# EDA Capstone

## Capstone Project

# NYC Taxi Trip Time Prediction

Vishakha Kumari

Shivam Kumar

Abhinav Akotkar

Durgesh Shukla

AI

# Content

- Introduction
- Data Description and cleaning
- New FeaturesCreating
- Univariate Analysis
- Bivariate Analysis
- Location Visualization On Map
- Correlation Heatmap
- Linear Regression
- Lasso Regression
- XGBoost Model
- LightGBM Model
- Model Summary
- Conclusion

# Introduction

- Our task is to build a model that predicts the total ride duration of taxi trips in New York City.
- Primary dataset contains NYC Taxi and Limousine Commission, which includes pickup time, geo-coordinates, number of passengers, and several other variables.

- 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, you should predict the duration of each trip in the test set.

# Data Description and cleaning

**SHAPE OF DATA**

Following is the shape of data

Rows: 1458644

Columns: 11

Null values - 0 (There are no NaN or null values in our dataset)

Data contains the training set (contains 1458644 trip records)

# Data Description and cleaning- Column description

- 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

# New FeaturesCreating

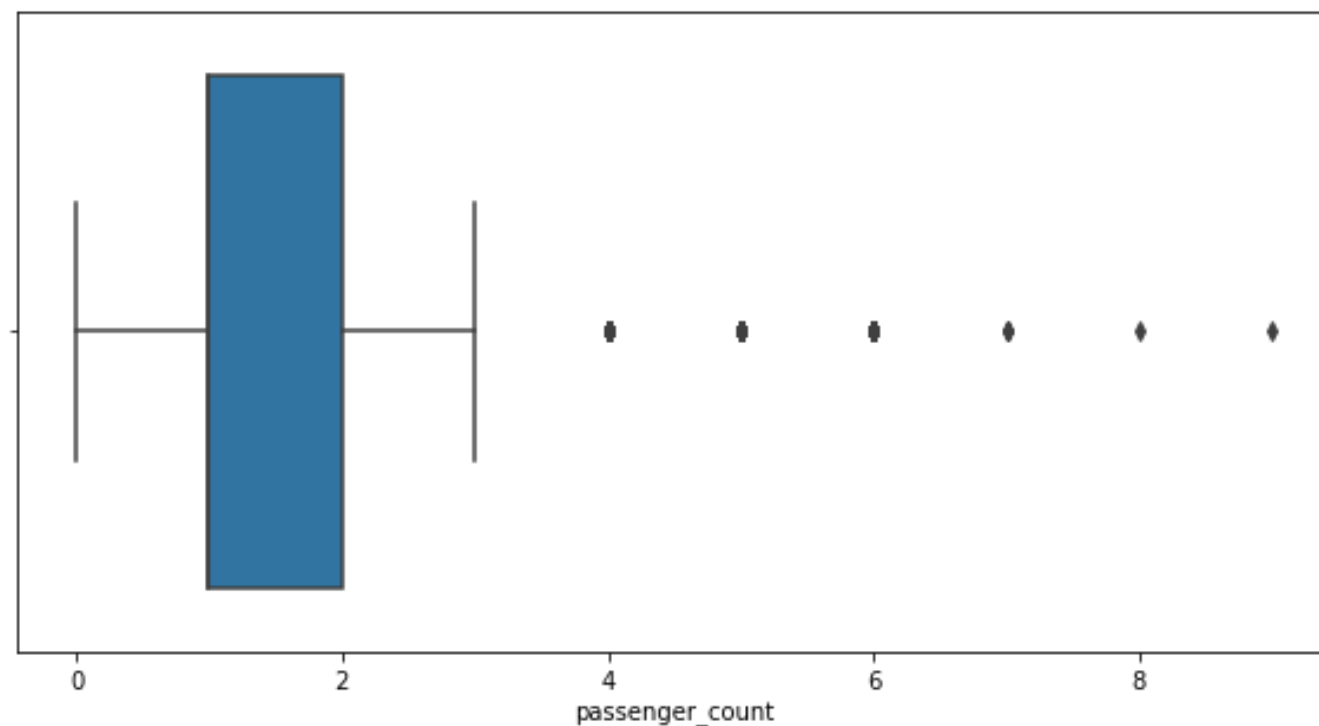Here we added some new features for our ease of processing.
- pickup_weekday: This will contain the day on which ride was taken.
- dropoff_weekday: This will contain the day on which ride was ended.
- pickup_hour: time at which ride was started.
- dropoff_hour: time at which ride was stoped.
- month: month of the ride
- pickup_weekday_num: number of the day on which ride was taken, taking Monday as 0.
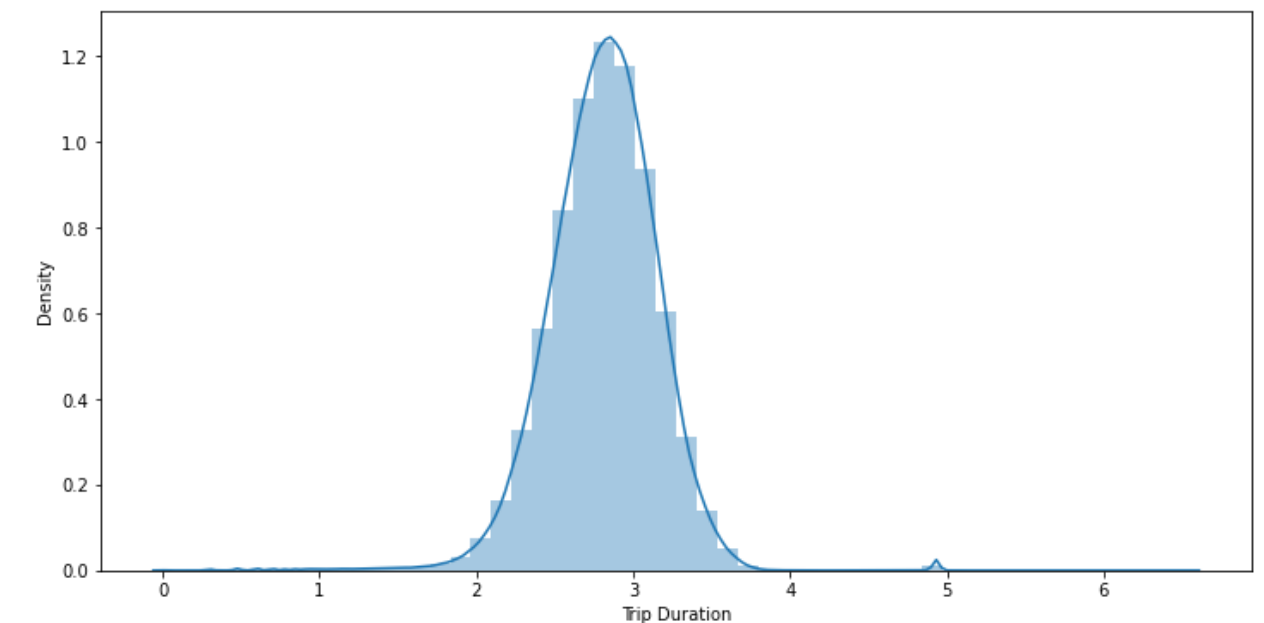
# Univariate Analysis

## Number of Passengers Per Cab

to check whether there are any outliers in our dataset



**From the above Box Plot we can clearly notice there are some outliers in our dataset. It shows number of passengers per taxi is more than 7.**

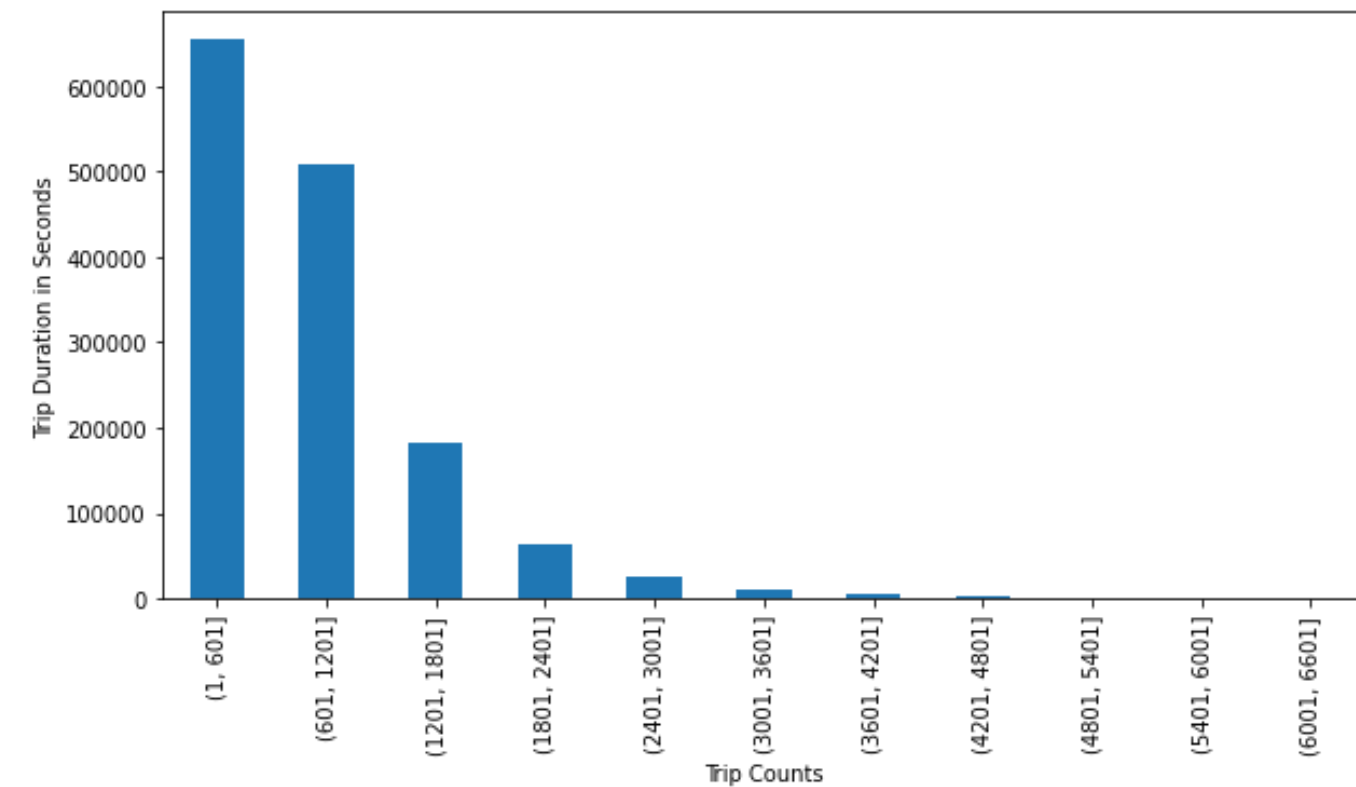**In most of the trips passenger number is between 1 or 2**

## Trip Duration



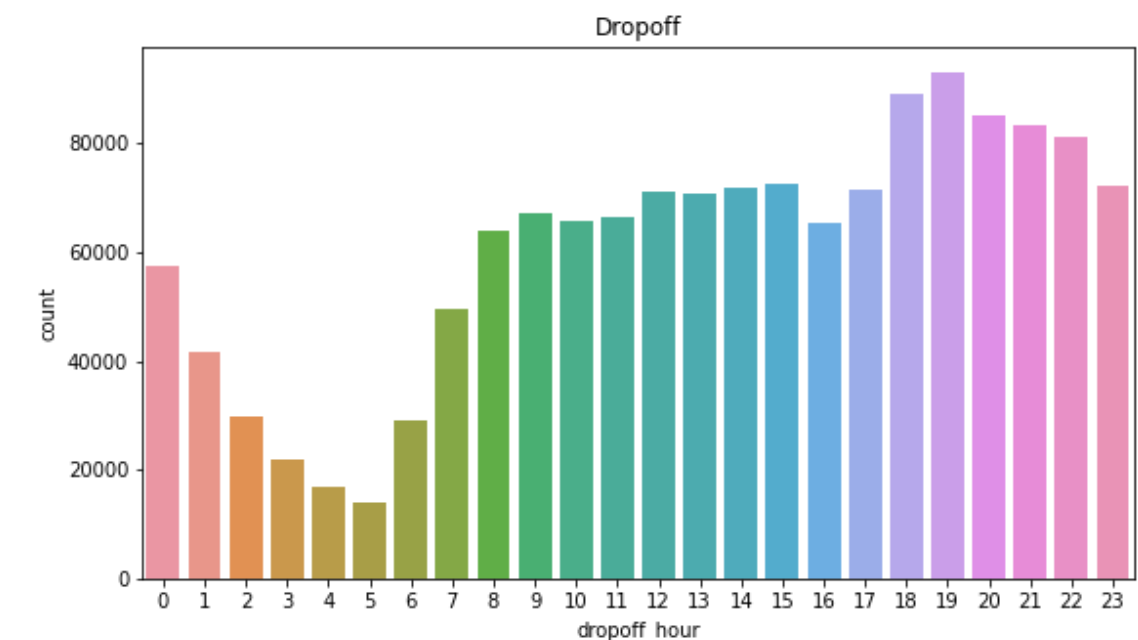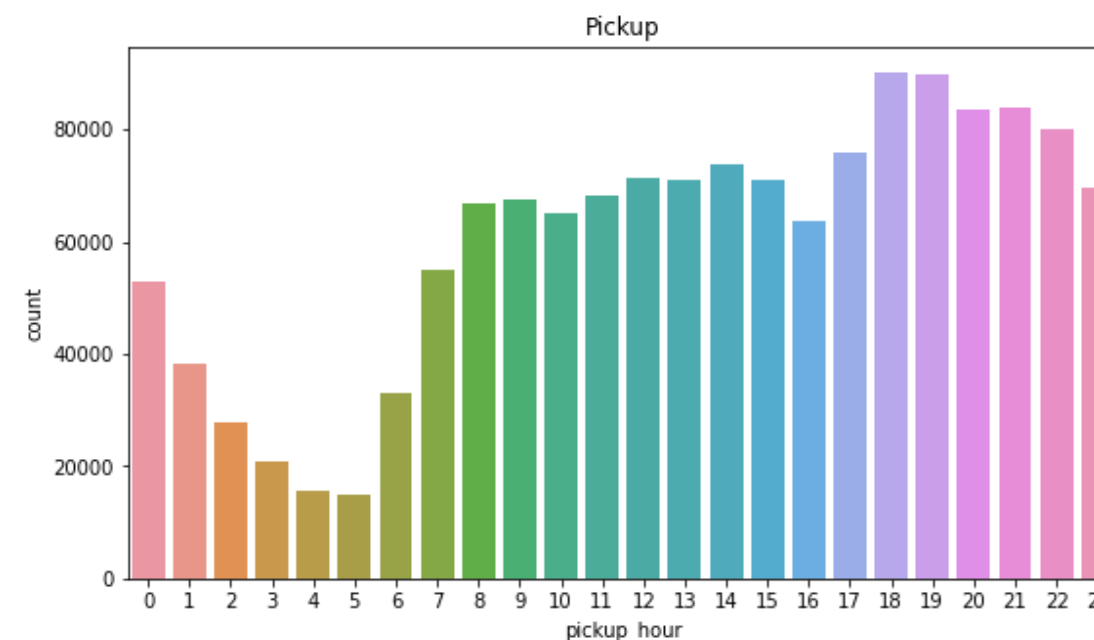**We can spot an outlier near 5 in the above plot.**

# Univariate Analysis

## Number of Trips Taken per Minute



**Most of the trips took around 1800 seconds**
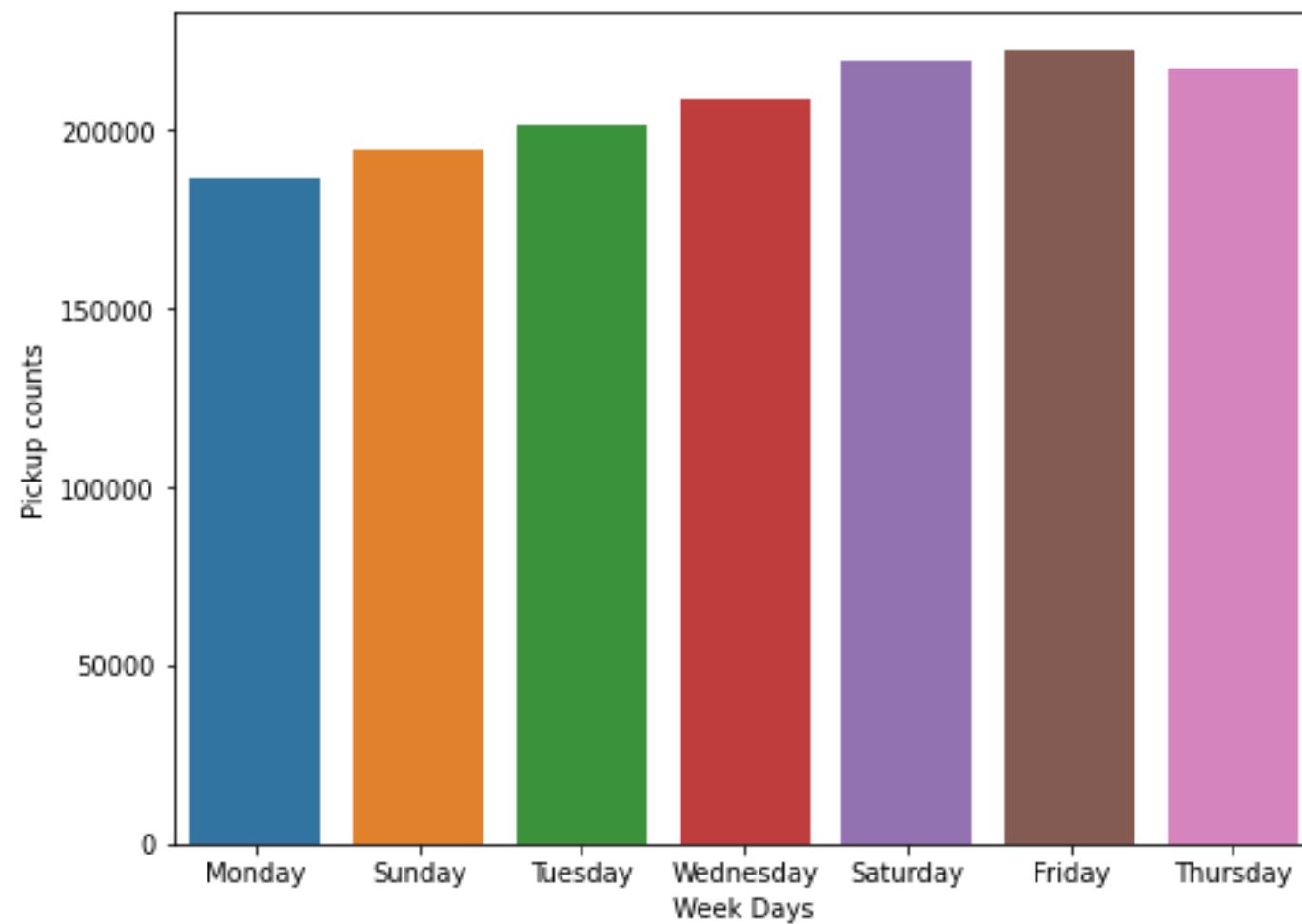
## Total Trips per Hour



**From the above two plots we can notice that most of the pickups and drops are between 6 PM to 7 PM. least number of pickups are between 4 AM to 5 AM.**
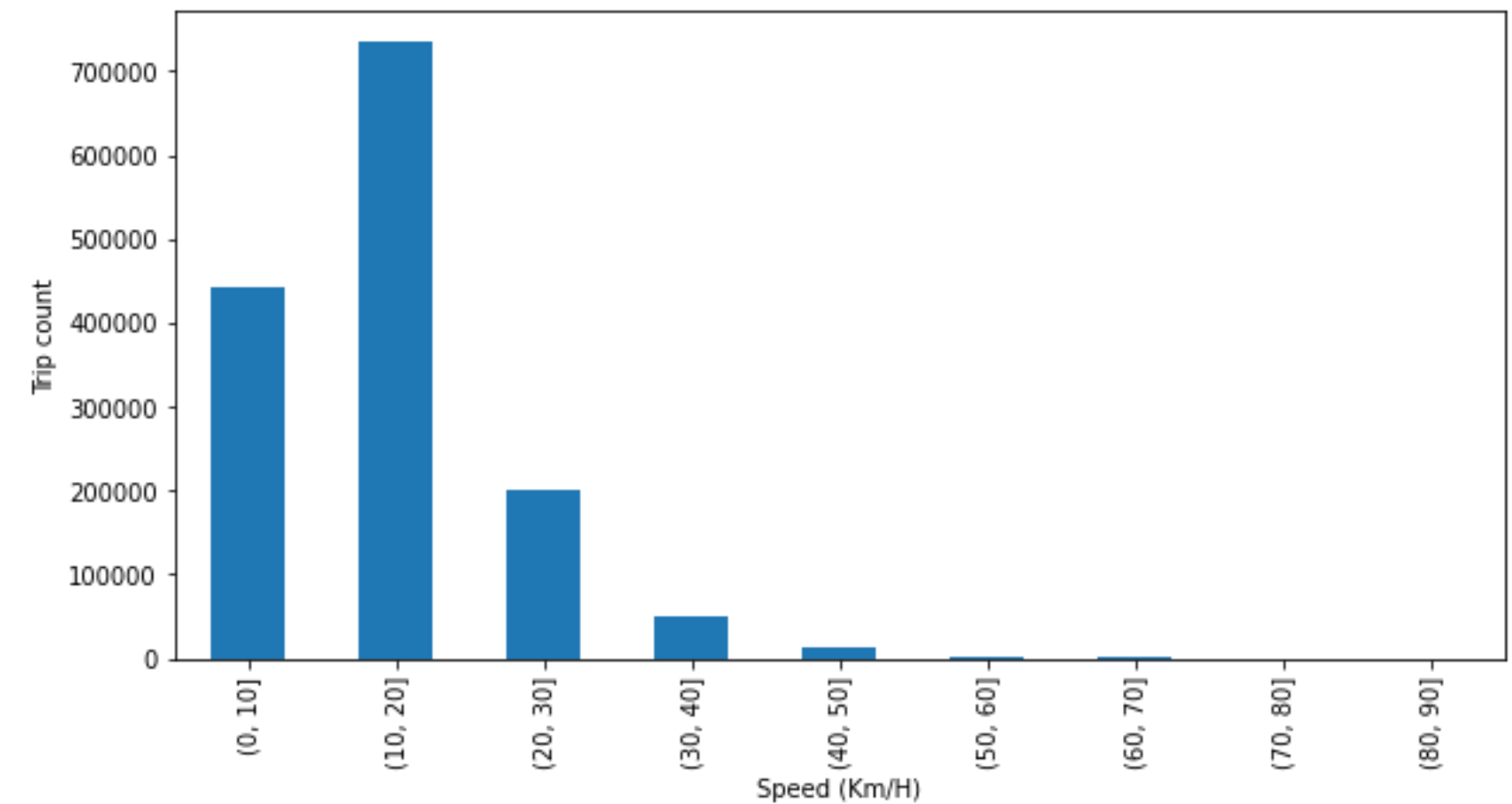
# Univariate Analysis

## Trips Per Day

sns.countplot(data.pickup_weekday)



We can see the number of pickups starts increasing from monday to friday, then starts decreasing on weekends

## Speed

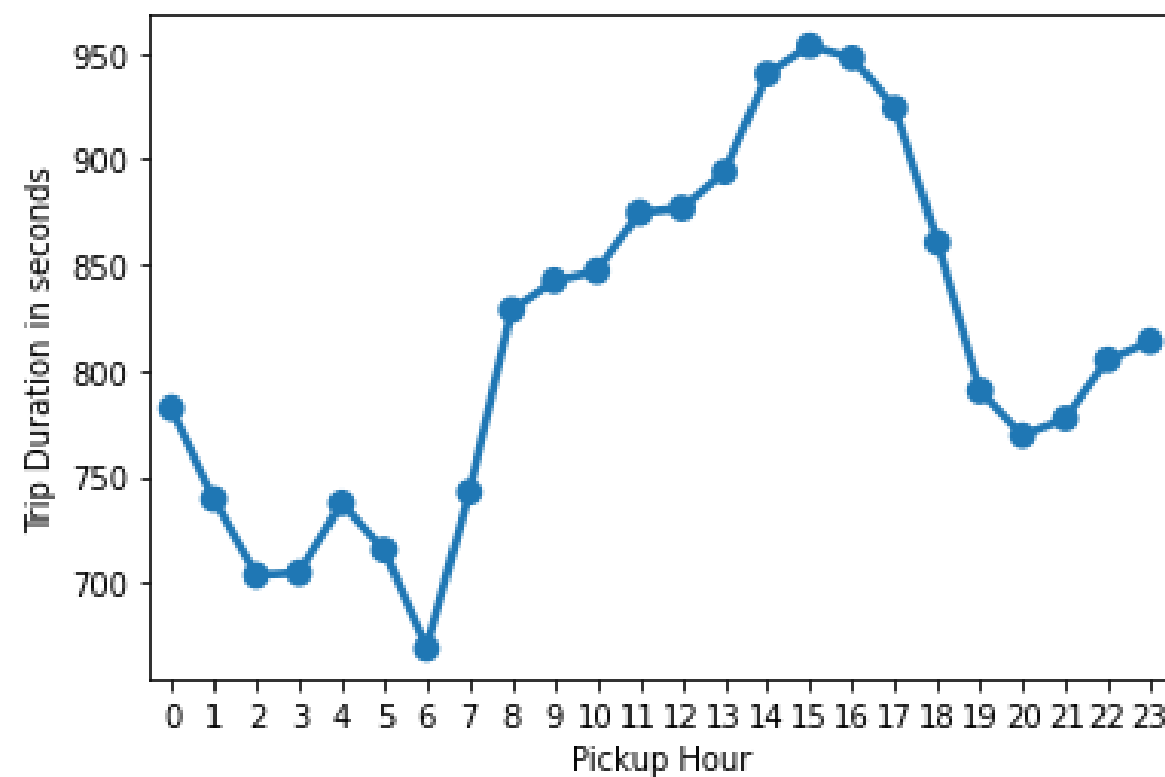data.speed.groupby(pd.cut(data.speed, np.arange(0,100,10))).count().plot(kind = 'bar')



Most of the trips are done at the range 10-20 Km/Hr
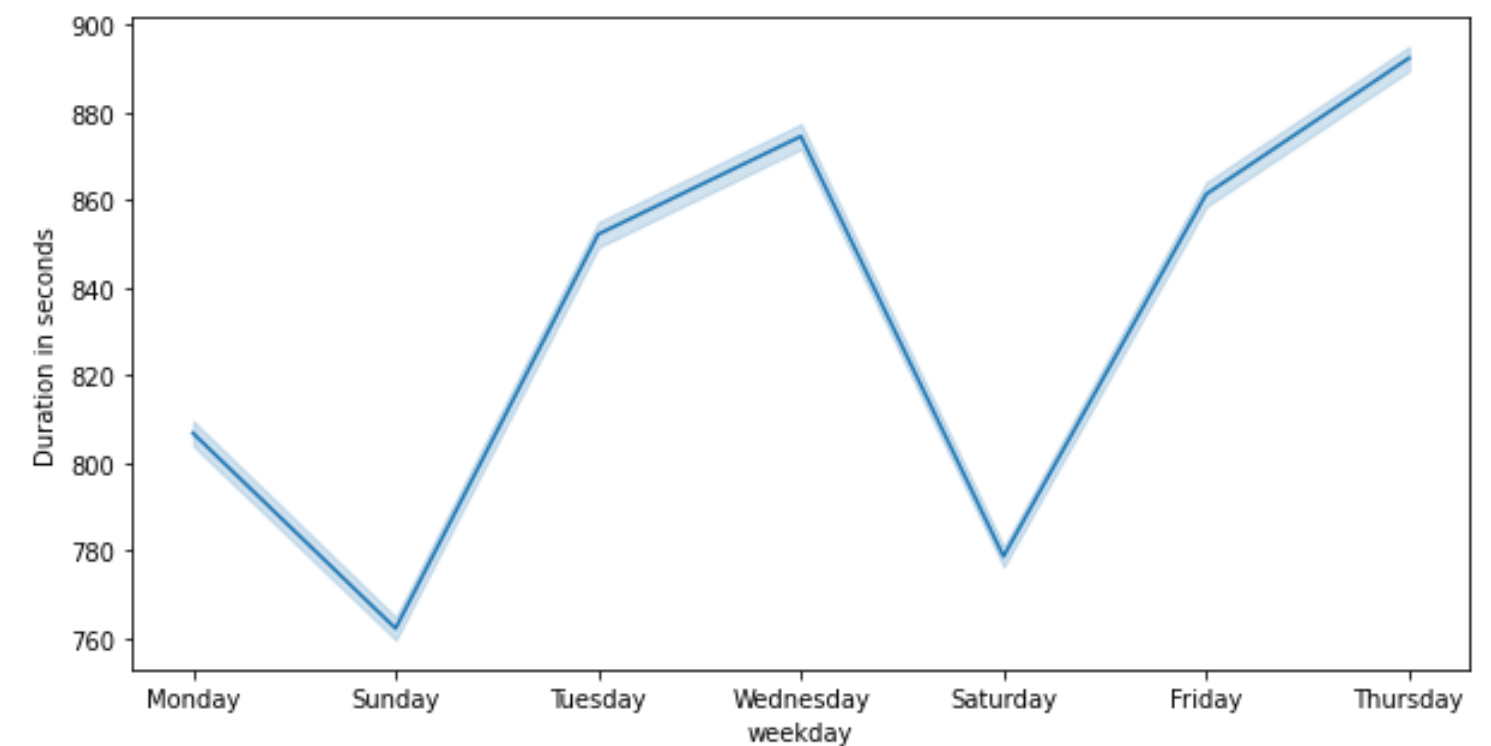
# Bivariate Analysis

## Trip Duaration Per Hour

duration = data.groupby('pickup_hour').trip_duration.mean()
sns.pointplot(duration.index, duration.values)



- From above plot we can see lowest mean trip duration is at 6 AM.
- Highest mean trip duration is at 3 PM. Most probaly due to high traffic

## Trip Duration per Weekday

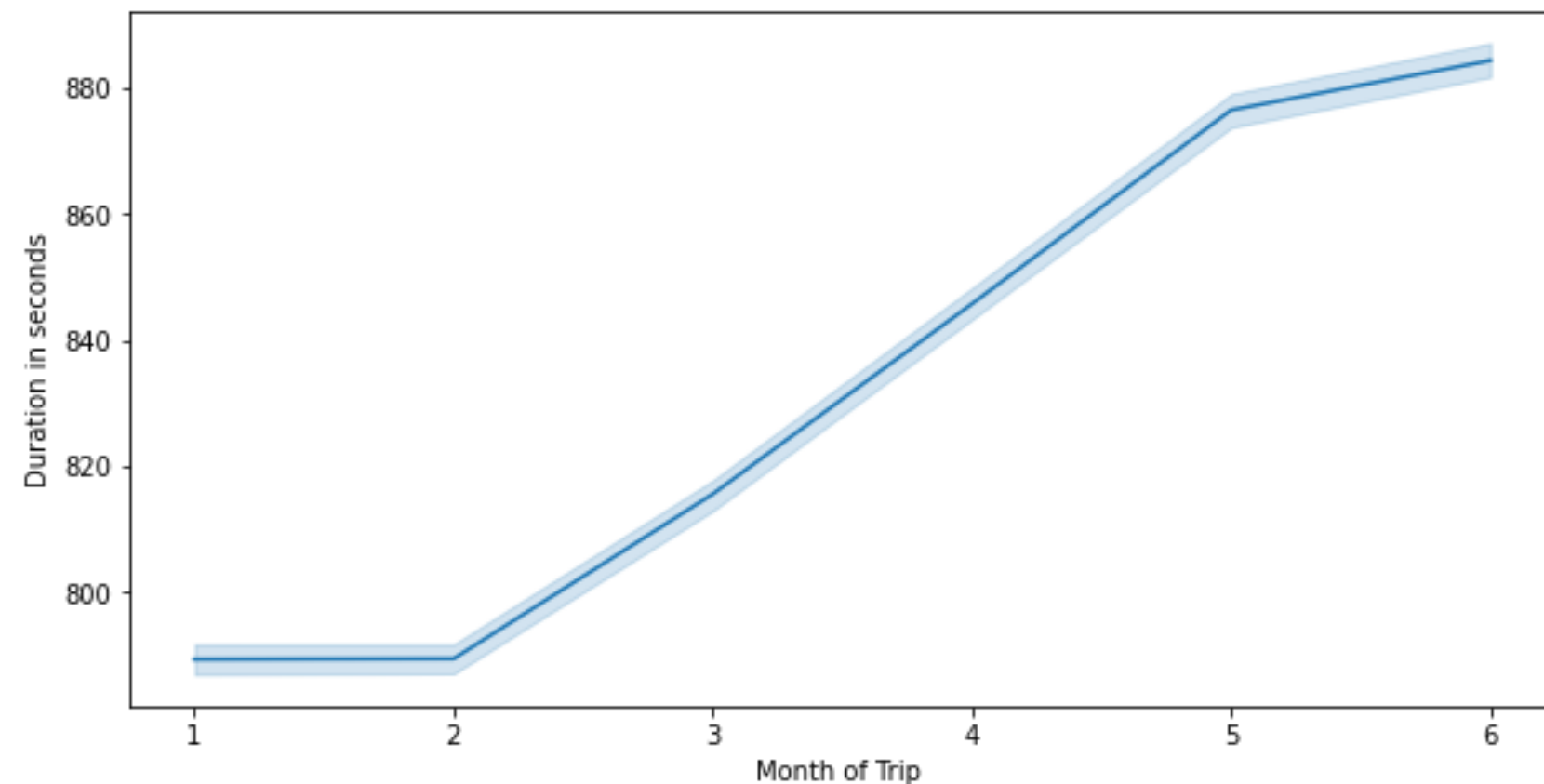sns.lineplot(x='pickup_weekday',y='trip_duration',data=data)



- Trip duration is highest on thursday
- Trip duration is lowest on weekends, because of holiday.
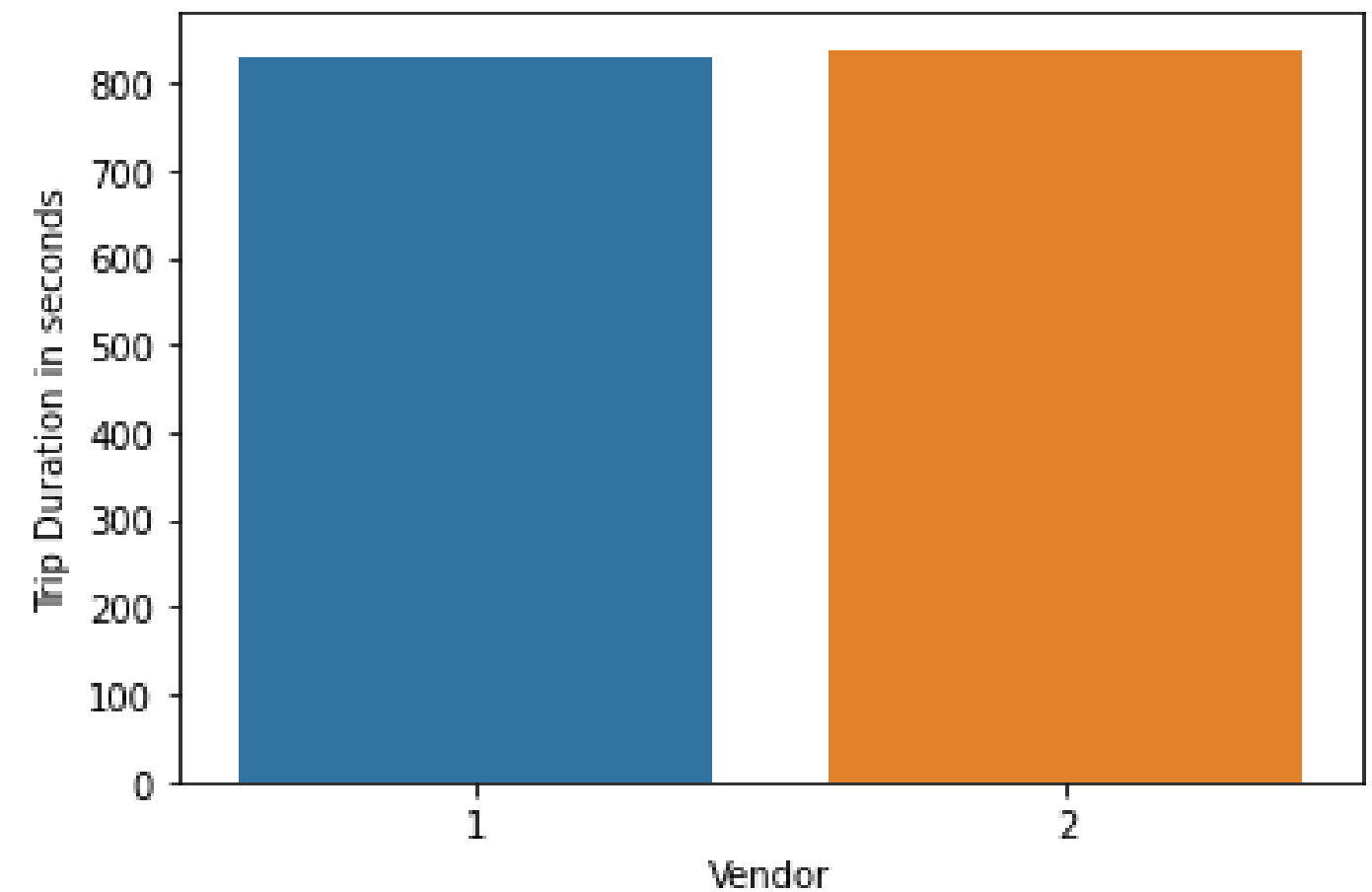
# Bivariate Analysis

## Trip Duaration Per Month

**sns.lineplot(x='month',y='trip_duration', data=data)**



## Trip Duration per Vendor

**sns.lineplot(x='pickup_weekday',y='trip_duration',data=data)**



- There is an increasing trend of trip duration as month increases.
- This might be due to some climatic conditions like rain and snow.
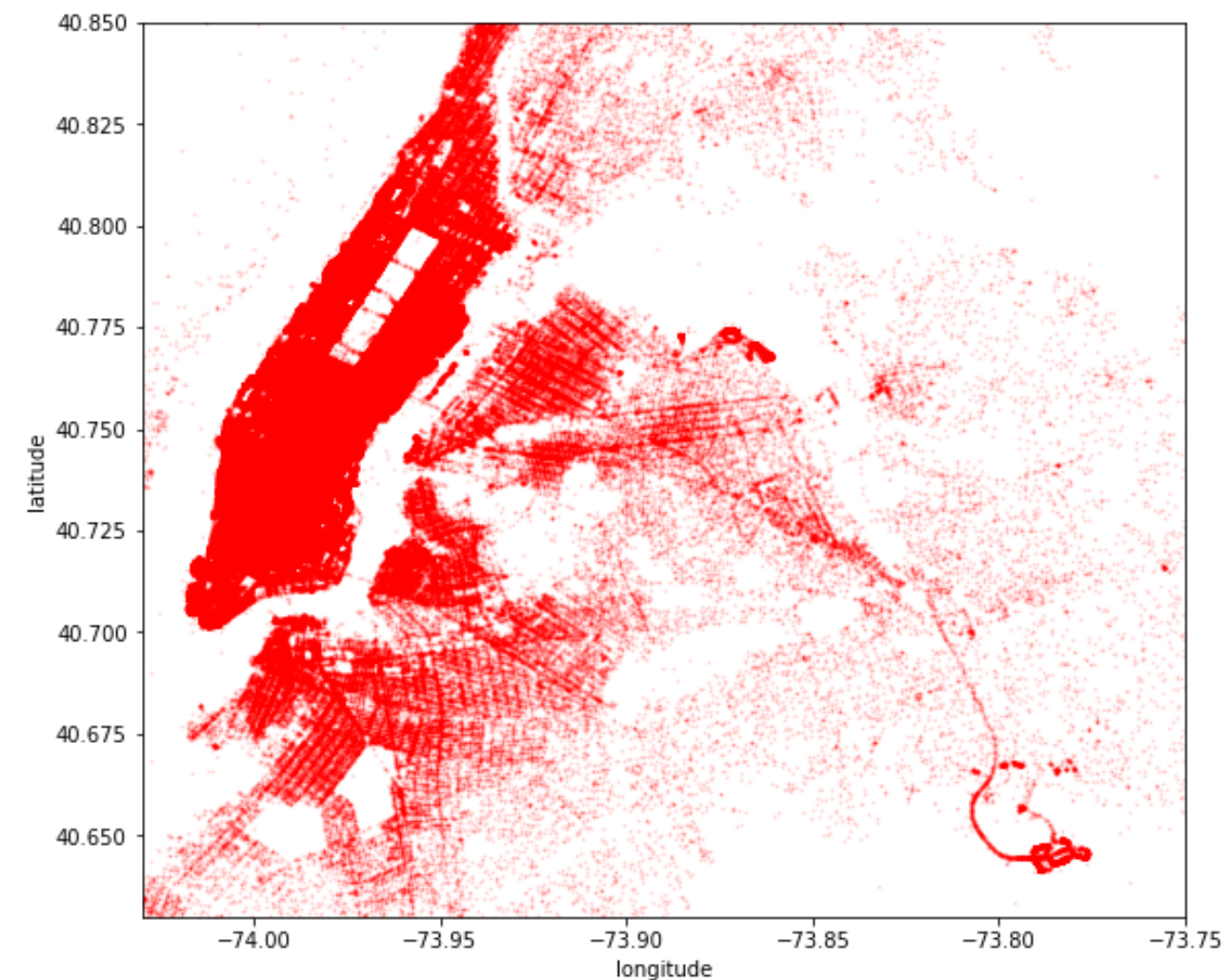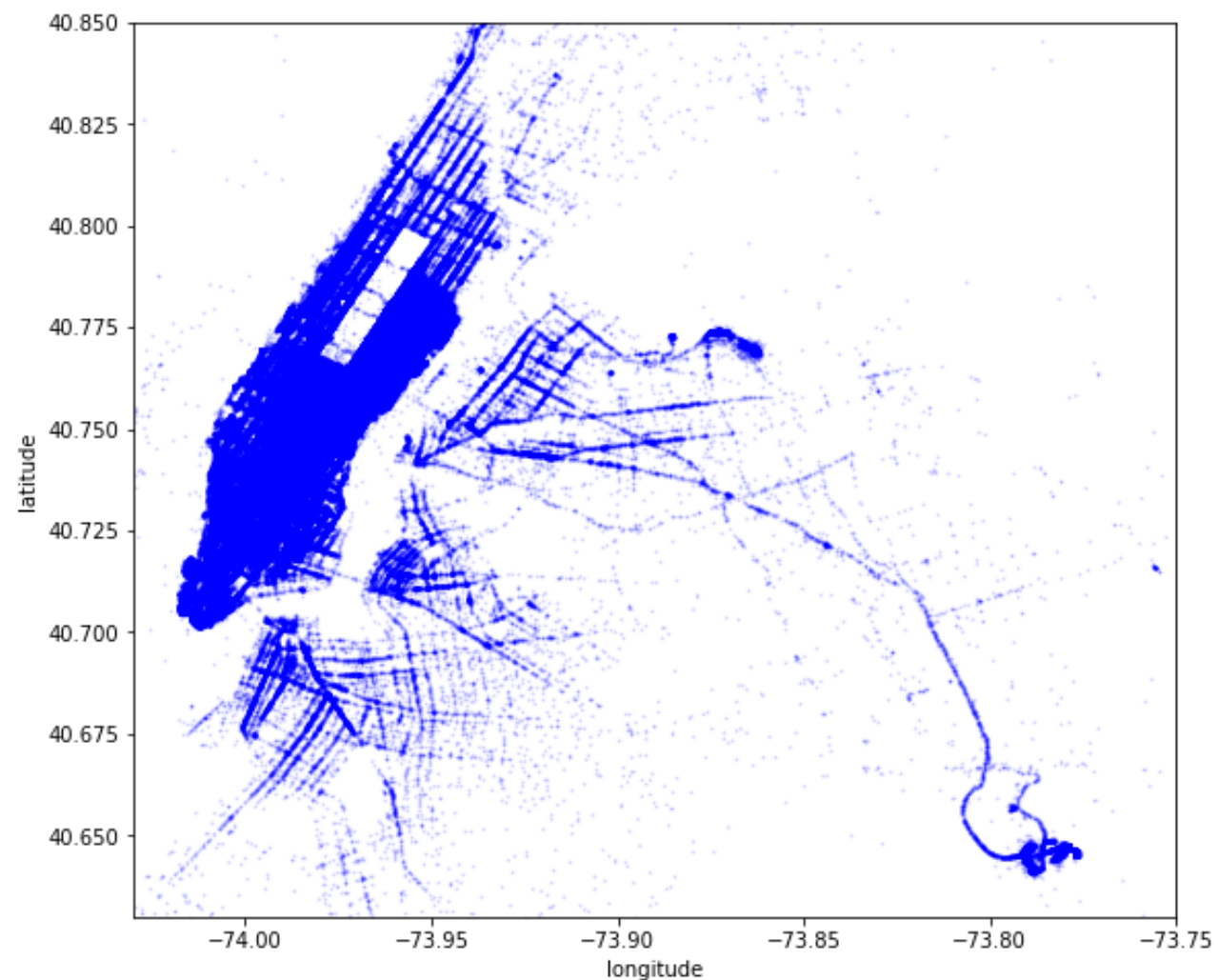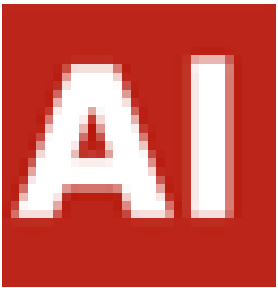- As in rainy season roads are more blocked and trip duration will increase

- Trip duration for vendor 2 is higher than vendor 1

# Location Visualization On Map

**Plotting pickup and dropoff location according to their lattitude and longitude. Here we can see the map of the streets of NYC very clearly.**
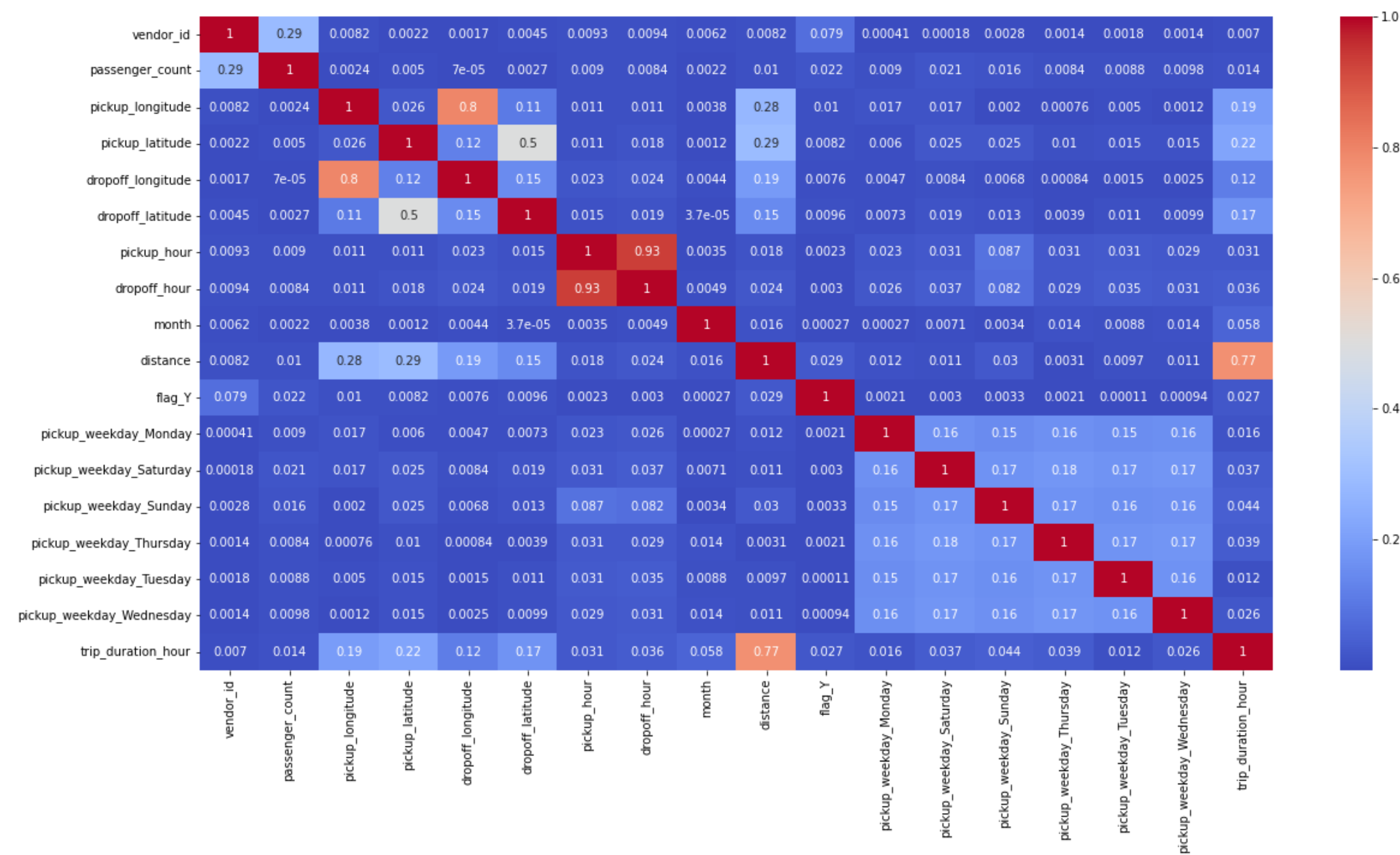
plt.scatter(data['pickup_longitude'].values, data['pickup_latitude'].values, color='blue', s=1, alpha=0.1)

# Correlation Heatmap

**sns.heatmap(abs(corelation), annot=True, cmap='coolwarm')**



we can see that there isn't much relation between the variables except distance-trip duration and dropoff longitude and pickup longitude.
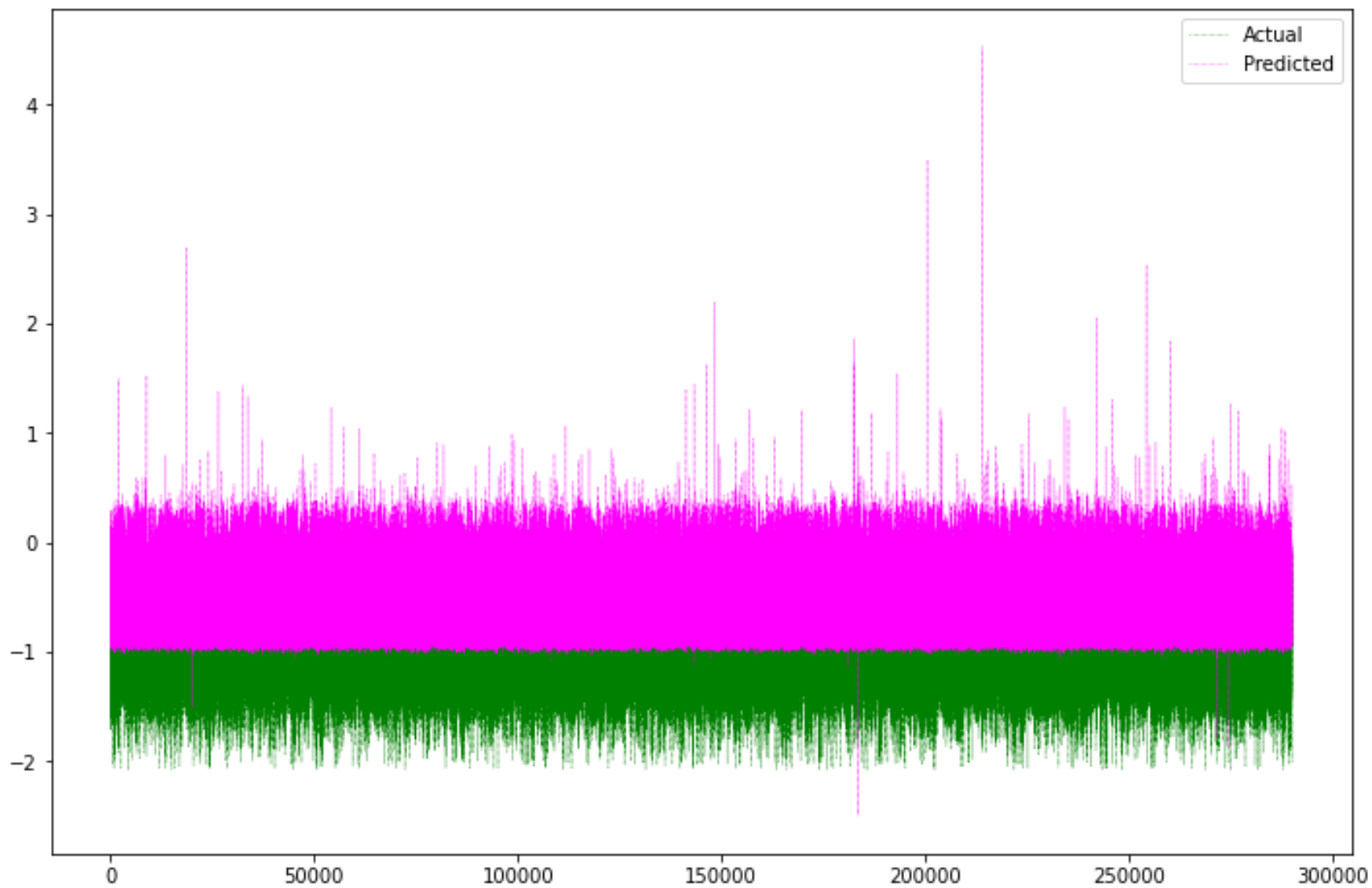
# Linear Regression

Model Evaluation

Visualization of Actual vs Predicted Data

```
plt.figure(figsize= (12,8))
c= [i for i in range(0, len(y_test))]
plt.plot(c, y_test, color='green', linewidth=0.5, linestyle=':')
plt.plot(c, y_pred_test, color='magenta', linewidth=0.5, linestyle=':')
plt.title('Actual vs Predicted Test Data', fontsize=20)
plt.legend(["Actual", "Predicted"])
```
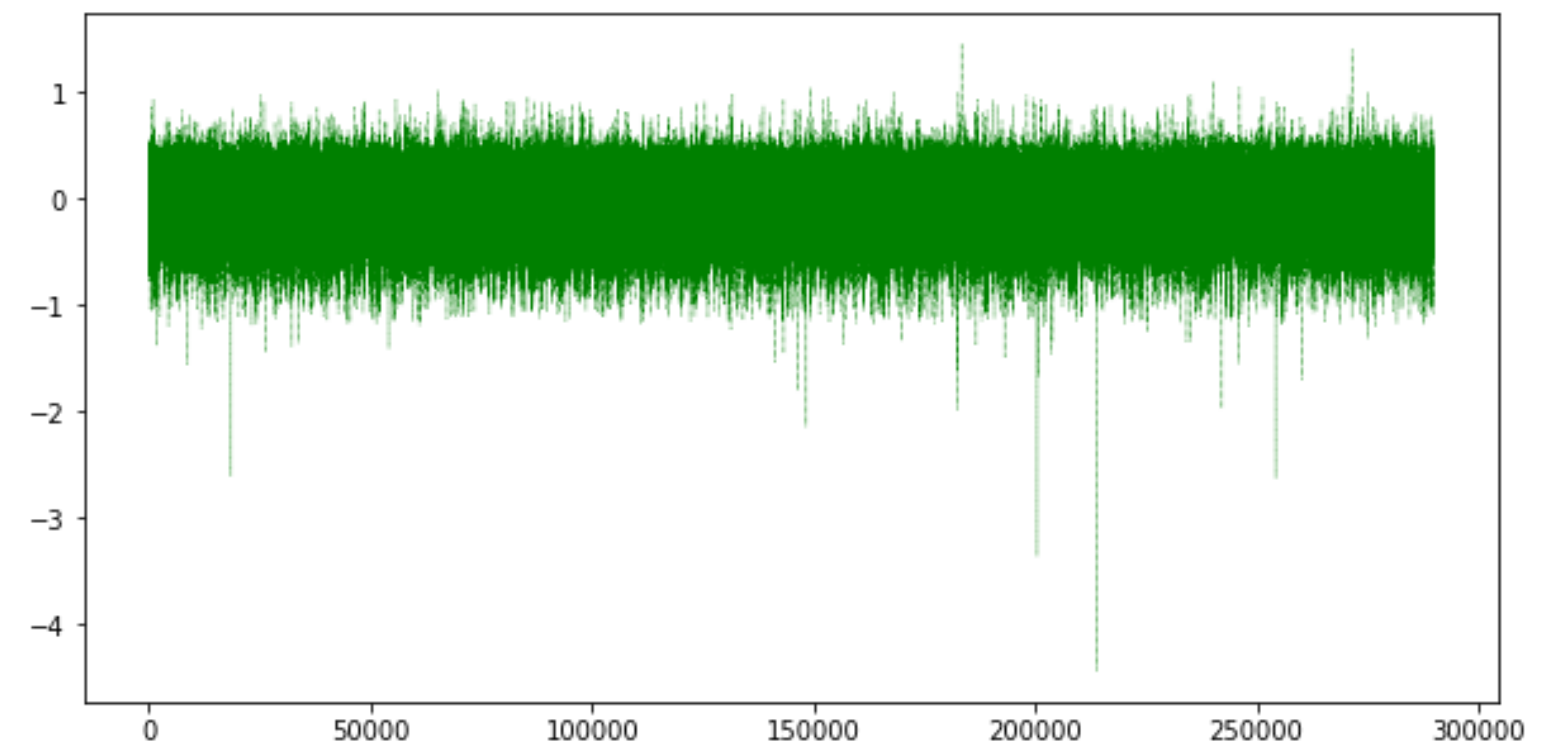
Ploting the error

```
plt.figure(figsize= (10,5))
c= [i for i in range(0, len(y_test))]
plt.plot(c, y_test-y_pred_test, color='green', linewidth=0.5, linestyle=':')
plt.title('Error', fontsize=20)
plt.show()
```

# Lasso Regression
## Model Evaluation

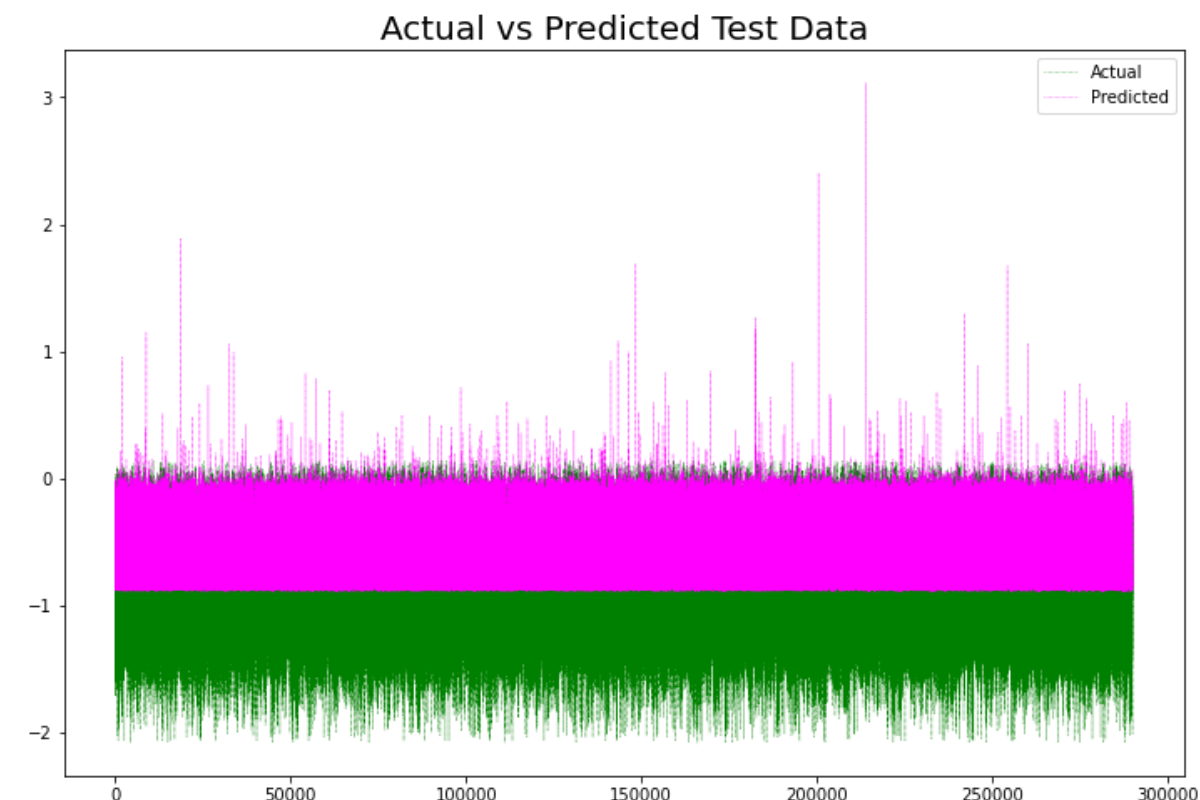**MSE RMSE R2 ADJUSTED R2 SCORE FOR TRAINING DATASET**

Train MSE : 0.06132640952105937
Train RMSE : 0.24764169584514512
Train R2 : 0.40403416099387746
Train Adjusted R2 :  0.40402645457228004

**MSE RMSE R2 ADJUSTED R2 SCORE FOR Test DATASET**

Test MSE : 0.06128689967297325
Test RMSE : 0.24756191078793452
Test R2 : 0.40381719983745323
Test Adjusted R2 :  0.40378636173268256

## Ploting the Actual vs Predicted data

# XGBoost Model
## Model Evaluation

**MSE RMSE R2 ADJUSTED R2 SCORE FOR TRAINING DATASET**

Train MSE : 0.017228304970182916
Train RMSE : 0.1312566378137994
Train R2 : 0.8325765146468818
Train Adjusted R2 : 0.8325743496973511

**MSE RMSE R2 ADJUSTED R2 SCORE FOR Test DATASET**

Test MSE : 0.01970132651011081
Test RMSE : 0.1403614138932449
Test R2 : 0.8083506904674747
Test Adjusted R2 : 0.8083407772303183

## Ploting graph for Actual vs Predicted data



Actual vs Predicted Test Data



Error

# LightGBM Model
## Model Evaluation

**MSE RMSE R2 ADJUSTED R2
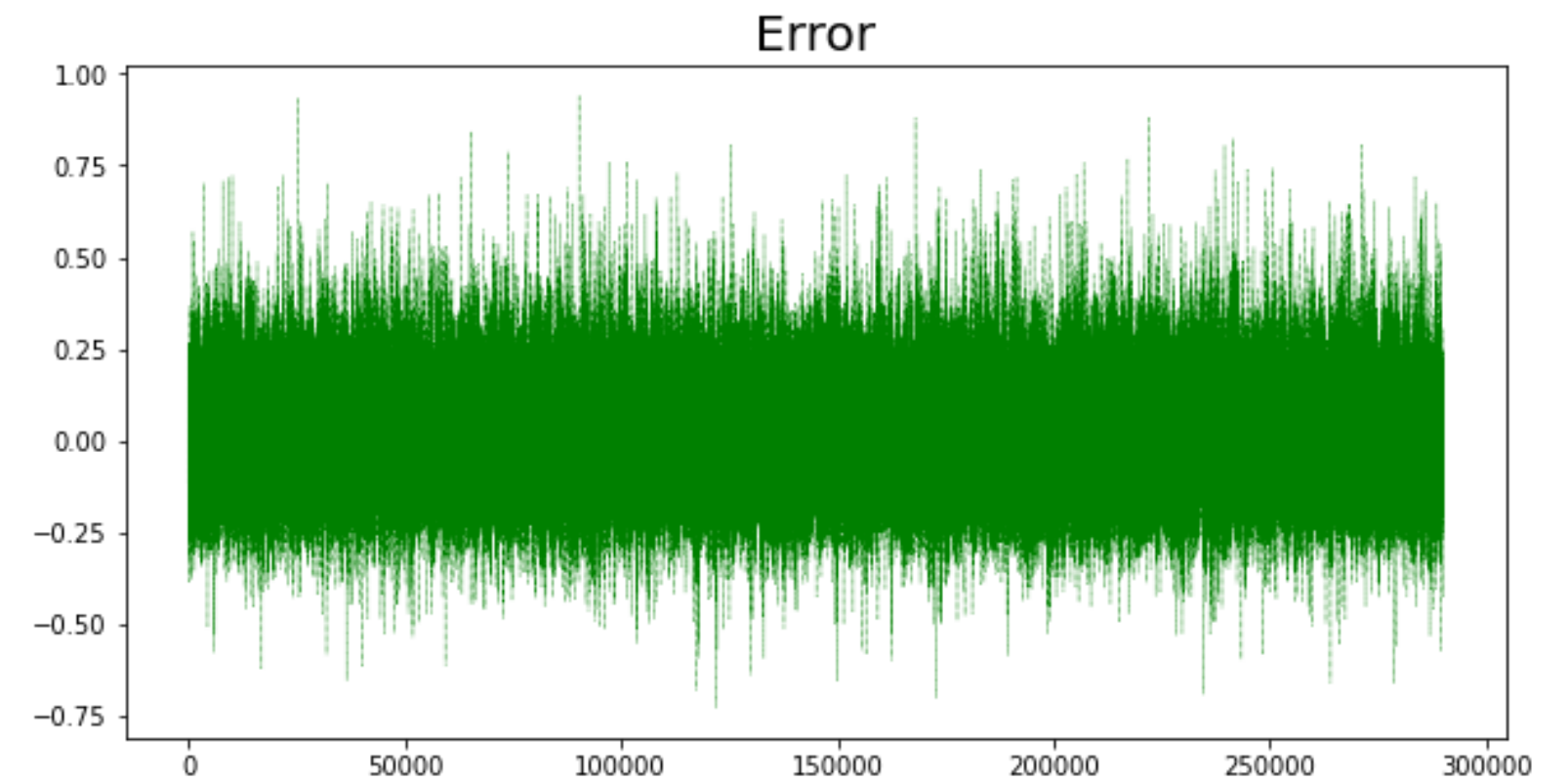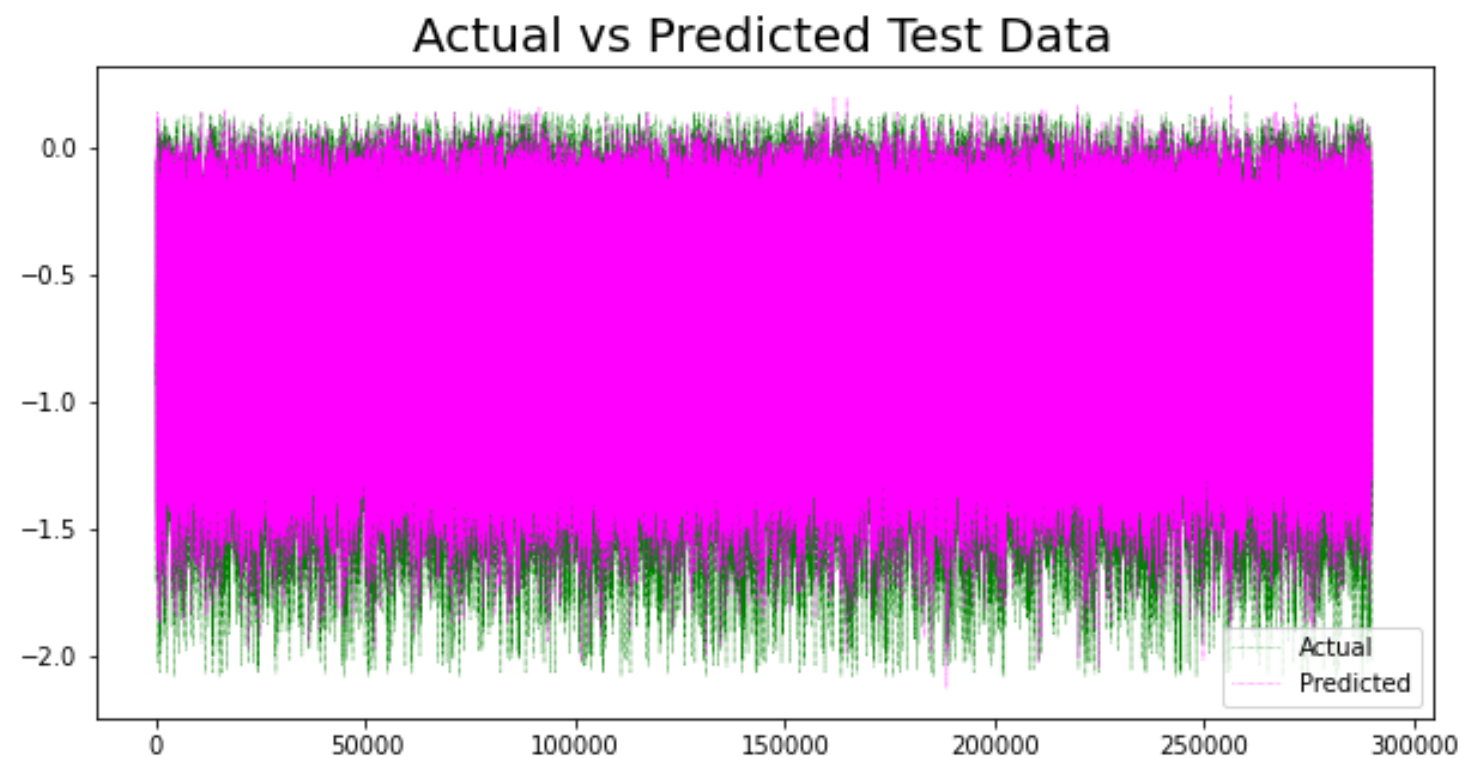SCORE FOR TRAINING DATASET**

Train MSE : 0.010872571045036155
Train RMSE : 0.10427162147505022
Train R2 : 0.8943411007490413 Train
Adjusted R2 : 0.8943397344760607

**MSE RMSE R2 ADJUSTED R2
SCORE FOR Test DATASET**

Test MSE : 0.010934883632925212
Test RMSE : 0.10456999394149935
Test R2 : 0.8936283352802094
Test Adjusted R2 : 0.8936228331078213

## Ploting graph for Actual vs Predicted data

# Model Summary

**Here we compare different parameters of all models.**

## Train_data_df

| | Model Name | Train MSE | Train RMSE | Train R^2 | Train Adjusted R^2 |
|---|---|---|---|---|---|
| 0 | Linear Regression | 0.056833 | 0.238396 | 0.447703 | 0.447696 |
| 1 | Lasso Regression | 0.061326 | 0.247642 | 0.404034 | 0.404026 |
| 2 | LightGBM Regression | 0.010873 | 0.104272 | 0.894340 | 0.894340 |
| 3 | XGBoost Regressor | 0.017228 | 0.131257 | 0.832577 | 0.832574 |

## Test_data_df

| | Model Name | Test MSE | Test RMSE | Test R^2 | Test Adjusted R^2 |
|---|---|---|---|---|---|
| 0 | Linear Regression | 0.056709 | 0.238137 | 0.448348 | 0.448319 |
| 1 | Lasso Regression | 0.061287 | 0.247562 | 0.403817 | 0.403786 |
| 2 | LightGBM Regression | 0.010935 | 0.104570 | 0.893628 | 0.893623 |
| 3 | XGBoost Regressor | 0.019701 | 0.140361 | 0.808351 | 0.808341 |

# Conclusion

## Comparing different parameters of all models.

- We compared MSE, RMSE and R2 for all four regression models, to find which is the best model to predict the NYC taxi trip duration.
- The Linear Regression and Lasso Regression didn't show any good prediction as compared to the other two.

- From he comparison table we can clearly see that XGBoost and LightGBM are the best models to predict trip duration of theNYC taxi. While LightGBM is fastest and more accurate than XGBoost. So, in between these two LightGBM is the best model.

- R Square: R2 score represents the coefficient of how well the values fit compared to the original values. The value from 0 to 1 interpreted as percentages. The higher the value is, the better the model is.

- RMSE (Root Mean Squared Error): RMSE is the standard deviation of the residuals. Residuals are a measure of how far from the regression line data points are; RMSE is a measure of how spread out these residuals are. In other words, it displays how concentrated the data is around the line of best fit.