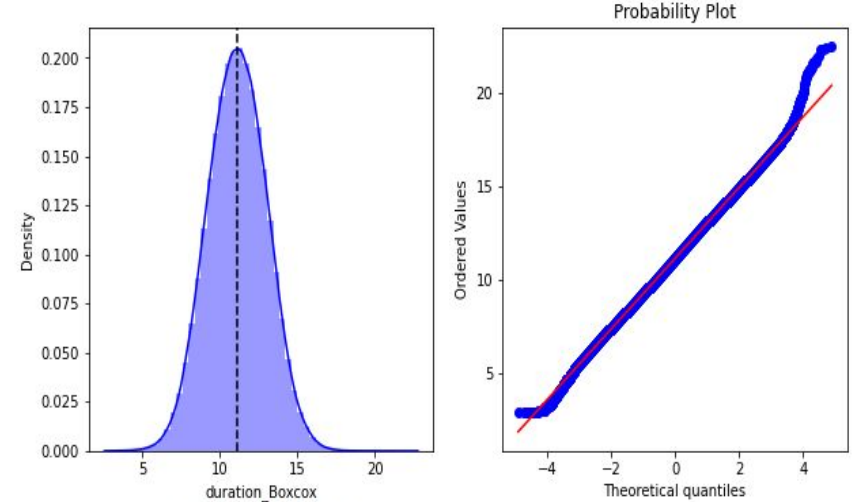- Distance has the highest correlation with trip duration

# Feature Engineering

- Features extractions from "Pickup Datetime"
  - Month
  - day
  - is_Weekend
  - Hour
  - Peak Hour

- New Derived features from Latitudes and Longitudes
  - Delta longitude / latitude
  - Direction
  - Distance_km
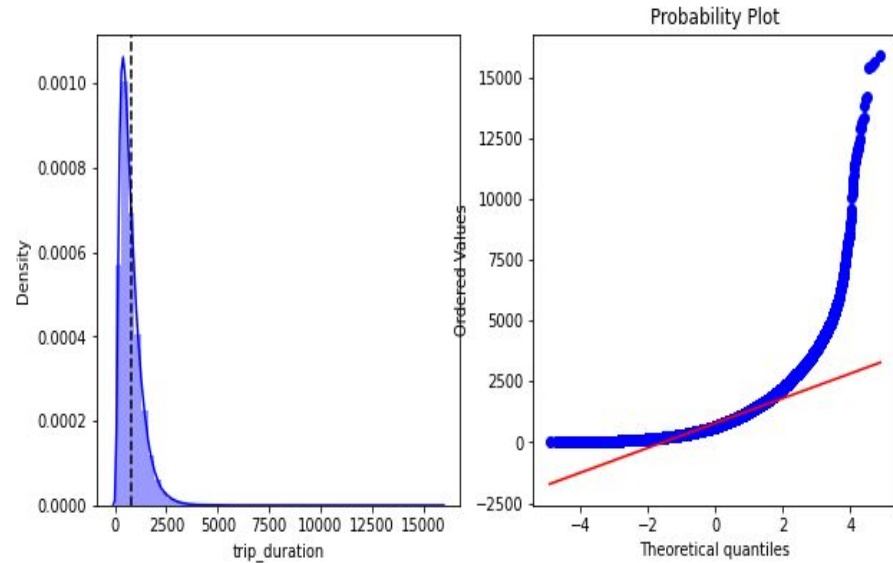
# Gaussian Distribution - distance
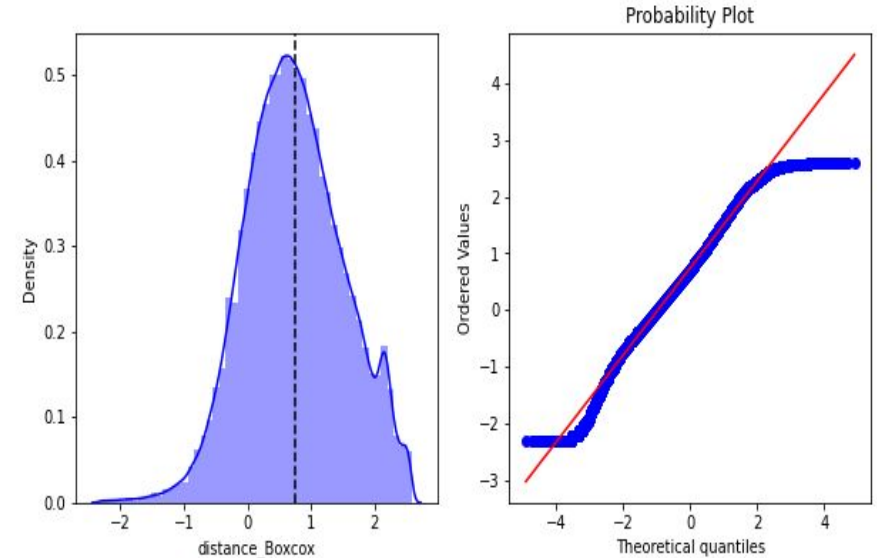


skewness : 1.9517243542550862

skewness : -0.0021837506737579933

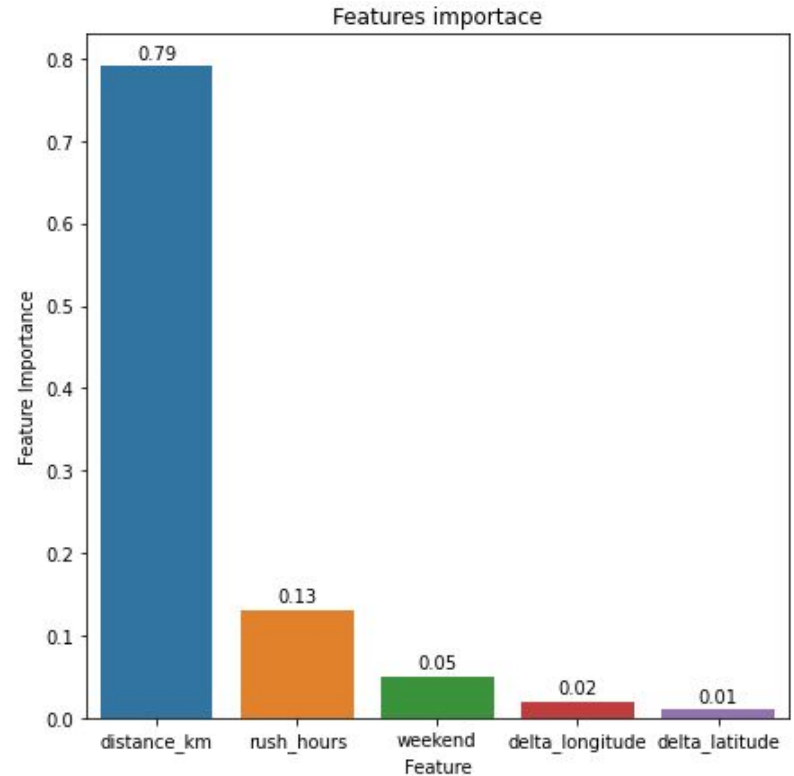# Gaussian Distribution - Trip duration



skewness  :  2.08118809062355

skewness  :  0.0012122340625846646

# Models

| Model Name | Evaluation Metrics | | |
|---|---|---|---|
| | R-Squared | RMSE | MAE |
| Linear Regression | 0.55 | 372 | 241 |
| Lasso Regression | 0.55 | 372 | 241 |
| Ridge Regression | 0.55 | 372 | 241 |
| Decision Tree | 0.61 | 346 | 230 |
| Random Forest | 0.6 | 348 | 232 |
| Xgboost | 0.62 | 341 | 224 |

# Parameters of importance

- Distance is most important parameter
- Rush hour is second most important parameter
- Delta latitude and Delta longitude is not that important as distance already captures the information from these parameters

# Conclusion

- Linear, Lasso, and Ridge regressions are giving similar results.

- In the Decision Tree regressor, the results are slightly improved.

- Out of all tried models, Random Xgboost Regressor is giving the best result.

# Future Challenges

- Distance and time of day are not sufficient to predict the time of ride it requires a thorough understanding of traffic and that the actual path of the ride is important.
- Deep Neural Networks can improve this performance.