# Capstone Project - Regression

## NYC Taxi Trip Time Prediction

Submitted by - Deepak Solanki

# Problem Statement

**Objective** - To 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.

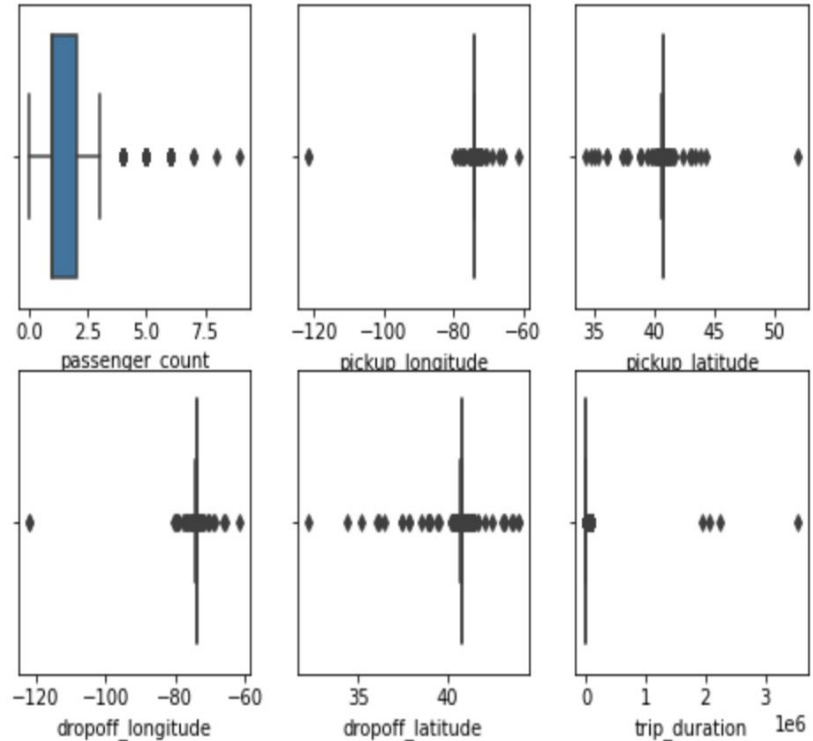**NYC Taxi Data.csv** - the dataset contains 1458644 trip records

.

# Data Fields

- **id** - a unique identifier for each trip
- **vendor_id** - a code indicating the provider associated with the trip record
- **pickup_datetime** - date and time when the meter was engaged
- **dropoff_datetime** - date and time when the meter was disengaged
- **passenger_count** - the number of passengers in the vehicle (driver entered value)
- **pickup_longitude** - the longitude where the meter was engaged
- **pickup_latitude** - the latitude where the meter was engaged
- **dropoff_longitude** - the longitude where the meter was disengaged
- **dropoff_latitude** - the latitude where the meter was disengaged
- **store_and_fwd_flag** - This flag indicates whether the trip record was held in vehicle memory before sending to the vendor because the vehicle did not have a connection to the server - Y=store and forward; N=not a store and forward trip
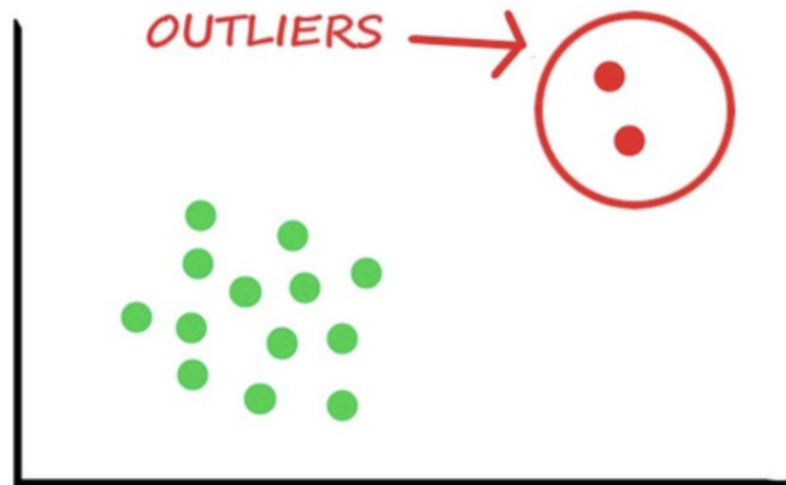- **trip_duration** - duration of the trip in seconds

# Outliers detection (Box Plots)

- Box plots shows the presence of outliers in a couple of features

# Missing values checking & Outliers removal

- We checked and found that there were no missing values in the whole dataset.
- We removed outliers from following features using various methods like predefined threshold and interquartile range.
  - passenger_count
  - pickup_longitude
  - pickup_latitude
  - dropoff_longitude
  - dropoff_latitude
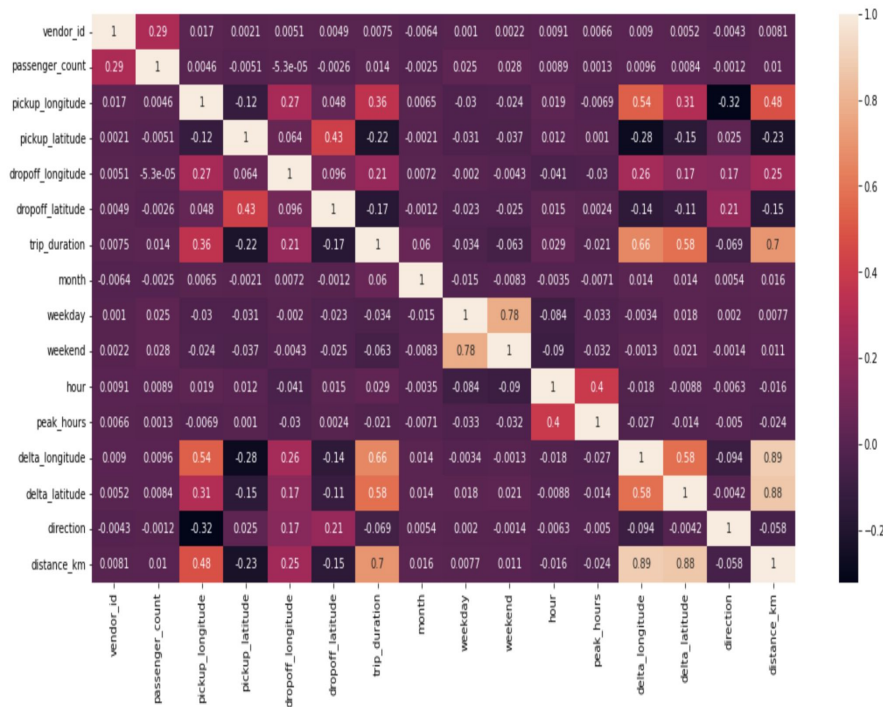  - trip_duration
  - distance_km

# Feature Engineering

- Features extractions from "Pickup Datetime"
  - Month
  - Weekday
  - Weekend
  - Hour
  - Peak_hours
- New Derived features from Latitudes and Longitudes
  - Delta longitude / latitude
  - Direction (in degrees)
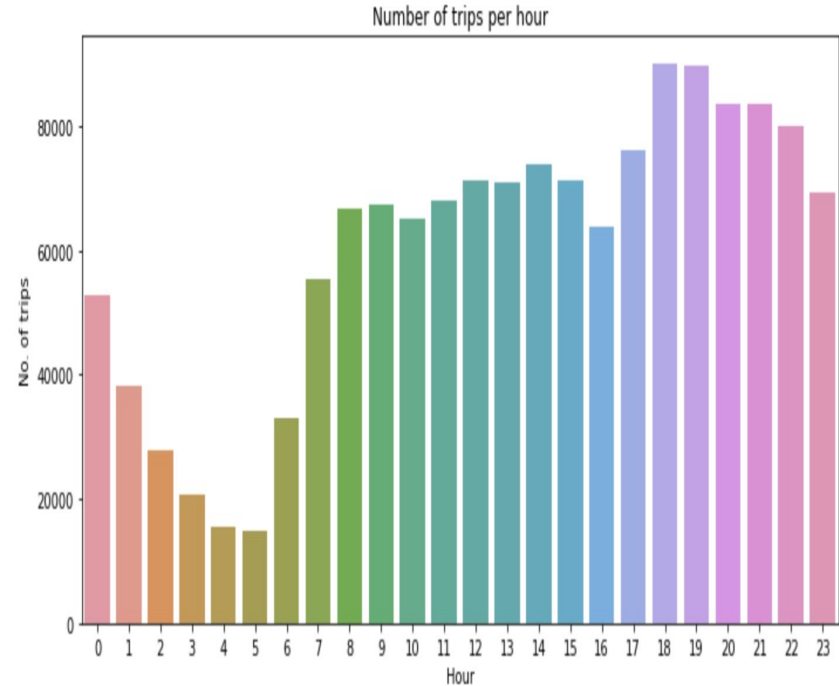  - Distance_km (haversine distance)

# EDA (Feature correlation matrix)

- There are few couplets of variables that are significantly correlated.
- For ex.
  - Corr(distance_km, trip_duration) = 0.7
  - Corr(delta_longitude, trip_duration) = 0.66
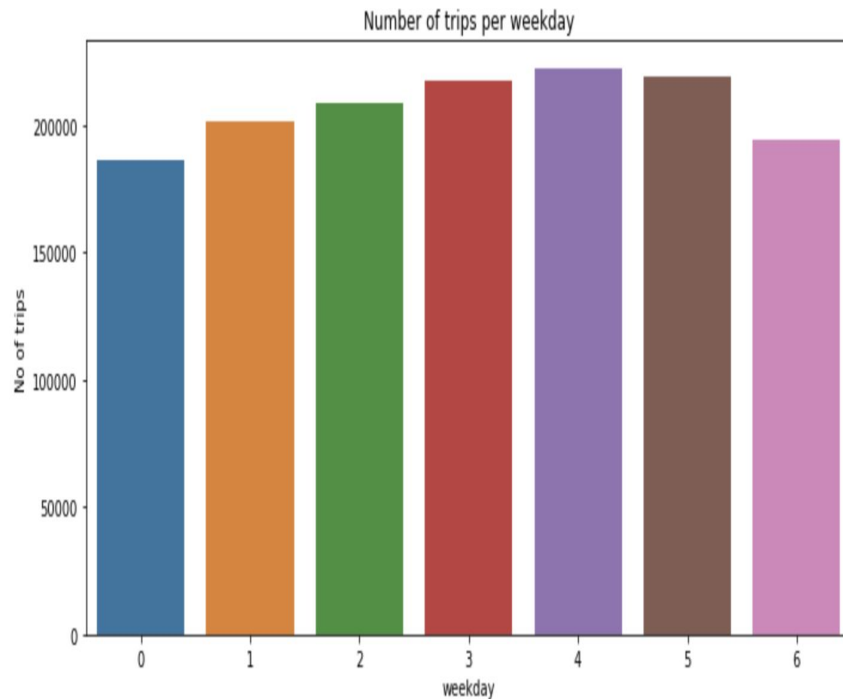  - Corr(delta_longitude, pickup_longitude) = 0.54

# EDA (Continue)

- Maximum no. of trips occurs in the evening time slot **(18:00-20:00 hours)**. Then it has an decreasing trend till 5:00. After that, It start increasing till evening.
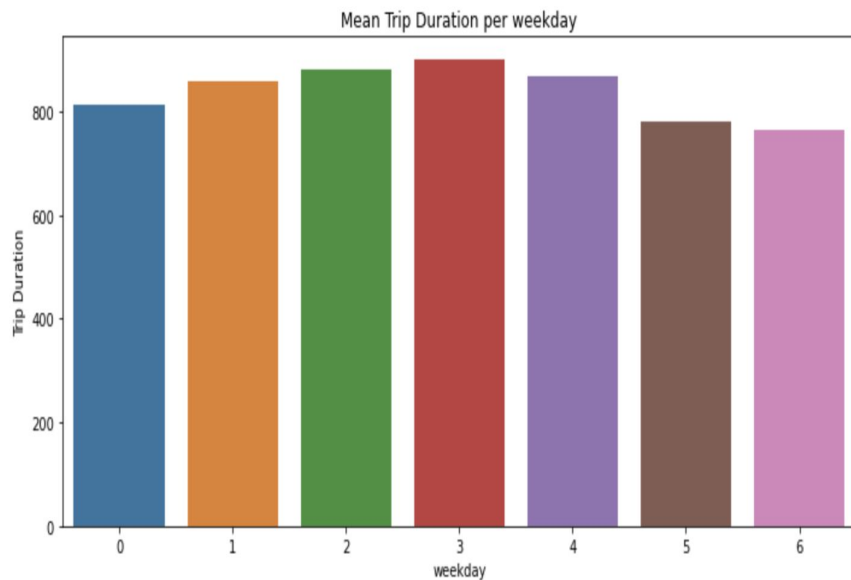


Number of trips per hour

# EDA (Continue)

- Number of trips has an increasing trend from Monday to Friday. And it's relatively low on Sunday as compared to other days.



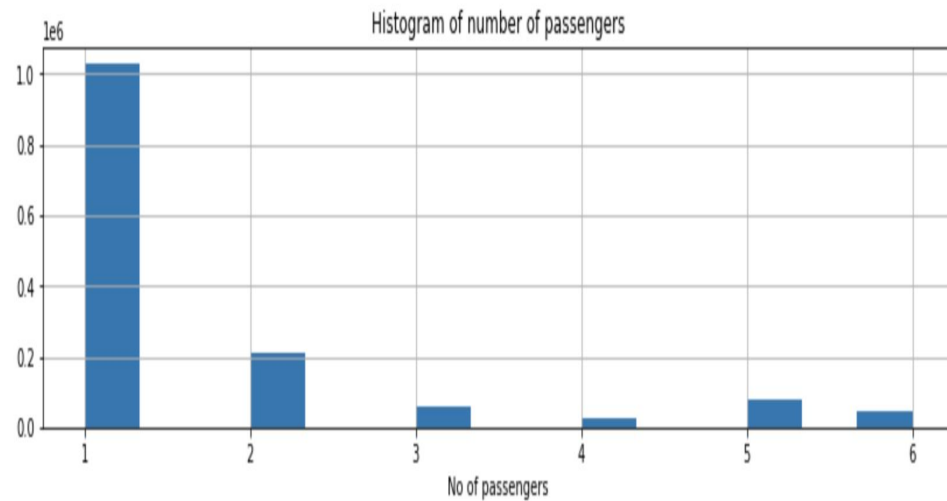Number of trips per weekday

# EDA (Continue)

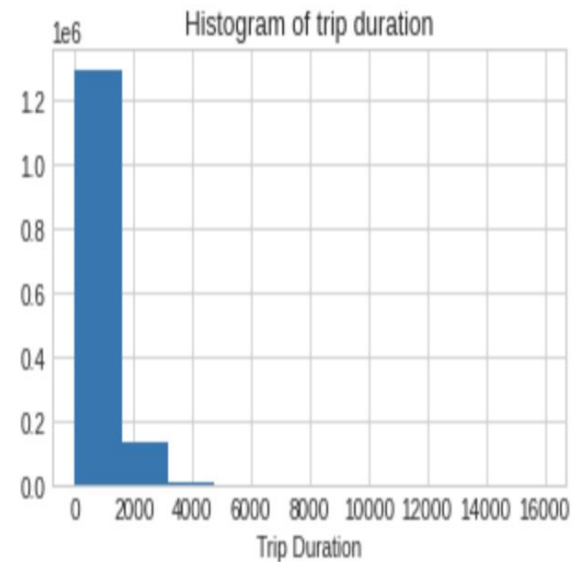- Mean trip duration on weekdays is higher than that of on weekends. It is highest on Thursday (weekday = 3).



Mean Trip Duration per weekday

# EDA (Continue)

- No. of passengers ranges between 1 and 6.
- **71%** of the trips are having only 1 passenger.
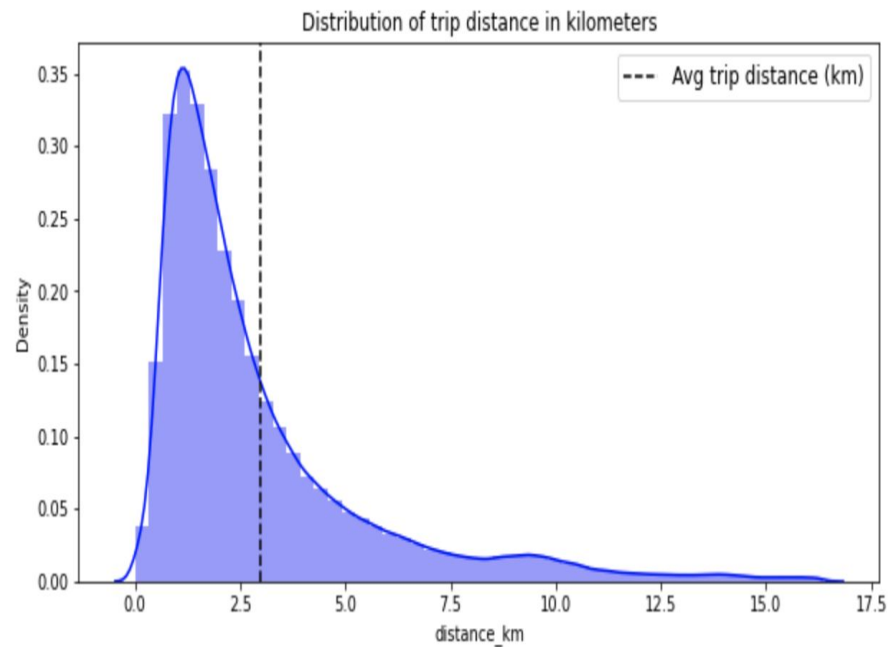


Histogram of number of passengers

# EDA (Continue)

- Majority of the trip durations are below 1600 seconds. Also, a few of trips are above 3000 seconds.

# EDA (Continue)

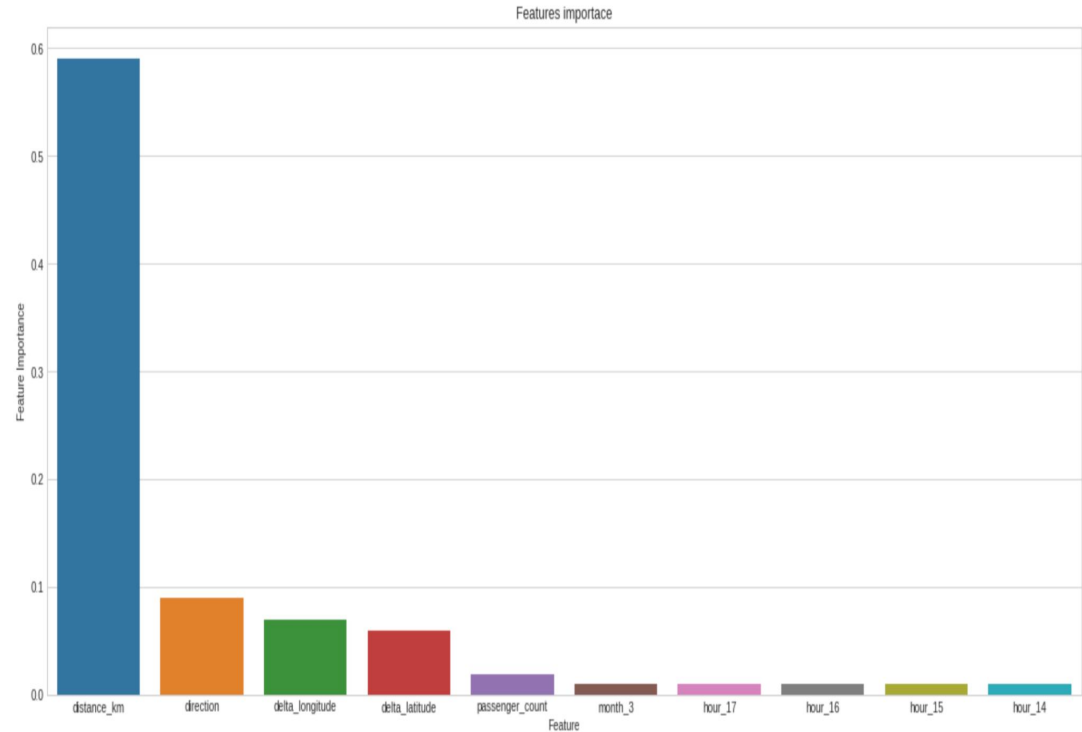- The distribution of **distance_km** is right skewed. And avg distance is ~3 km.



Distribution of trip distance in kilometers

# Modeling

| Model Name | Evaluation Metrics | | | |
|---|---|---|---|---|
| | R-Squared | MSE | RMSE | MAE |
| Baseline | - | - | - | 413.33 |
| Linear Regression | 0.57 | 133915.42 | 365.94 | 248.37 |
| Lasso Regression | 0.57 | 133914.8 | 365.94 | 248.38 |
| Ridge Regression | 0.57 | 133914.77 | 365.94 | 248.38 |
| Decision Tree regressor | 0.56 | 138504.49 | 372.16 | 250.61 |
| Random Forest Regressor | 0.66 | 107773.48 | 328.29 | 213.14 |

# Random Forest Regressor (Best Model)

- We used **GridSearchCV** cross validation technique for Hyperparameter tuning.
- param_grid
  **{'n_estimators':** [10, 25],
   **'max_features'**: [5, 10],
   **'max_depth'**: [10, 50, None],
   **'bootstrap'**: [True, False]}
- Random Forest **best parameters** are
  **{'bootstrap'**: True,
  **'max_depth'**: 50,
  **'max_features'**: 10,
  **'n_estimators'**: 25}

**AI**

# Random Forest Top 10 important features

- **distance_km** is the most important feature contributing around 58%.
- Other important features are direction (9%), delta_longitude (7%) and delta_latitude (6%).



Features importace

# Conclusions

1. Linear, Lasso and Ridge regressions are giving the similar results. And They are performing better than the baseline model.
2. In Decision Tree regressor, the results are slightly poorer than that of linear models.
3. Out of all tried models, Random Forest is giving the best result.

# Challenges Faced

1. Deriving new features from geospatial data (using given Latitude and Longitude coordinates).
2. Extracting useful features from pickup datetime column.
3. High computation time in hyperparameter tuning.

# Further scope for improvement

- **Feature Engineering** - More features can be derived from the Geospatial data which can improve the models performance.
- Neural Network based regression model can be built to improve the results further.

# Thank You