

CAR FARE PREDICTION

Himachala Venkata

16-11-2019

CONTENTS

CHAPTERS	Page No:
1. INTRODUCTION:.....	3
1.1. PROBLEM STATEMENT:.....	3
1.2. DATA:	3
2. METHODOLOGY:.....	4
2.1. Pre-Processing:.....	4
2.1.1. Removing values which are not within desired range(outlier) depending upon basic understanding of dataset:.....	9
2.1.2. Missing Value Analysis:	12
2.1.3. OUTLIER ANALYSIS:	14
2.1.4. FEATURE ENGINEERING:	15
2.1.5. FEATURE SELECTION:	18
1. Correlation Analysis:	18
2. Chi-Square Test of Independence:.....	20
3. Analysis Of Variance (ANOVA) Test:	20
4. Multicollinearity:	22
5. Feature Scaling:.....	23
3. Splitting train and Validation Dataset:	26
4. Hyperparameter Optimization:.....	26
5. Model Development	27
5.1. Model Performance:	28
6. Improving Accuracy	31
7. Finalize model	33
REFERENCE	34

CHAPTER-1

1. INTRODUCTION:

1.1. PROBLEM STATEMENT:

You are a cab rental start-up company. You have successfully run the pilot project and now want to launch your cab service across the country. You have collected the historical data from your pilot project and now have a requirement to apply analytics for fare prediction. You need to design a system that predicts the fare amount for a cab ride in the city.

1.2. DATA:

Number of attributes:

- pickup_datetime - timestamp value indicating when the cab ride started.
- pickup_longitude – float for longitude coordinate of where the cab ride started.
- pickup_latitude - float for latitude coordinate of where the cab ride started.
- dropoff_longitude - float for longitude coordinate of where the cab ride ended.
- dropoff_latitude - float for latitude coordinate of where the cab ride ended.
- passenger_count - an integer indicating the number of passengers in the cab ride.

CHAPTER-2

2. METHODOLOGY:

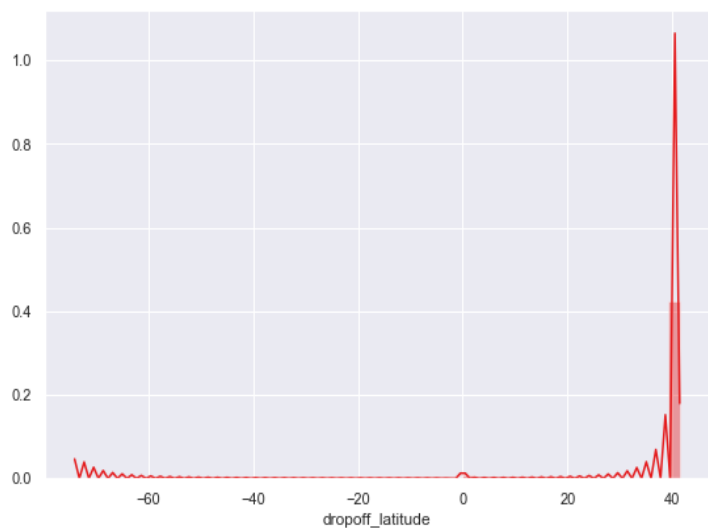
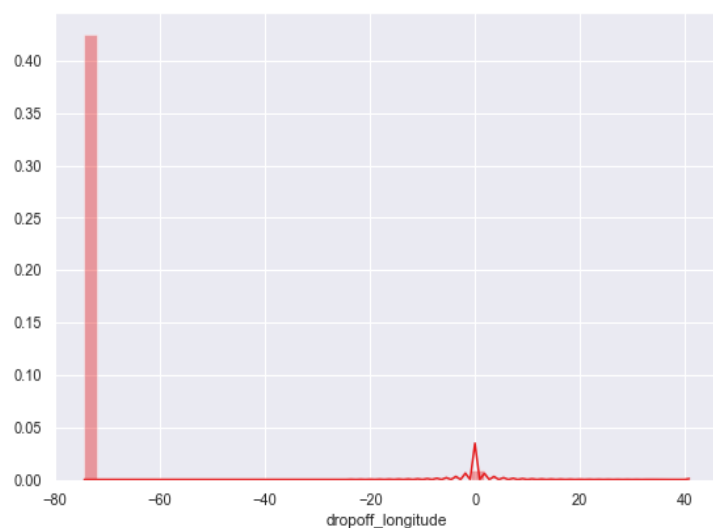
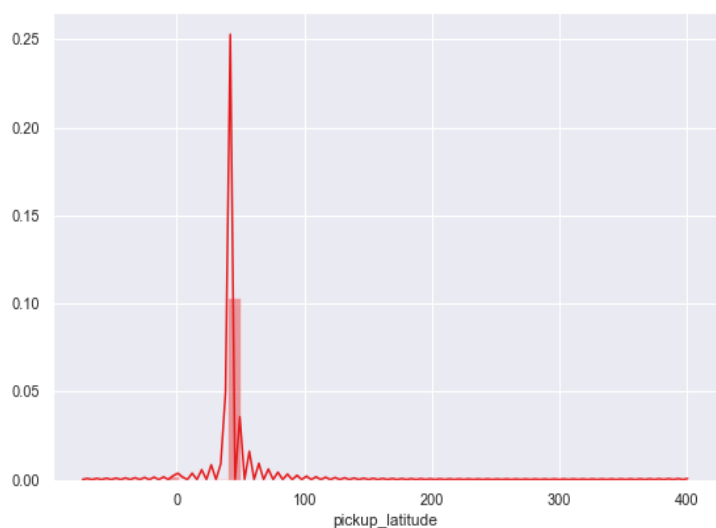
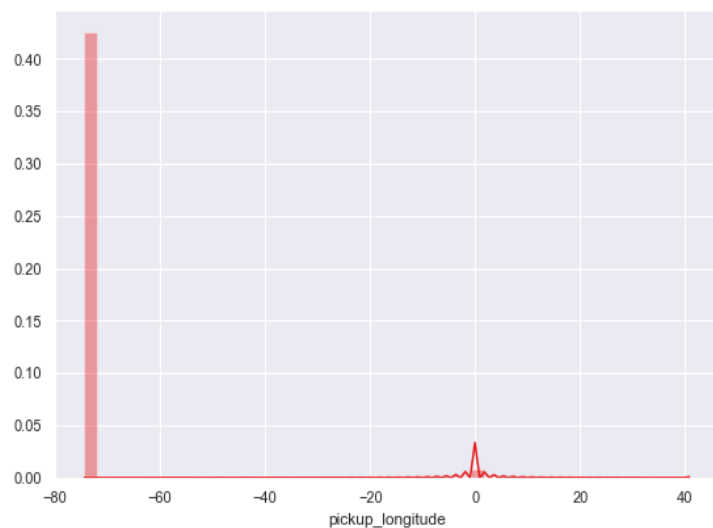
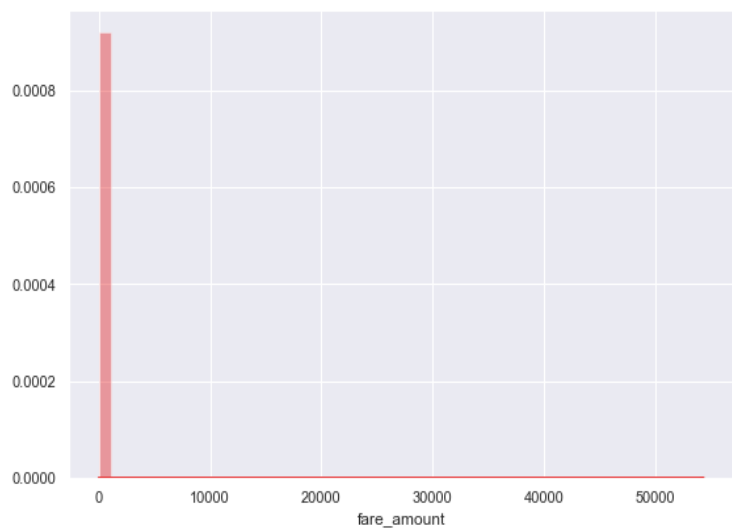
2.1. Pre-Processing:

- Data pre-processing is the first stage of any type of project.
- In this stage we get the feel of the data.
- We do this by looking at plots of independent variables vs target variables.
- If the data is messy, we try to improve it by sorting deleting extra rows and columns. This stage is called as Exploratory Data Analysis.
- This stage generally involves:
 - data cleaning
 - merging
 - Sorting
 - looking for outlier analysis,
 - looking for missing values in the data,
 - imputing missing values if found by various methods such as:
 - mean
 - median,
 - Mode
 - KNN imputation, etc.

Further we will look into what Pre-Processing steps do this project was involved in.

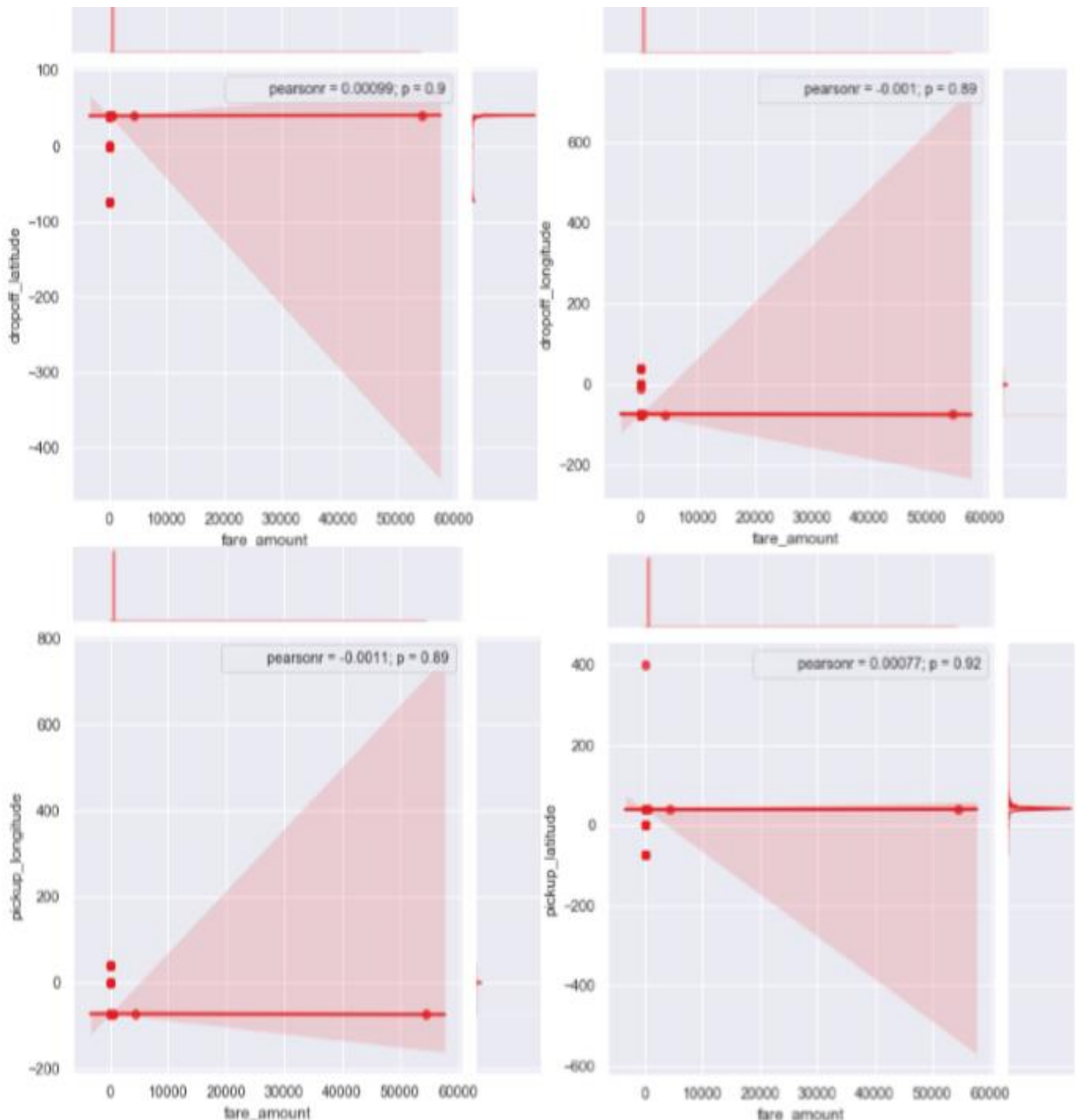
Getting feel of data via visualization:

Some Histogram plots from seaborn library for each individual variable created using `distplot()` method.

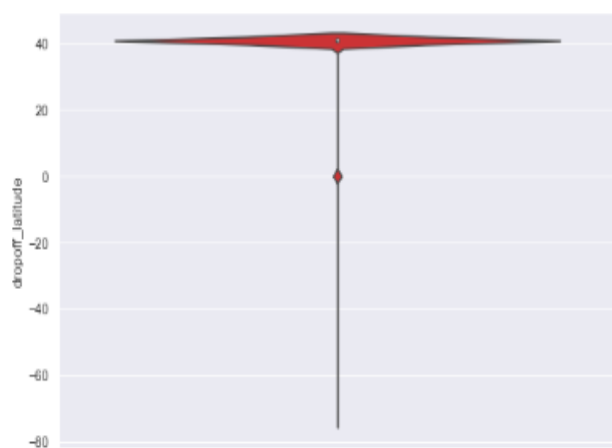
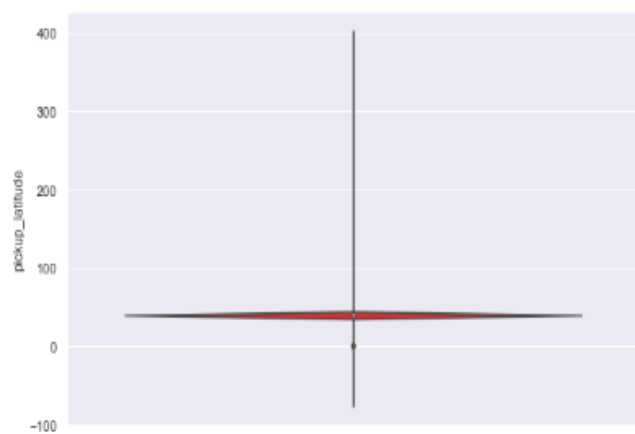
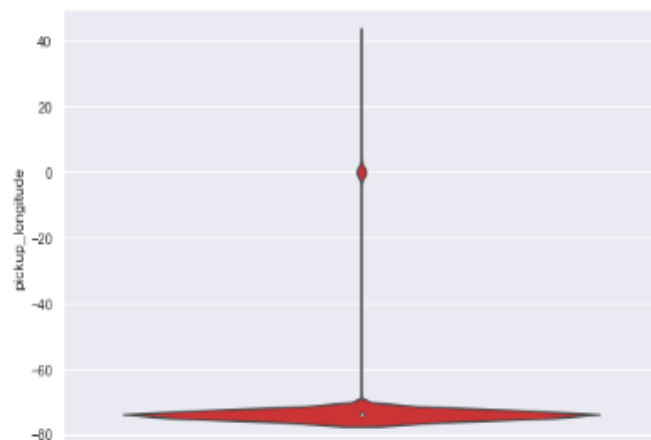
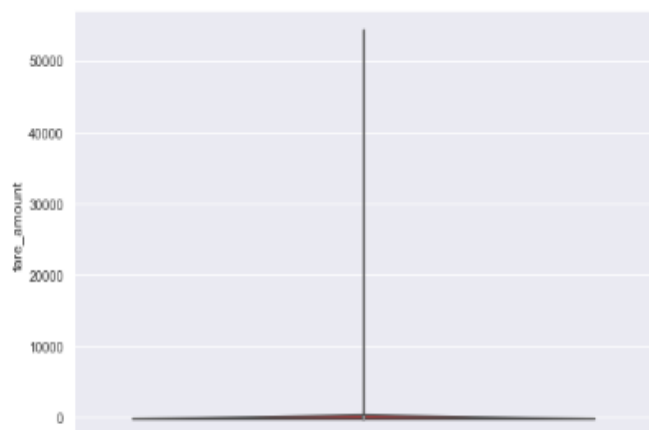


Some Jointplots:

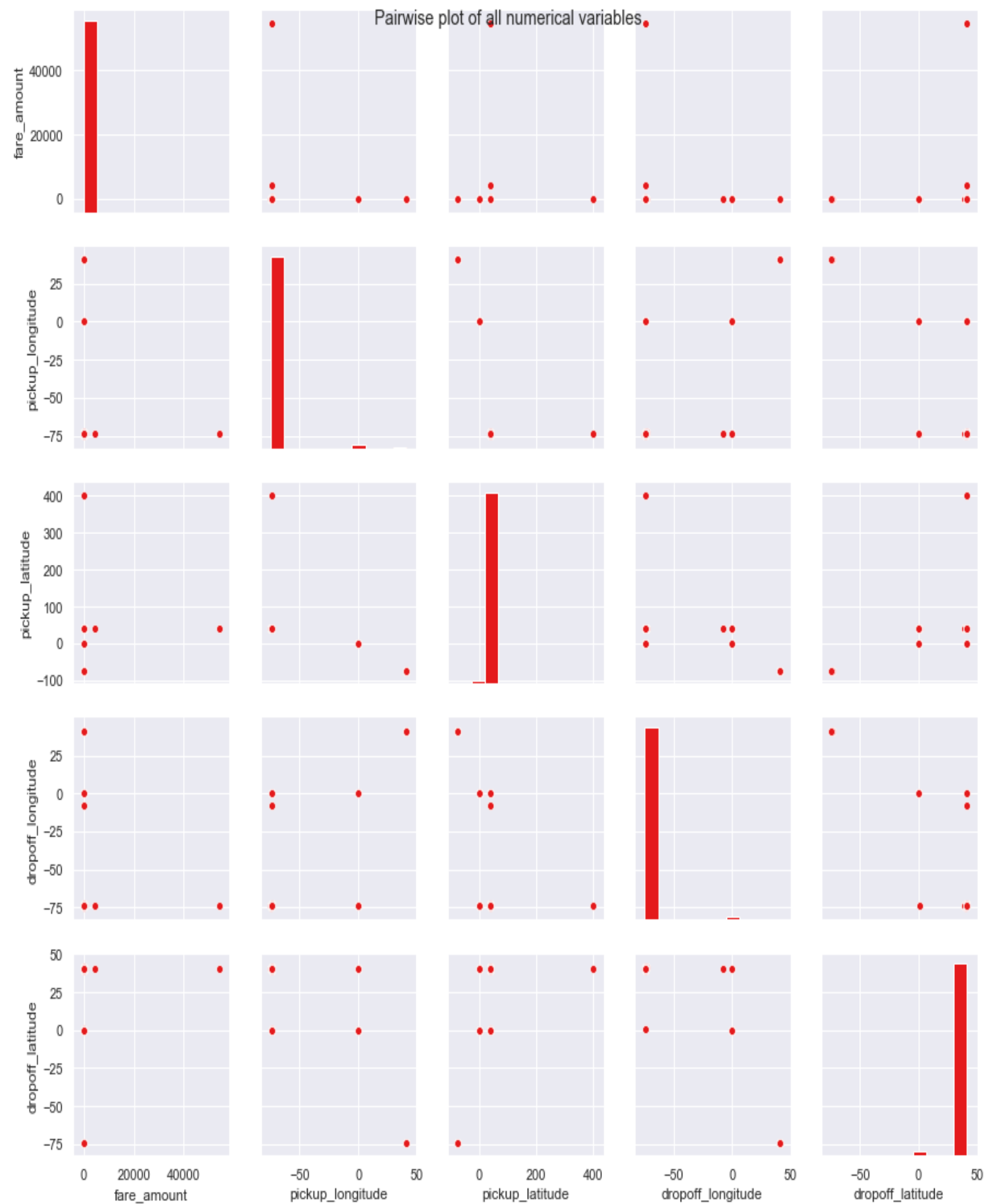
- They are used for Bivariate Analysis.
- Here we have plotted Scatter plot with Regression line between 2 variables along with separate Bar plots of both variables.
- Also, we have annotated Pearson correlation coefficient and p value.
- Plotted only for numerical/continuous variables
- Target variable 'fare_amount' Vs each numerical variable.



Some Violin Plots to get the idea about till what range is the variables is spread.



Pairwise Plots for all Numerical variables:



2.1.1. Removing values which are not within desired range(outlier) depending upon basic understanding of dataset:

In this step we will remove values in each variable which are not within desired range and we will consider them as outliers depending upon basic understanding of all the variables. You would think why haven't made those values NA instead of removing them well I did made them NA but it turned out to be a lot of missing values(NA's) in the dataset. Missing values percentage becomes very much high and then there will be no point of using that imputed data. Take a look at below 3 scenarios—

- If everything beyond range is made nan also except latitudes and longitudes then:

Variables	Missing_percentage	
0	Passenger_count	29.563702
1	Pickup_longitude	1.966764
2	Pickup_latitude	1.960540
3	Dropoff_longitude	1.954316
4	Dropoff_latitude	1.941868
5	Fare_amount	0.186718
6	Pickup_datetime	0.006224

- After imputing above mentioned missing values kNN algorithm imputes every value to 0 at a particular row which was made nan using np.nan method:

```
fare_amount      0.0
pickup_longitude  0.0
pickup_latitude  0.0
dropoff_longitude 0.0
dropoff_latitude  0.0
passenger_count  0.0
Name: 1000, dtype: float64
```

- And If everything is dropped which are beyond range then below are the missing percentages for each variable:

Variables	Missing_percentage	
0	Passenger_count	0.351191
1	Pickup_longitude	0.140476
2	Pickup_latitude	0.006385
3	Dropoff_longitude	0.000000
4	Dropoff_latitude	0.000000
5	Fare_amount	0.000000
6	Pickup_datetime	0.000000

- After imputing above mentioned missing values kNN algorithm imputes every value to 0 at a particular row which was made nan using np.nan method:

fare_amount 0.0

pickup_longitude 0.0

pickup_latitude 0.0

dropoff_longitude 0.0

dropoff_latitude 0.0

passenger_count 0.0

Name: 1000, dtype: float64

- And If everything is dropped which are beyond range then below are the missing percentages for each variable:

Variables	Missing_percentage	
0	Passenger_count	0.351191
1	Pickup_longitude	0.140476

2	Pickup_latitude	0.006385
3	Dropoff_longitude	0.000000
4	Dropoff_latitude	0.000000
5	Fare_amount	0.000000
6	Pickup_datetime	0.000000

- After imputing above mentioned missing values kNN algorithm values at a particular row which was made nan using np.nan method

fare_amount 7.3698

pickup_longitude -73.9954

pickup_latitude 40.7597

dropoff_longitude -73.9876

dropoff_latitude 40.7512

passenger_count 2

Name: 1000, dtype: object

- If everything beyond range is made nan except passenger_count:

Variables	Missing_percentage	
0	Passenger_count	1.951342
1	Pickup_longitude	1.951342
2	Pickup_latitude	1.945087
3	Dropoff_longitude	1.938833
4	Dropoff_latitude	0.343986
5	Fare_amount	0.181375
6	Pickup_datetime	0.006254

- After imputing above mentioned missing values kNN algorithm imputes every value to 0 at a particular row which was made nan using np.nan method:

```
fare_amount      0.0
pickup_longitude 0.0
pickup_latitude  0.0
dropoff_longitude 0.0
dropoff_latitude 0.0
passenger_count  0.0
Name: 1000, dtype: float64
```

2.1.2. Missing Value Analysis:

In this step we look for missing values in the dataset like empty row column cell which was left after removing special characters and punctuation marks. Some missing values are in form of NA. missing values left behind after outlier analysis; missing values can be in any form. Unfortunately, in this dataset we have found some missing values. Therefore, we will do some missing value analysis. Before imputed we selected random row no-1000 and made it NA, so that we will compare original value with imputed value and choose best method which will impute value closer to actual value.

	Index	
0	Fare_amount	22
1	Pickup_datetime	1
2	Pickup_longitude	0
3	Pickup_latitude	0
4	Dropoff_longitude	0
5	Dropoff_latitude	0
6	Passenger_count	55

We will impute values for fare_amount and passenger_count both of them has missing values 22 and 55 respectively. We will drop 1 value in pickup_datetime i.e it will be an entire row to drop.

Below are the missing value percentage for each variable:

Variables	Missing_percentage	
0	Passenger_count	0.351191
1	Pickup_longitude	0.140476
2	Pickup_latitude	0.006385
3	Dropoff_longitude	0.000000
4	Dropoff_latitude	0.000000
5	Fare_amount	0.000000
6	Pickup_datetime	0.000000

And below is the Standard deviation of particular variable which has missing values in them:

fare_amount 435.982171

passenger_count 1.266096

dtype: float64

We'd tried central statistical methods and algorithmic method--KNN to impute missing values in the dataset.

1. For Passenger_count:

Actual value = 1

Mode = 1

KNN = 2

We will choose the KNN method here because it maintains the standard deviation of variable. We will not use Mode method because whole variable will be more biased towards 1 passenger_count also passenger_count has maximum value equals to 1

2. For fare_amount:

Actual value = 7.0,

Mean = 15.117,

Median = 8.5,

KNN = 7.369801

We will Choose KNN method here because it imputes value closest to actual value also it maintains the Standard deviation of the variable.

Standard deviation for passenger_count and fare_amount after KNN imputation:

fare_amount 435.661995

passenger_count 1.264322

dtype: float64

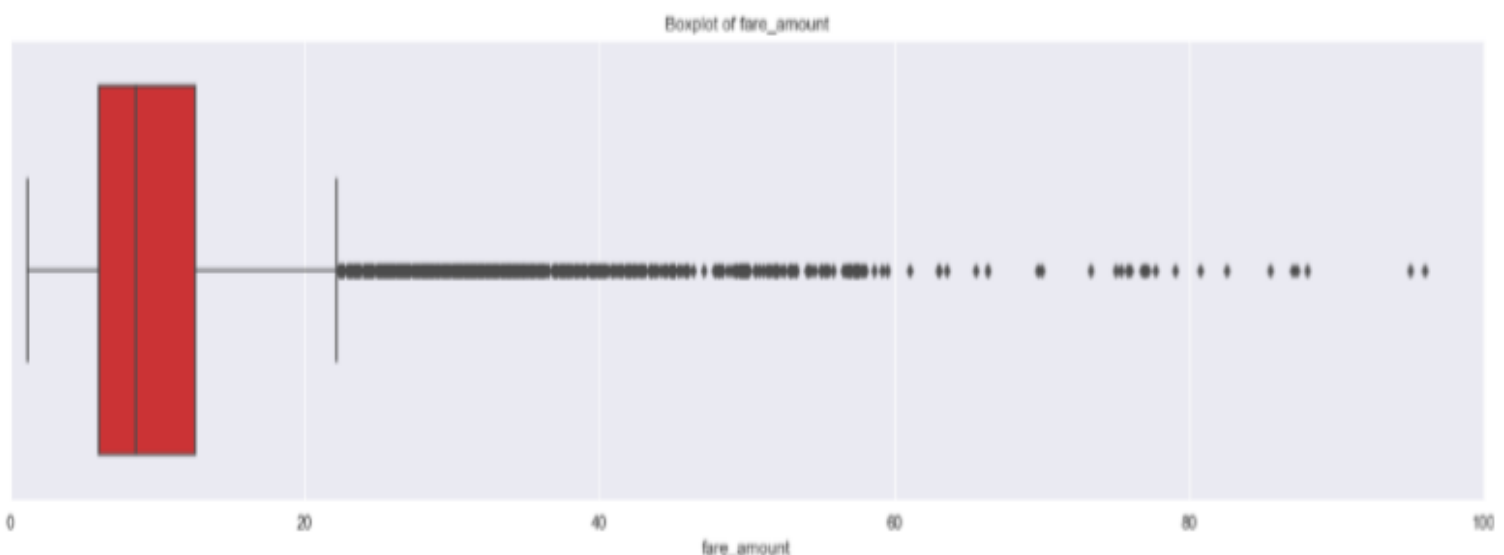
2.1.3. OUTLIER ANALYSIS:

We look for outlier in the dataset by plotting Boxplots. There are outliers present in the data. we have removed these outliers. This is how we done,

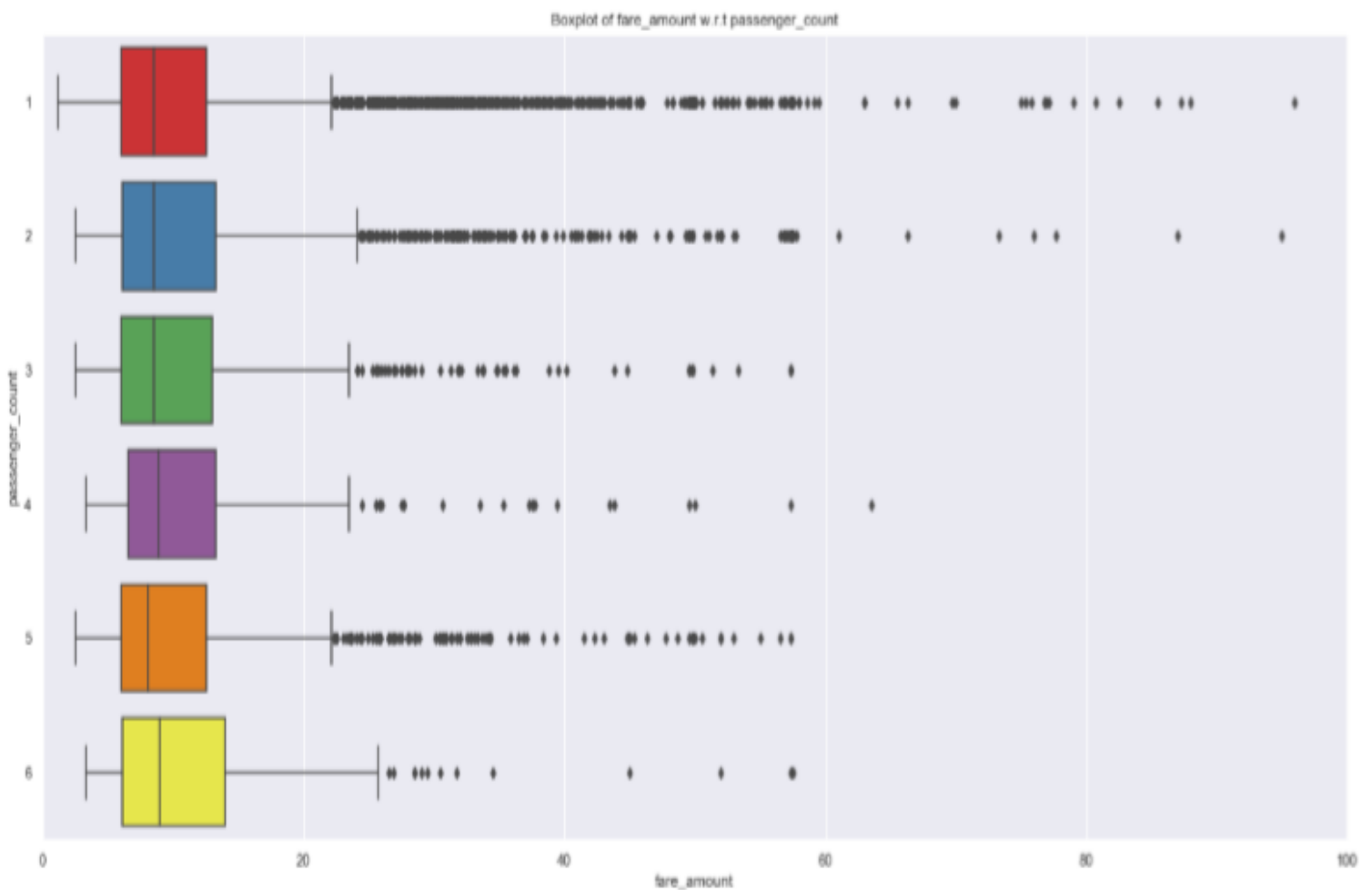
- ✓ We replaced them with Nan values or we can say created missing values.
- ✓ Then we imputed those missing values with KNN method.

We Will do Outlier Analysis only on Fare_amount just for now and we will do outlier analysis after feature engineering latitudes and longitudes.

Univariate Boxplots: Boxplots for all Numerical Variables including target variable.



Bivariate Boxplots: Boxplot for Numerical Variable Vs Categorical Variable.



From above Boxplots we see that 'fare_amount' have outliers in it:

'fare_amount' has 1359 outliers.

We successfully imputed these outliers with KNN and K value is 3

2.1.4. FEATURE ENGINEERING:

Feature Engineering is used to drive new features from existing features.

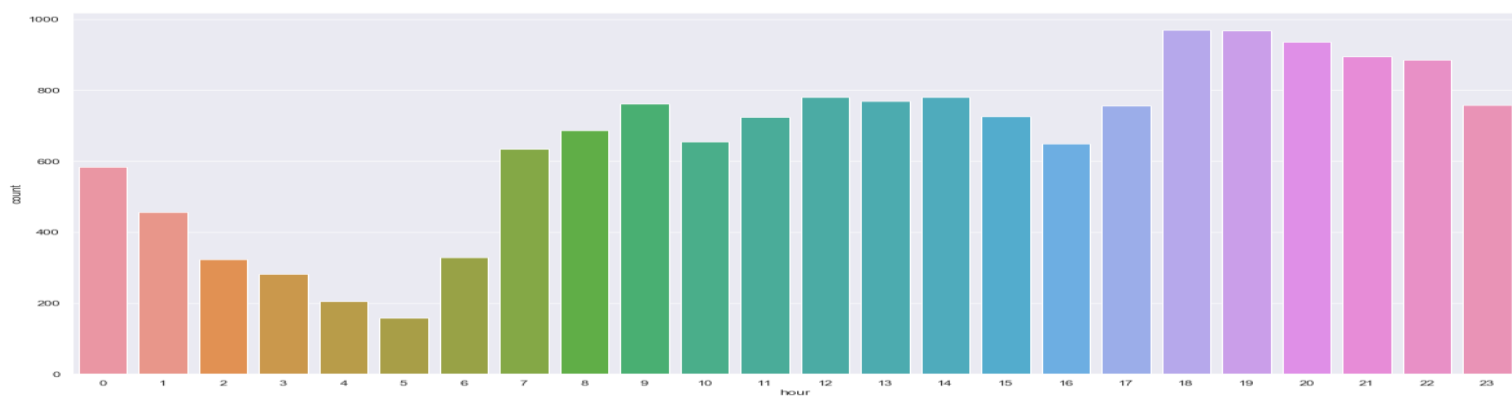
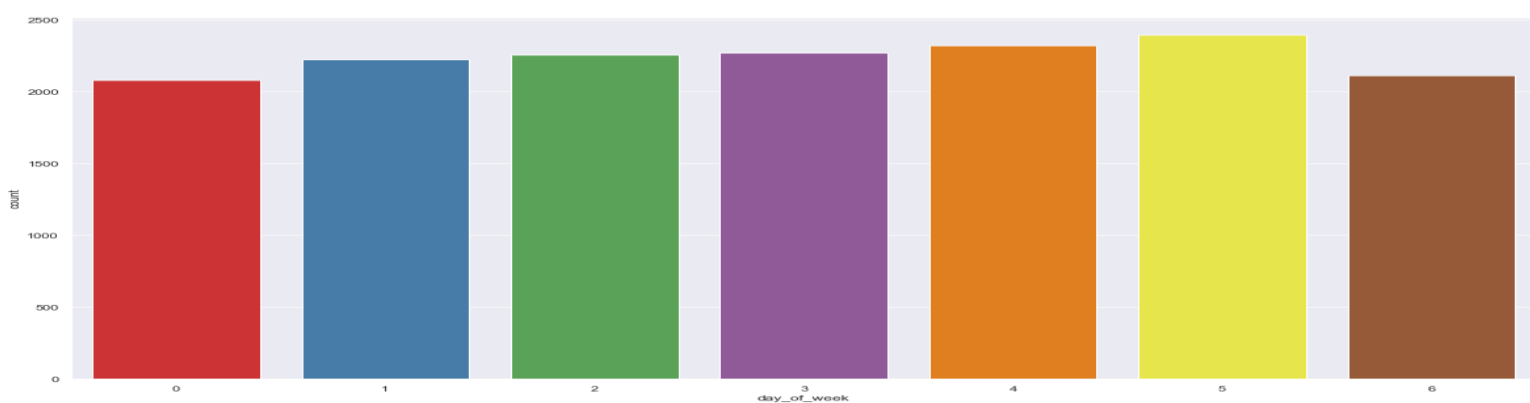
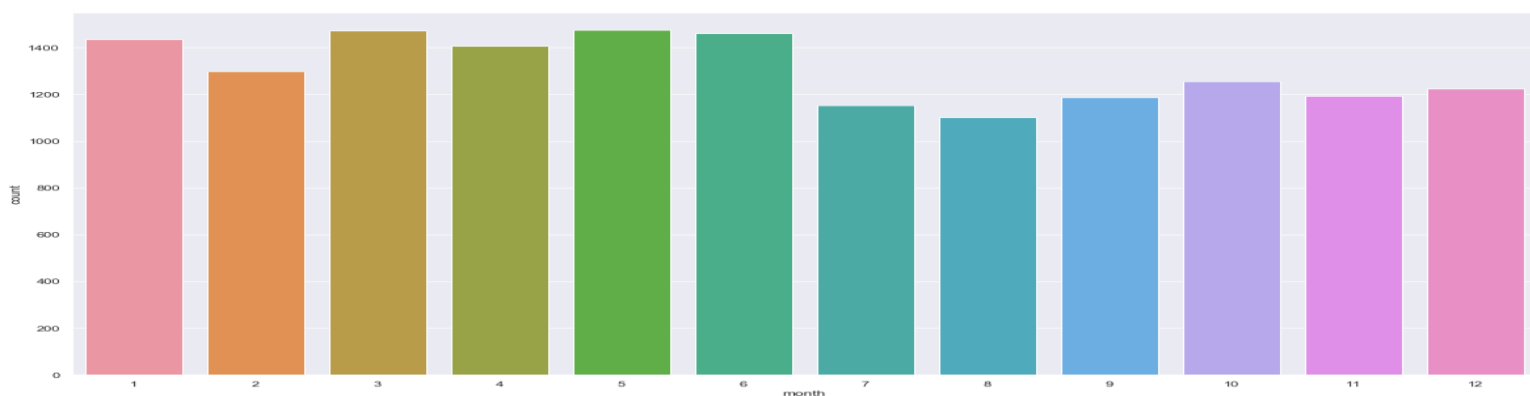
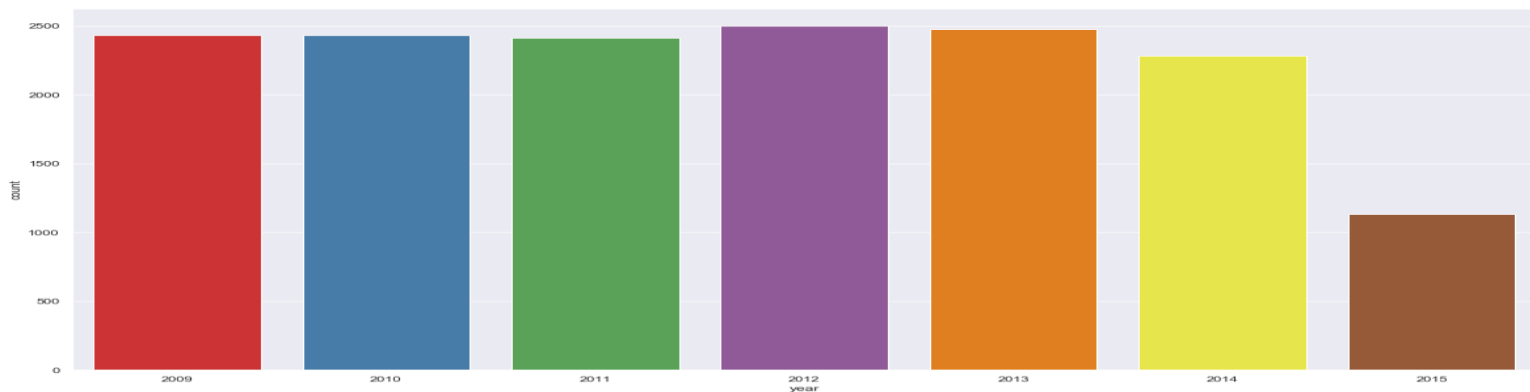
1. For 'pickup_datetime' variable:

We will use this timestamp variable to create new variables.

New features will be year, month, day_of_week, hour.

- 'year' will contain only years from pickup_datetime. For ex. 2009, 2010, 2011, etc.

- 'month' will contain only months from pickup_datetime. For ex. 1 for January, 2 for February, etc.
- 'day_of_week' will contain only week from pickup_datetime. For ex. 1 which is for Monday, 2 for Tuesday, etc.
- 'hour' will contain only hours from pickup_datetime. For ex. 1, 2, 3, etc.



As we have now these new variables we will categorize them to new variables like Session from hour column, seasons from month column, week: weekday/weekend from day_of_week variable.

So, session variable which will contain categories—morning, afternoon, evening, night_PM, night_AM.

Seasons variable will contain categories—spring, summer, fall, winter.

Week will contain categories—weekday, weekend.

We will one-hot-encode session, seasons, week variable.

2. For 'passenger_count' variable:

As passenger_count is a categorical variable we will one-hot-encode it.

3. For 'Latitudes' and 'Longitudes' variables:

As we have latitude and longitude data for pickup and dropoff, we will find the distance the cab travelled from pickup and dropoff location.

We will use both haversine and vincenty methods to calculate distance. For haversine, variable name will be 'great_circle' and for vincenty, new variable name will be 'geodesic'. As Vincenty is more accurate than haversine. Also, vincenty is preferred for short distances.

Therefore, we will drop great_circle. Columns in training data after feature engineering:

```
Index(['fare_amount', 'passenger_count_2', 'passenger_count_3', 'passenger_count_4',
      'passenger_count_5', 'passenger_count_6', 'season_spring', 'season_summer',
      'season_winter', 'week_weekend', 'session_evening', 'session_morning',
      'session_night_AM', 'session_night_PM', 'year_2010', 'year_2011', 'year_2012',
      'year_2013', 'year_2014', 'year_2015', 'geodesic'], dtype='object')
```

Columns in testing data after feature engineering:

```
Index(['passenger_count_2', 'passenger_count_3', 'passenger_count_4',
      'passenger_count_5', 'passenger_count_6', 'season_spring', 'season_summer',
      'season_winter', 'week_weekend', 'session_evening', 'session_morning',
      'session_night_AM', 'session_night_PM', 'year_2010', 'year_2011', 'year_2012',
      'year_2013', 'year_2014', 'year_2015', 'geodesic'], dtype='object')
```

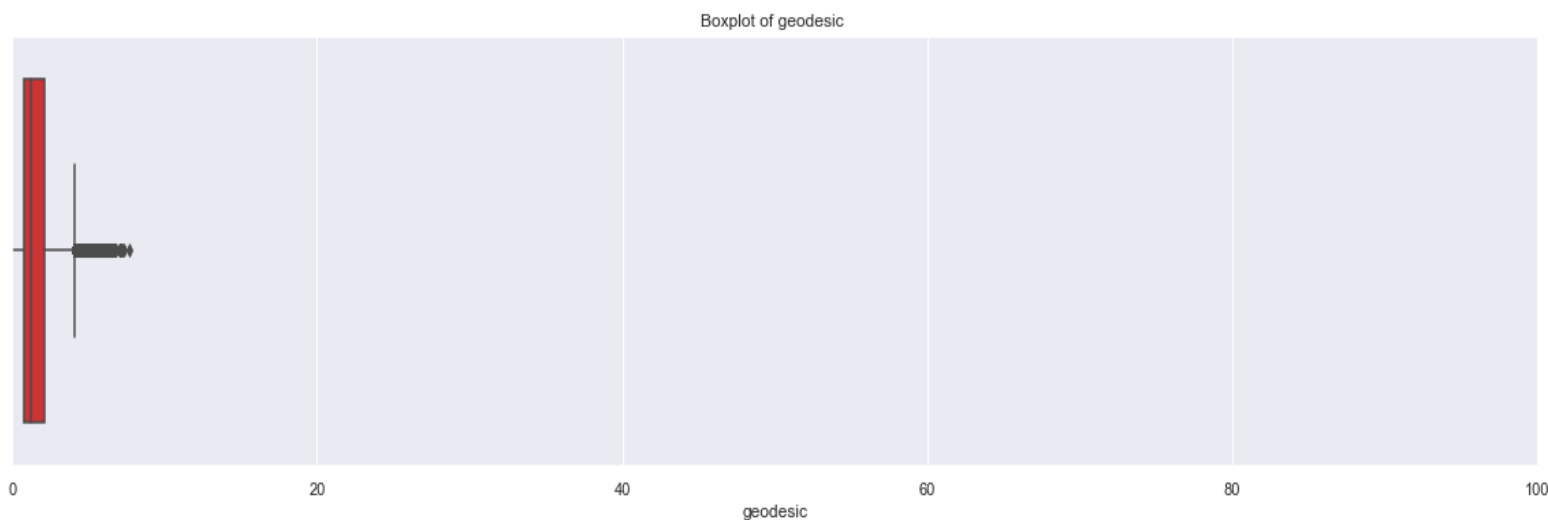
we will plot boxplot for our new variable 'geodesic':

Boxplot of geodesic



We see that there are outliers in 'geodesic'

Boxplot of 'geodesic' for range 0 to 100 miles



We will treat these outliers like we previously did.

2.1.5. FEATURE SELECTION:

In this step we would allow only to pass relevant features to further steps. We remove irrelevant features from the dataset. We do this by some statistical techniques, like we look for features which will not be helpful in predicting the target variables. In this dataset we have to predict the fare_amount.

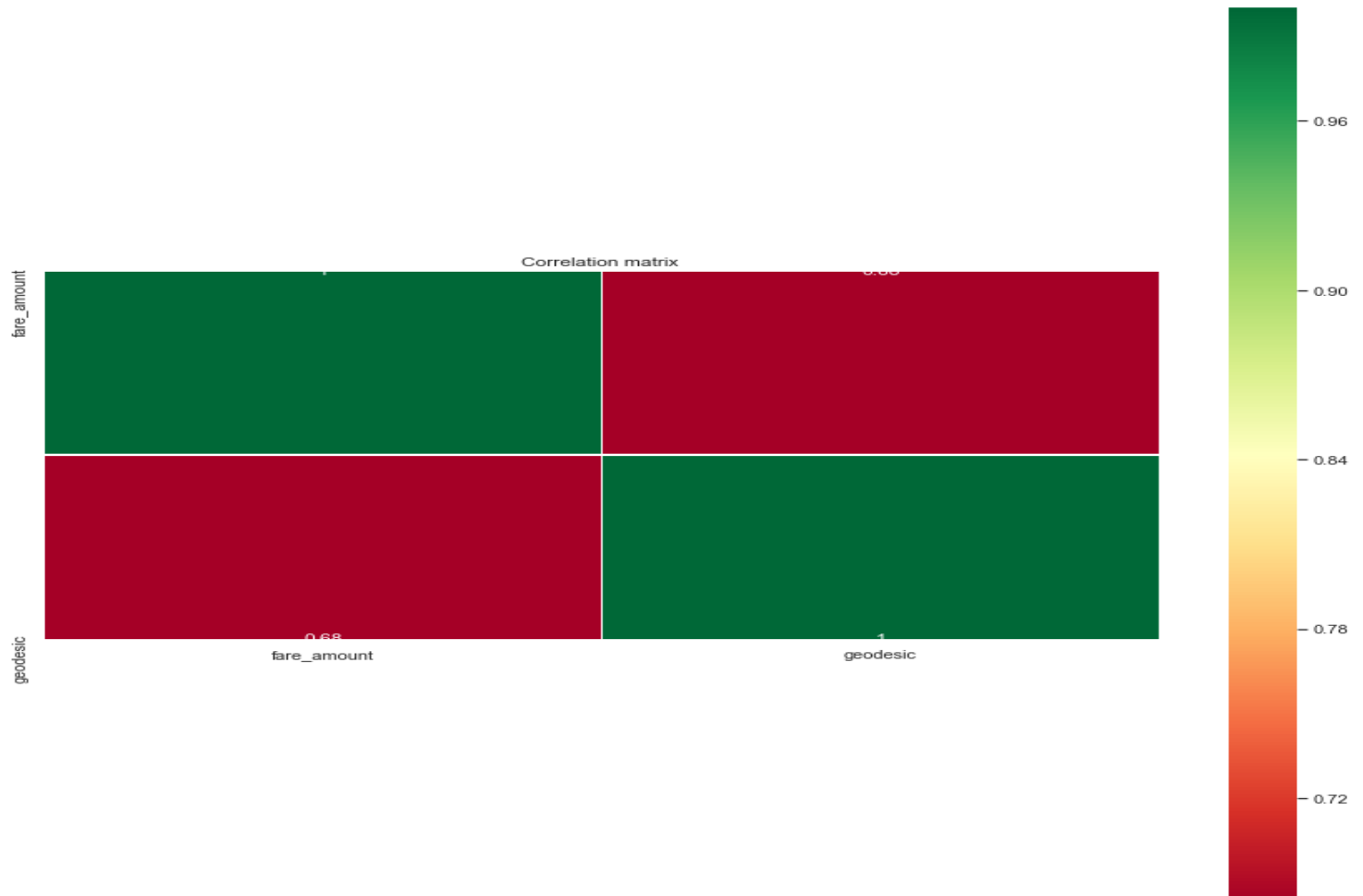
Further below are some types of test involved for feature selection:

1. Correlation Analysis:

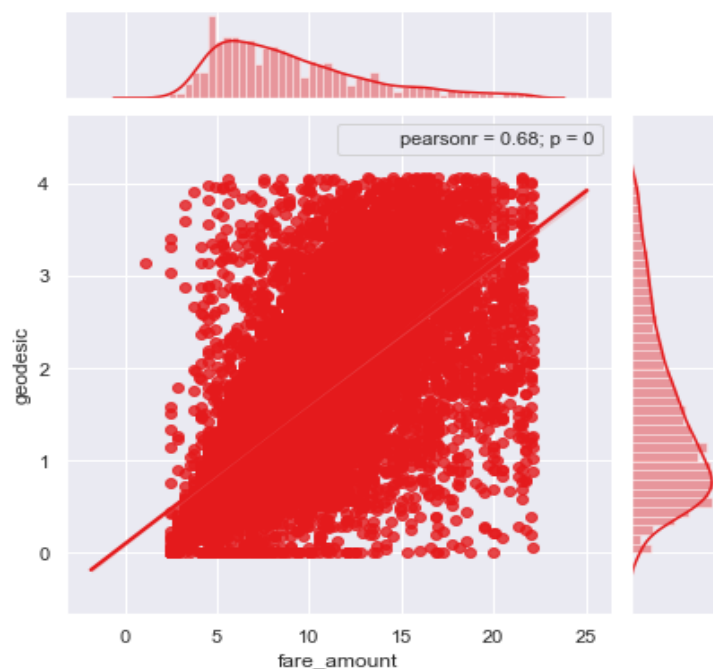
This requires only numerical variables. Therefore, we will filter out only numerical variables and feed it to correlation analysis. We do this by plotting correlation plot for all numerical variables. There should be no correlation between independent variables but there should be high correlation between independent variable and dependent variable. So, we plot the correlation plot. we can see that in correlation plot faded colour like skin colour indicates that 2 variables are highly correlated with each other. As the colour fades correlation values increases.

From below correlation plot we see that:

- 'fare_amount' and 'geodesic' are very highly correlated with each other.
- As fare_amount is the target variable and 'geodesic' is independent variable we will keep 'geodesic' because it will help to explain variation in fare_amount.



Joint plot between 'geodesic' and 'fare_amount':



2. Chi-Square Test of Independence:

Unlike correlation analysis we will filter out only categorical variables and pass it to Chi-Square test. Chi-square test compares 2 categorical variables in a contingency table to see if they are related or not.

- I. Assumption for chi-square test: Dependency between Independent variable and dependent variable should be high and there should be no dependency among independent variables.
- II. Before proceeding to calculate chi-square statistic, we do the hypothesis testing:
 - a. Null hypothesis: 2 variables are independent.
 - b. Alternate hypothesis: 2 variables are not independent.

The interpretation of chi-square test:

- ✓ For theoretical or excel sheet purpose: If chi-square statistics is greater than critical value then reject the null hypothesis saying that 2 variables are dependent and if it's less, then accept the null hypothesis saying that 2 variables are independent.
- ✓ While programming: If p-value is less than 0.05 then we reject the null hypothesis saying that 2 variables are dependent and if p-value is greater than 0.05 then we accept the null hypothesis saying that 2 variables are independent.

Here we did the test between categorical independent variables pairwise.

- If $p\text{-value} < 0.05$ then remove the variable
- If $p\text{-value} > 0.05$ then keep the variable

3. Analysis Of Variance (ANOVA) Test:

- It is carried out to compare between each group in a categorical variable.
- ANOVA only lets us know the means for different groups are same or not. It doesn't help us identify which mean is different.

Hypothesis testing:

- Null Hypothesis: mean of all categories in a variable are same.
 - Alternate Hypothesis: mean of at least one category in a variable is different.
- If p-value is less than 0.05 then we reject the null hypothesis.
 - And if p-value is greater than 0.05 then we accept the null hypothesis.

Below is the anova analysis table for each categorical variable:

	df	sum_sq	mean_sq	F	PR(>F)
C(passenger_count_2)	1.0	3.434995	3.434995	0.213220	6.442619e-01
C(passenger_count_3)	1.0	132.234249	132.234249	8.208148	4.175825e-03
C(passenger_count_4)	1.0	182.677382	182.677382	11.339294	7.606710e-04
C(passenger_count_5)	1.0	147.344023	147.344023	9.146054	2.496594e-03
C(passenger_count_6)	1.0	200.407176	200.407176	12.439831	4.214776e-04
C(season_spring)	1.0	42.743011	42.743011	2.653178	1.033633e-01
C(season_summer)	1.0	13.706072	13.706072	0.850774	3.563477e-01
C(season_winter)	1.0	290.635488	290.635488	18.040554	2.174990e-05
C(week_weekend)	1.0	19.015324	19.015324	1.180334	2.773047e-01
C(session_night_AM)	1.0	971.255609	971.255609	60.288539	8.697668e-15
C(session_night_PM)	1.0	165.207739	165.207739	10.254904	1.365945e-03
C(session_evening)	1.0	37.207511	37.207511	2.309574	1.285995e-01
C(session_morning)	1.0	83.968705	83.968705	5.212171	2.244266e-02
C(year_2010)	1.0	1061.365927	1061.365927	65.881938	5.141274e-16
C(year_2011)	1.0	871.782895	871.782895	54.113991	1.985740e-13
C(year_2012)	1.0	297.639977	297.639977	18.475342	1.731580e-05
C(year_2013)	1.0	276.057693	276.057693	17.135670	3.498491e-05
C(year_2014)	1.0	1007.664465	1007.664465	62.548539	2.772370e-15
C(year_2015)	1.0	1717.554314	1717.554314	106.613378	6.505331e-25
Residual	15640.0	251962.277484	16.110120	NaN	NaN

Looking at above table every variable has p value less than 0.05 so reject the null hypothesis.

4. Multicollinearity:

In regression, "multicollinearity" refers to predictors that are correlated with other predictors. Multicollinearity occurs when your model include multiple factors that are correlated not just to your response variable, but also to each other.

- I. Multicollinearity increases the standard errors of the coefficients.
- II. Increased standard errors in turn means that coefficients for some independent variables may be found not to be significantly different from 0.
- III. In other words, by overinflating the standard errors, multicollinearity makes some variables statistically insignificant when they should be significant. Without multicollinearity (and thus, with lower standard errors), those coefficients might be significant.
- IV. VIF is always greater or equal to 1. if VIF is 1 --- Not correlated to any of the variables. if VIF is between 1-5 --- Moderately correlated. if VIF is above 5 --- Highly correlated. If there are multiple variables with VIF greater than 5, only remove the variable with the highest VIF.
- V. And if the VIF goes above 10, you can assume that the regression coefficients are poorly estimated due to multicollinearity.

Below is the table for VIF analysis for each independent variable:

	VIF	features
0	15.597787	Intercept
1	1.040767	passenger_count_2[T.1.0]
2	1.019517	passenger_count_3[T.1.0]
3	1.011818	passenger_count_4[T.1.0]
4	1.024938	passenger_count_5[T.1.0]
5	1.017374	passenger_count_6[T.1.0]
6	1.642308	season_spring[T.1.0]
7	1.552431	season_summer[T.1.0]
8	1.587593	season_winter[T.1.0]
9	1.051211	week_weekend[T.1.0]
10	1.364003	session_night_AM[T.1.0]
11	1.421734	session_night_PM[T.1.0]

12	1.526639	session_evening[T.1.0]
13	1.559132	session_morning[T.1.0]
14	1.691580	year_2010[T.1.0]
15	1.687637	year_2011[T.1.0]
16	1.711577	year_2012[T.1.0]
17	1.709746	year_2013[T.1.0]
18	1.665111	year_2014[T.1.0]
19	1.406894	year_2015[T.1.0]
20	1.015210	geodesic

We have checked for multicollinearity in our Dataset and all VIF values are below 5.

5. Feature Scaling:

Data Scaling methods are used when we want our variables in data to scaled on common ground. It is performed only on continuous variables.

- **Normalization:** Normalization refer to the dividing of a vector by its length. normalization normalizes the data in the range of 0 to 1. It is generally used when we are planning to use distance method for our model development purpose such as KNN. Normalizing the data improves convergence of such algorithms. Normalisation of data scales the data to a very small interval, where outliers can be loosed.
- **Standardization:** Standardization refers to the subtraction of mean from individual point and then dividing by its SD. Z is negative when the raw score is below the mean and Z is positive when above mean. When the data is distributed normally you should go for standardization.

Linear Models assume that the data you are feeding are related in a linear fashion, or can be measured with a linear distance metric.

Also, our independent numerical variable 'geodesic' is not distributed normally so we had chosen normalization over standardization.

We have checked variance for each column in dataset before Normalisation.

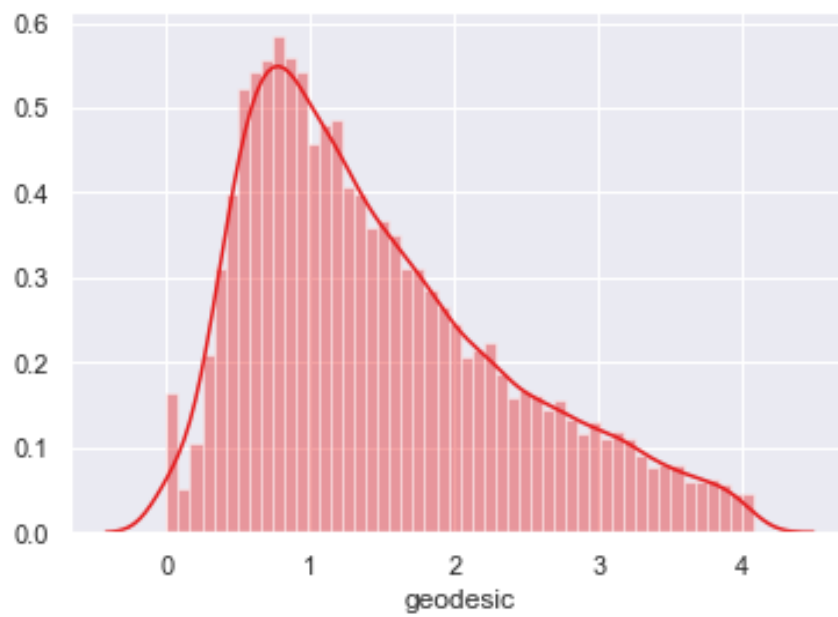
High variance will affect the accuracy of the model.

So, we want to normalise that variance.

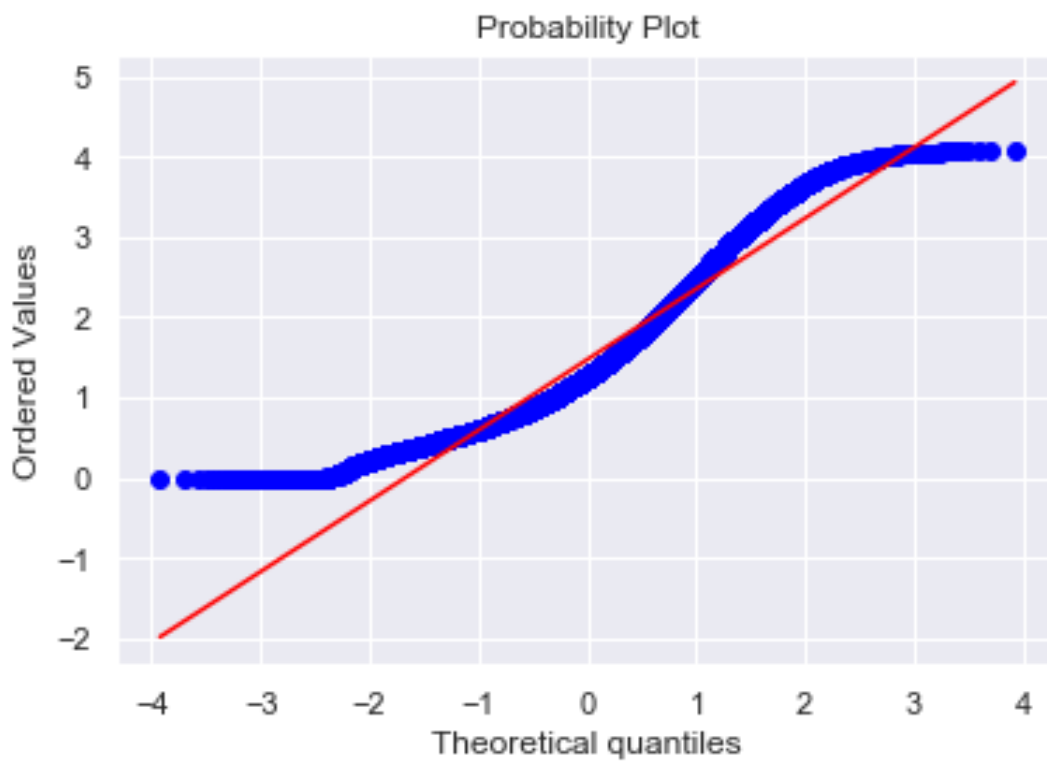
Graphs based on which standardization was chosen:

Note: It is performed only on Continuous variables.

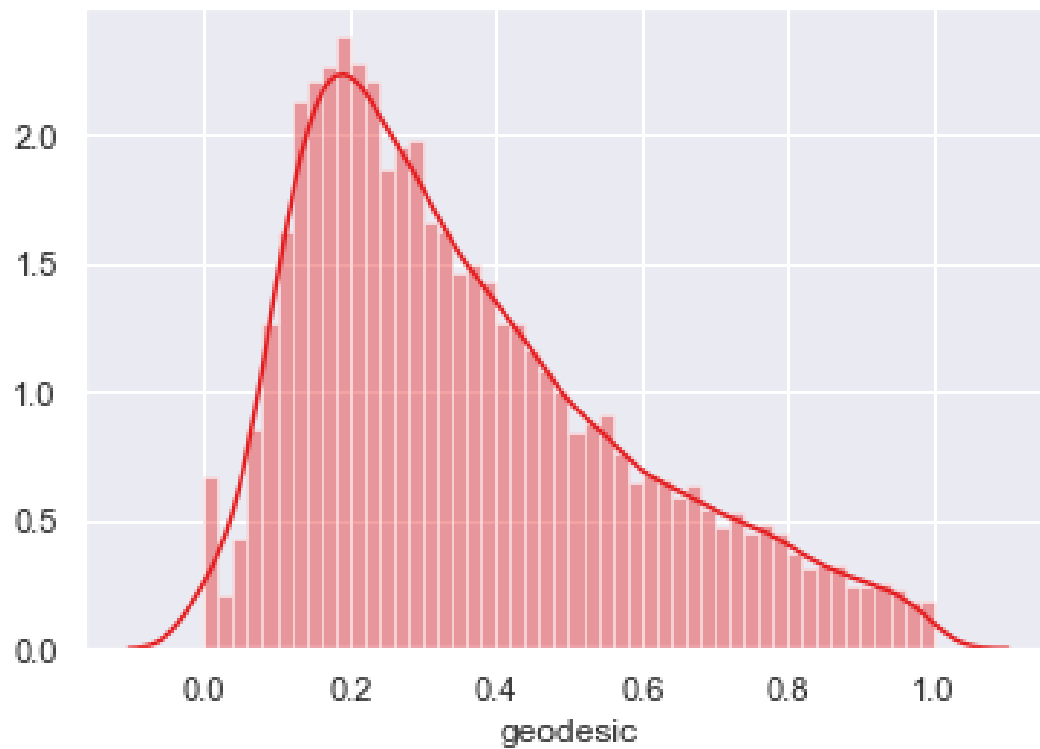
distplot() for 'geodesic' feature before normalization:



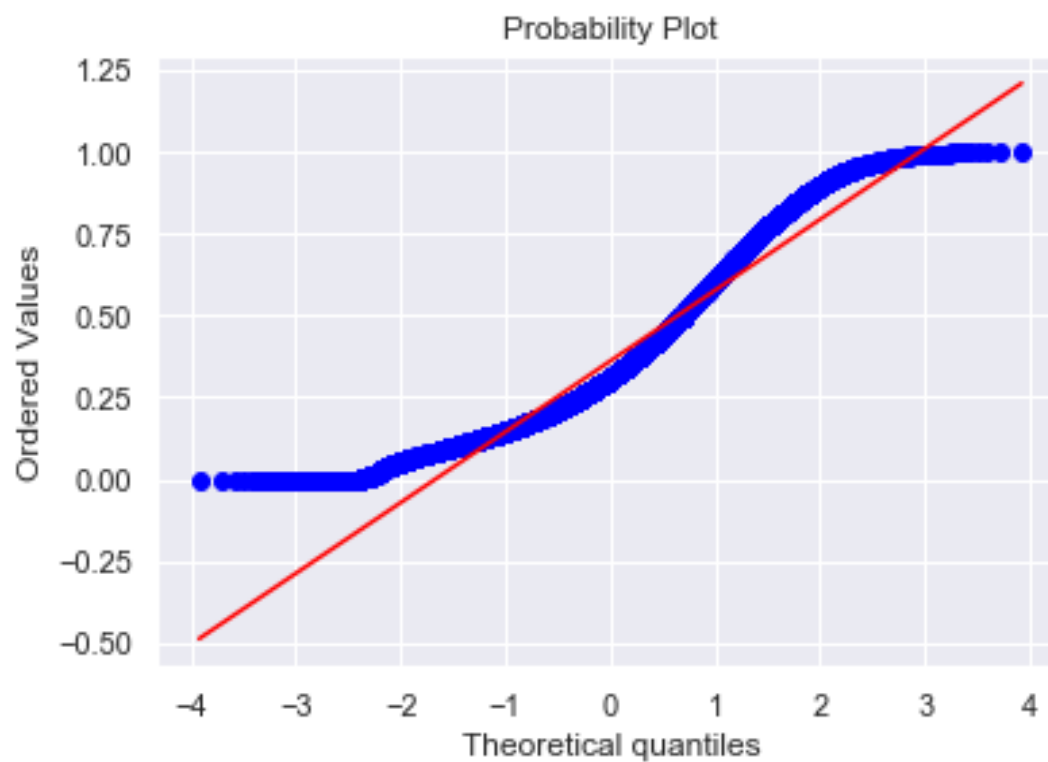
qq probability plot before normalization:



distplot() for 'geodesic' feature after normalization:



qq probability plot after normalisation:



CHAPTER 3

3. Splitting train and Validation Dataset:

- a. We have used sklearn's `train_test_split()` method to divide whole Dataset into train and validation dataset.
- b. 25% is in validation dataset and 75% is in training data.
- c. 11745 observations in training and 3915 observations in validation dataset.
- d. We will test the performance of model on validation dataset.
- e. The model which performs best will be chosen to perform on test dataset provided along with original train dataset.
- f. `X_train y_train`--are train subset. g) `X_test y_test`--are validation subset.

CHAPTER 4

4. Hyperparameter Optimization:

- To find the optimal hyperparameter we have used:
 - `sklearn.model_selection.GridSearchCV`.
 - `sklearn.model_selection.RandomizedSearchCV`
- `GridSearchCV` tries all the parameters that we provide it and then returns the best suited parameter for data.
- We gave parameter dictionary to `GridSearchCV` which contains keys which are parameter names and values are the values of parameters which we want to try for.

Below are best hyperparameter we found for different models:

Multi Linear Regression:

Tuned Decision reg Parameters: {'copy_X': True, 'fit_intercept': True} Best score is 0.4897614985052175

Ridge Regression:

Tuned Decision ridge Parameters: {'alpha': 0.0007906043210907702, 'max_iter': 500, 'normalize': True} Best score is 0.4897618771904997

Lasso Regression:

Tuned Decision lasso Parameters: {'alpha': 0.00014563484775012445, 'max_iter': 500, 'normalize': False} Best score is 0.48976204538973467

Decision Tree Regression:

Tuned Decision Tree Parameters: {'max_depth': 6, 'min_samples_split': 8} Best score is 0.47898878199556655

Random Forest Regression:

Tuned Random Forest Parameters: {'n_estimators': 300, 'min_samples_split': 2, 'min_samples_leaf': 4, 'max_features': 'sqrt', 'max_depth': 14, 'bootstrap': True} Best score is 0.49435114082662474

XgBoost Regression:

Tuned Xgboost Parameters: {'subsample': 0.1, 'reg_alpha': 0.08685113737513521, 'n_estimators': 200, 'max_depth': 3, 'learning_rate': 0.05, 'colsample_bytree': 0.7000000000000001, 'colsample_bynode': 0.7000000000000001, 'colsample_bylevel': 0.9000000000000001} Best score is 0.501889441417865

CHAPTER 5

5. Model Development

Our problem statement wants us to predict the fare_amount. This is a Regression problem. So, we are going to build regression models on training data and predict it on test data. In this project I have built models using 5 Regression Algorithms:

- I. Linear Regression
- II. Ridge Regression
- III. Lasso Regression
- IV. Decision Tree
- V. Random Forest
- VI. Xgboost Regression

We will evaluate performance on validation dataset which was generated using Sampling. We will deal with specific error metrics like:

Regression metrics for our Models: -

- r square
- Adjusted r square
- MAPE(Mean Absolute Percentage Error)
- MSE(Mean square Error)
- RMSE(Root Mean Square Error)
- RMSLE(Root Mean Squared Log Error)

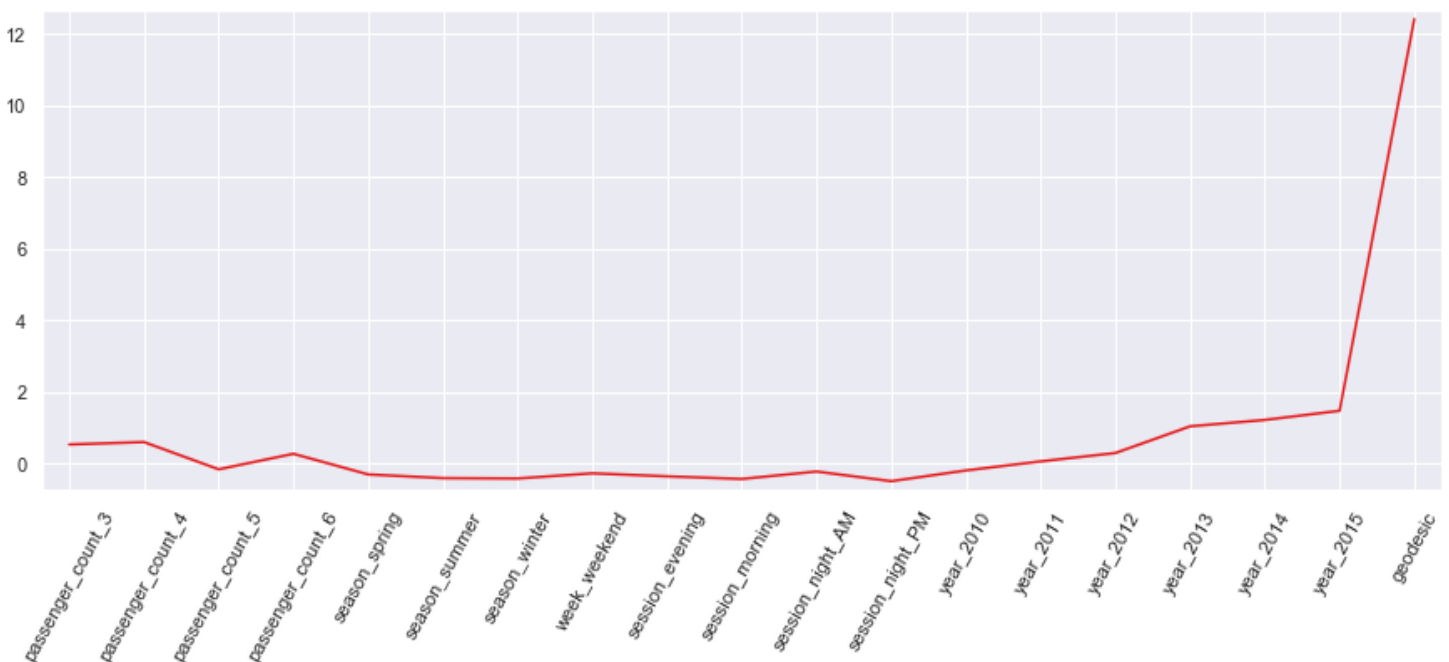
5.1. Model Performance:

Here, we will evaluate the performance of different Regression models based on different Error Metrics

Multiple Linear Regression:

Error Metrics	r-square	Adjusted r-square	MAPE	MSE	RMSE	RMSLE
Train	0.492	0.491	24.124	8.411	2.900	0.267
Test	0.484	0.482	24.602	8.538	2.922	0.267

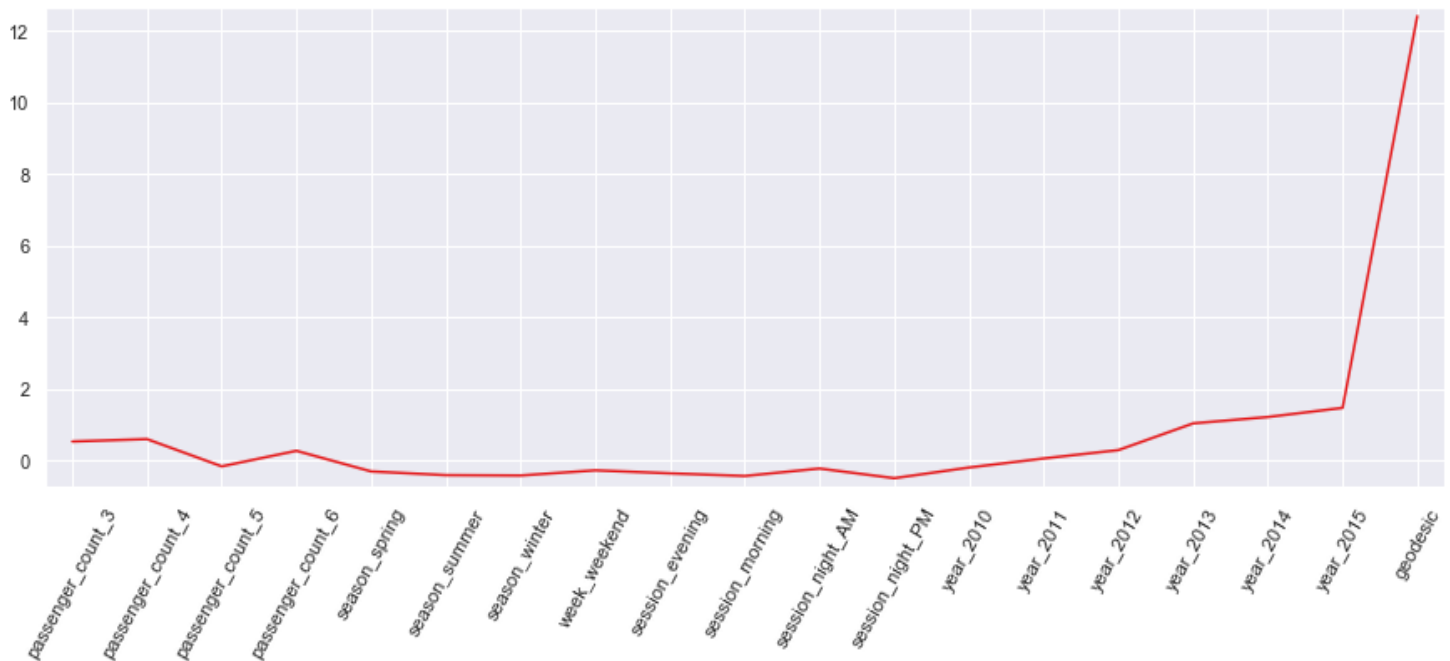
Line Plot for Coefficients of Multiple Linear regression:



Ridge Regression:

Error Metrics	r-square	Adjusted r-square	MAPE	MSE	RMSE	RMSLE
Train	0.492	0.491	24.130	8.411	2.900	0.267
Test	0.484	0.482	24.608	8.538	2.922	0.269

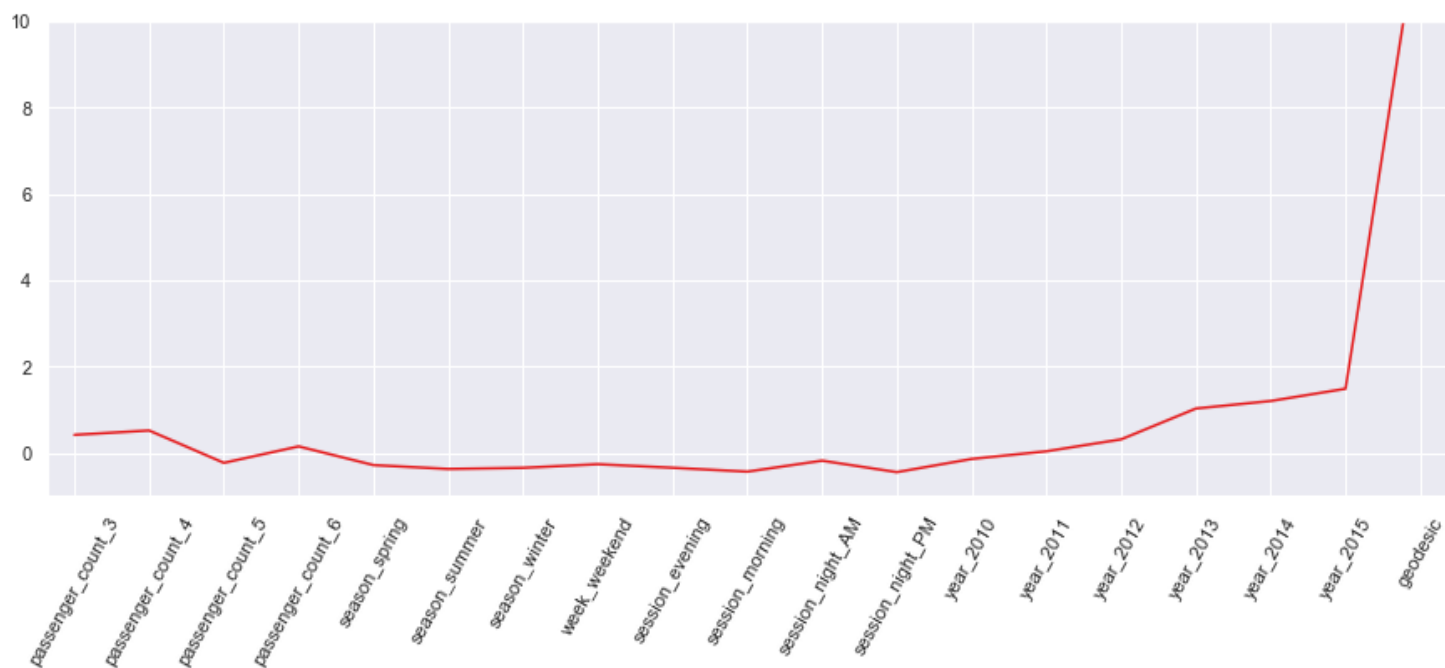
Line Plot for Coefficients of Ridge regression:



Lasso Regression:

Error Metrics	r-square	Adjusted r-square	MAPE	MSE	RMSE	RMSLE
Train	0.492	0.491	24.114	8.414	2.900	0.267
Test	0.486	0.483	24.587	8.518	2.918	0.269

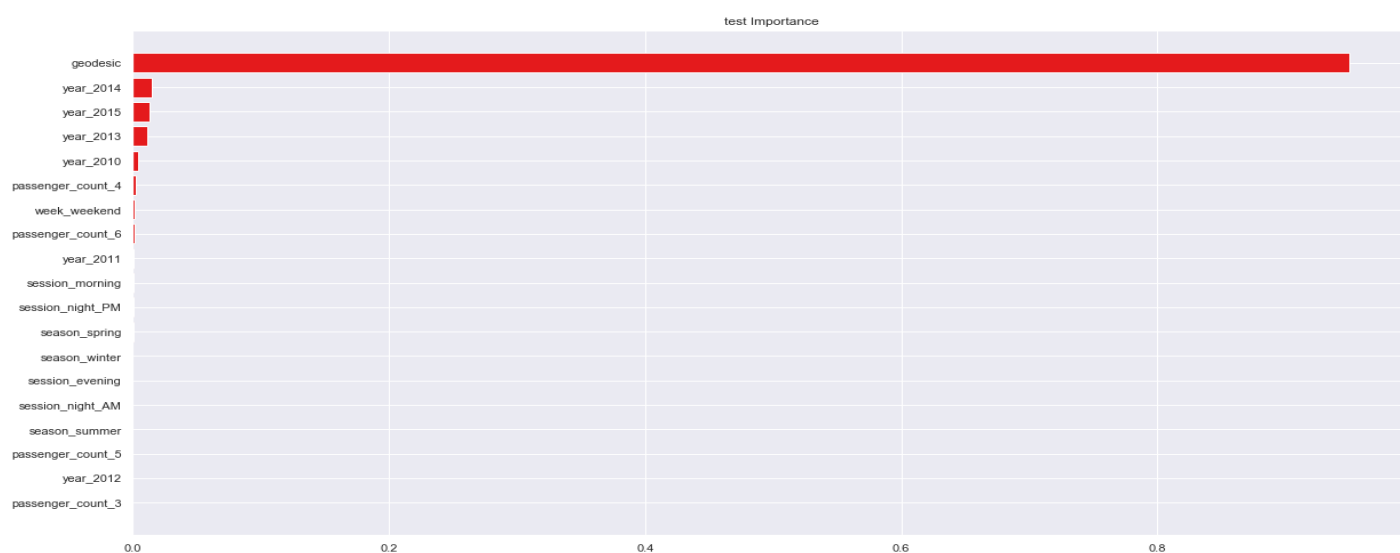
Line Plot for Coefficients of Ridge regression:



Decision Tree Regression:

Error Metrics	r-square	Adjusted r-square	MAPE	MSE	RMSE	RMSLE
Train	0.515	0.514	23.402	8.027	2.833	0.259
Test	0.480	0.477	24.449	8.614	2.935	0.268

Bar Plot of Decision tree Feature Importance:



Random Forest Regression:

Error Metrics	r-square	Adjusted r-square	MAPE	MSE	RMSE	RMSLE
Train	0.584	0.583	21.733	6.889	2.624	0.240
Test	0.495	0.493	24.062	8.362	2.891	0.265

Bar Plot of Random Forest Feature Importance:

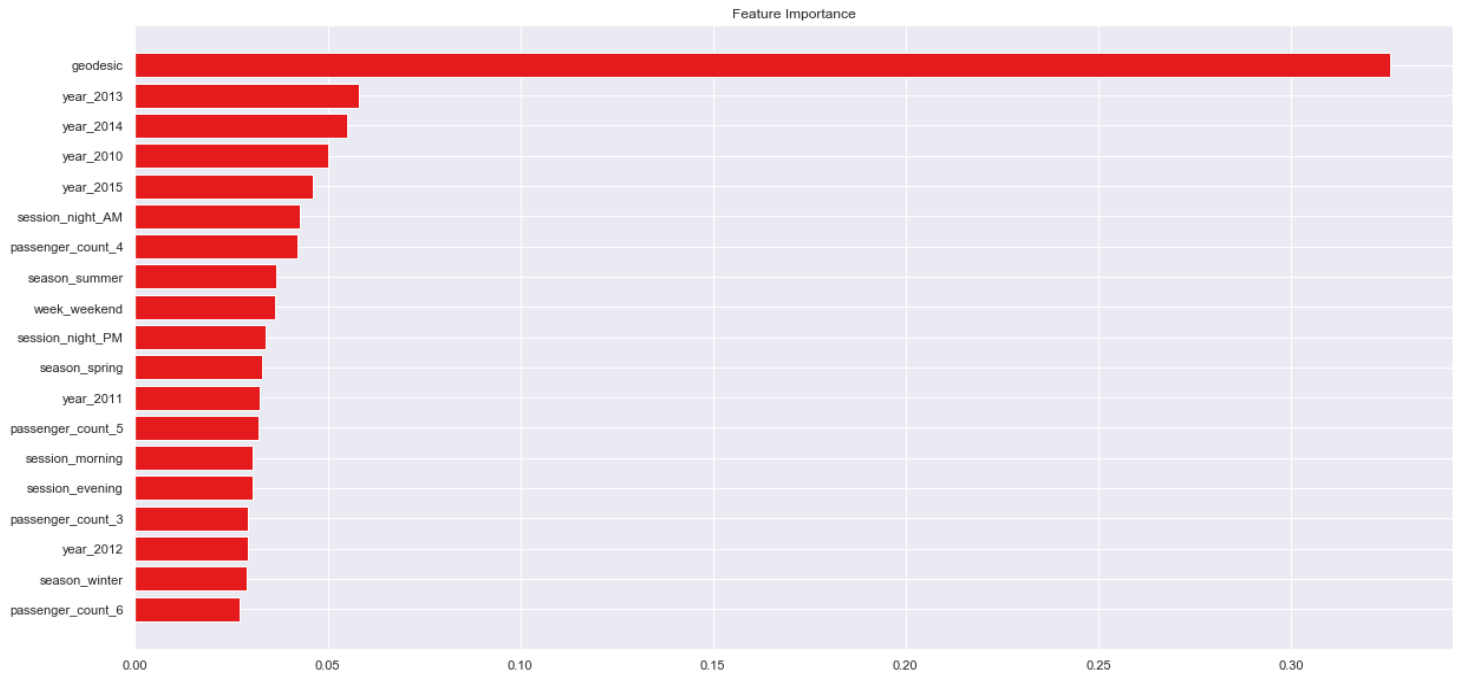
CHAPTER 6

6. Improving Accuracy

- Improve Accuracy
 - Algorithm Tuning
 - Ensembles
- We have used XgBoost as an ensemble technique.

Xgboost hyperparameters tuned parameters: Tuned Xgboost Parameters: {'subsample': 0.1, 'reg_alpha': 0.08685113737513521, 'n_estimators': 200, 'max_depth': 3, 'learning_rate': 0.05, 'colsample_bytree': 0.7000000000000001, 'colsample_bynode': 0.7000000000000001, 'colsample_bylevel': 0.9000000000000001} Best score is 0.501889441417865

Error Metrics	r-square	Adjusted r-square	MAPE	MSE	RMSE	RMSLE
Train	0.5197	0.5190	23.178	7.954	2.820	0.258
Test	0.498	0.495	23.892	8.316	2.883	0.264



CHAPTER 7

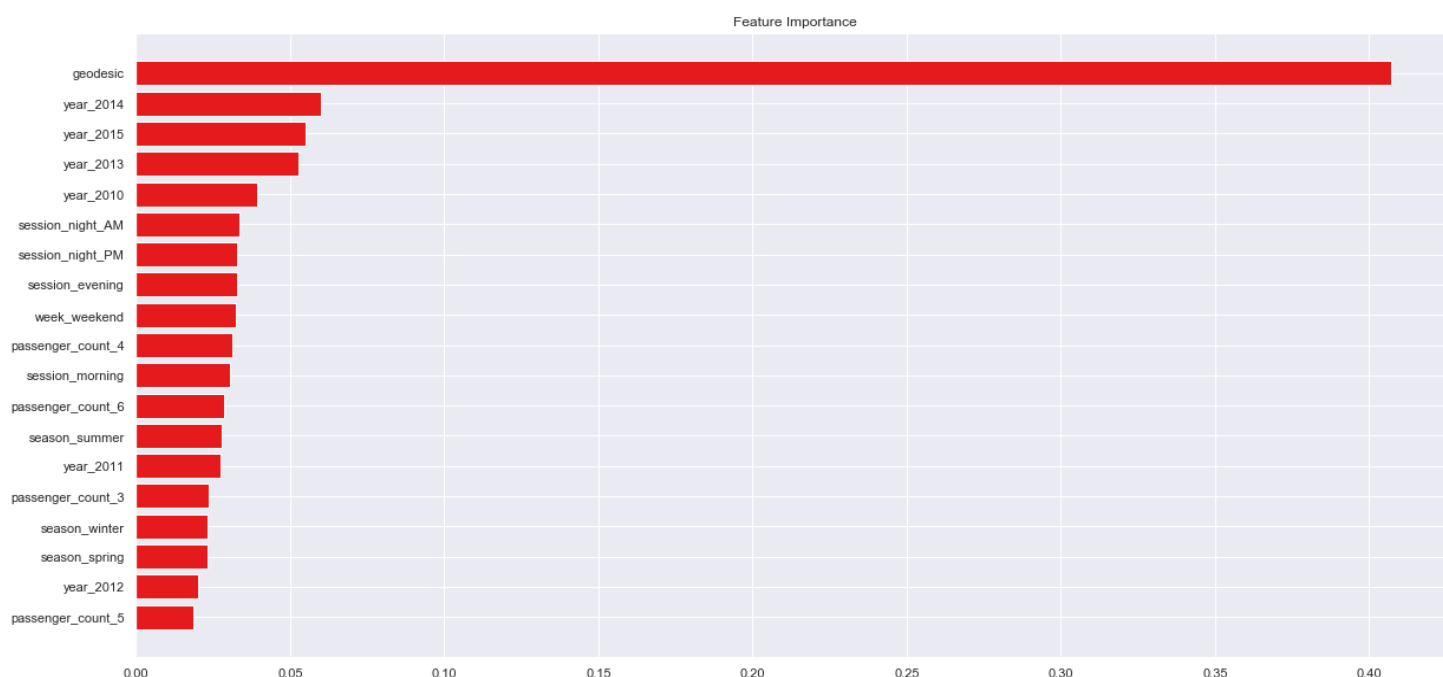
7. Finalize model

- Create standalone model on entire training dataset
- Save model for later use

We have trained a Xgboost model on entire training dataset and used that model to predict on test data. Also, we have saved model for later use.

```
<<<----- Training Data Score ----->
```

```
r_square    0.5158857617560766
Adjusted r square:0.5152976434359593
MAPE:23.45713846150603
MSE: 8.021710482552812
RMSE: 2.832262431794203
RMSLE: 0.2695533001410995
RMSLE: 0.2599544456743638
```



REFERENCE

James, Gareth, Daniela Witten, Trevor Hastie, and Robert Tibshirani. 2013. An Introduction to Statistical Learning. Vol. 6. Springer.

Wickham, Hadley. 2009. Ggplot2: Elegant Graphics for Data Analysis. Springer Science & Business Media.