

In [1]:

```

1 #Loading Libraries
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 import warnings
6 warnings.filterwarnings('ignore')
7 test['simple_mean'] = train['trip_duration'].mean()
8 from sklearn.metrics import mean_squared_error as MSE

```

In [2]:

```

1 #Loading data
2 data=pd.read_csv('nyc_taxi_trip_duration.csv')

```

In [ ]:

```

1 #Evaluation Metric used for the models is Mean Squared Error as it is most

```

In [3]:

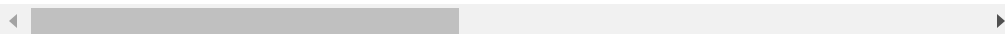
```

1 data.head()

```

Out[3]:

	id	vendor_id	pickup_datetime	dropoff_datetime	passenger_count	pi
0	id1080784	2	2016-02-29 16:40:21	2016-02-29 16:47:01	1	
1	id0889885	1	2016-03-11 23:35:37	2016-03-11 23:53:57	2	
2	id0857912	2	2016-02-21 17:59:33	2016-02-21 18:26:48	2	
3	id3744273	2	2016-01-05 09:44:31	2016-01-05 10:03:32	6	
4	id0232939	1	2016-02-17 06:42:23	2016-02-17 06:56:31	1	



In [4]:

```
1 #checking of null values
2 data.isnull().sum()
```

Out[4]:

```
id                0
vendor_id         0
pickup_datetime   0
dropoff_datetime  0
passenger_count   0
pickup_longitude  0
pickup_latitude   0
dropoff_longitude 0
dropoff_latitude  0
store_and_fwd_flag 0
trip_duration     0
dtype: int64
```

In [5]:

```
1 #dropping unnecessary coloumns
2 data.drop('id',inplace=True,axis=1)
```

In [6]:

```
1 data.drop('dropoff_datetime',inplace=True,axis=1)
```

In [7]:

```
1 data.head()
```

Out[7]:

	vendor_id	pickup_datetime	passenger_count	pickup_longitude	pickup_latitude
0	2	2016-02-29 16:40:21	1	-73.953918	40.77887
1	1	2016-03-11 23:35:37	2	-73.988312	40.73174
2	2	2016-02-21 17:59:33	2	-73.997314	40.72145
3	2	2016-01-05 09:44:31	6	-73.961670	40.75972
4	1	2016-02-17 06:42:23	1	-74.017120	40.70846

In [8]:

```
1 #Taking log for easier calculation
2 data['trip_duration']=np.log(data['trip_duration'])
```

In [9]:

```
1 data['passenger_count'].unique()
```

Out[9]:

```
array([1, 2, 6, 3, 4, 5, 0, 7, 9], dtype=int64)
```

In [10]:

```
1 data.drop(data[data['passenger_count'] ==0].index, inplace = True)
```

In [11]:

```
1 data.head()
```

Out[11]:

	vendor_id	pickup_datetime	passenger_count	pickup_longitude	pickup_latitude
0	2	2016-02-29 16:40:21	1	-73.953918	40.77887
1	1	2016-03-11 23:35:37	2	-73.988312	40.73174
2	2	2016-02-21 17:59:33	2	-73.997314	40.72145
3	2	2016-01-05 09:44:31	6	-73.961670	40.75972
4	1	2016-02-17 06:42:23	1	-74.017120	40.70846

In [12]:

```
1 data['passenger_count'].value_counts()
```

Out[12]:

```
1    517415
2    105097
5     38926
3     29692
6     24107
4     14050
7           1
9           1
```

Name: passenger\_count, dtype: int64

In [13]:

```
1 data["pickup_datetime"]=pd.to_datetime(data["pickup_datetime"])
```

In [14]:

```
1  #Extracting Month
2  data['Month'] = data['pickup_datetime'].dt.month
3
4  #Extracting day of month
5  data['DayofMonth'] = data['pickup_datetime'].dt.day
6
7  #Extracting day of week
8  data['dayofweek'] =data['pickup_datetime'].dt.dayofweek
9
10
11 #Extracting hour of the day
12 data['Hour'] = data['pickup_datetime'].dt.hour
```

In [15]:

```
1
2  # Converting yes/no flag to 1 and 0
3
4
5  data['store_and_fwd_flag']=data['store_and_fwd_flag'].map({'N':0,'Y':1})
6
```

In [16]:

```
1  import math
```

In [17]:

```
1
2 def haversine(lat1, lon1, lat2, lon2):
3
4     # distance between latitudes
5     # and longitudes
6     dLat = (lat2 - lat1) * math.pi / 180.0
7     dLon = (lon2 - lon1) * math.pi / 180.0
8
9     # convert to radians
10    lat1 = (lat1) * math.pi / 180.0
11    lat2 = (lat2) * math.pi / 180.0
12
13    # apply formulae
14    a = (pow(math.sin(dLat / 2), 2) +
15         pow(math.sin(dLon / 2), 2) *
16         math.cos(lat1) * math.cos(lat2));
17    rad = 6371
18    c = 2 * math.asin(math.sqrt(a))
19    return rad * c
```

In [18]:

```
1 #changing longitude and latitude to distance
2 data['distance']=data.apply(lambda row:haversine(row['pickup_latitude'],row
3 data['distance']=data['distance'].astype(float)
```

In [19]:

```
1 data.head()
```

Out[19]:

	vendor_id	pickup_datetime	passenger_count	pickup_longitude	pickup_latitude
0	2	2016-02-29 16:40:21	1	-73.953918	40.77887
1	1	2016-03-11 23:35:37	2	-73.988312	40.73174
2	2	2016-02-21 17:59:33	2	-73.997314	40.72145
3	2	2016-01-05 09:44:31	6	-73.961670	40.75972
4	1	2016-02-17 06:42:23	1	-74.017120	40.70846

In [20]:

```
1 #droppiong Longitude and Latitude
2 data.drop('pickup_longitude',inplace=True,axis=1)
3 data.drop('pickup_latitude',inplace=True,axis=1)
4 data.drop('dropoff_longitude',inplace=True,axis=1)
5 data.drop('dropoff_latitude',inplace=True,axis=1)
6 data.drop('pickup_datetime',inplace=True,axis=1)
```

In [21]:

```
1 data.head()
```

Out[21]:

	vendor_id	passenger_count	store_and_fwd_flag	trip_duration	Month	Dayoff
0	2	1	0	5.991465	2	
1	1	2	0	7.003065	3	
2	2	2	0	7.399398	2	
3	2	6	0	7.039660	1	
4	1	1	0	6.742881	2	

In [22]:

```
1 #getting rid of outliers
2 data.drop(data[data['distance'] > 200].index, inplace = True)
```

In [23]:

```
1 data.drop('store_and_fwd_flag',inplace=True,axis=1)
```

In [24]:

```
1 data.head()
```

Out[24]:

	vendor_id	passenger_count	trip_duration	Month	DayofMonth	dayofweek	H
0	2	1	5.991465	2	29	0	
1	1	2	7.003065	3	11	4	
2	2	2	7.399398	2	21	6	
3	2	6	7.039660	1	5	1	
4	1	1	6.742881	2	17	2	

In [25]:

```
1 data.shape
```

Out[25]:

(729283, 8)

In [26]:

```
1 from sklearn.utils import shuffle
2
3 # Shuffling the Dataset
4 data = shuffle(data, random_state = 42)
5
6 #creating 4 divisions
7 div = int(data.shape[0]/4)
8
9 # 3 parts to train set and 1 part to test set
10 train = data.loc[:3*div+1,:]
11 test = data.loc[3*div+1:]
```




In [27]:

```
1 train.head()
```

Out[27]:

	vendor_id	passenger_count	trip_duration	Month	DayofMonth	dayofwe
<b>267686</b>	1	1	6.588926	6	3	
<b>217626</b>	2	2	6.386879	6	25	
<b>540092</b>	1	1	5.407172	4	19	
<b>488315</b>	2	1	6.614726	5	27	
<b>529358</b>	2	1	5.583496	6	28	



In [28]:

```
1 train.shape
```

Out[28]:


(570718, 8)

In [29]:

```
1 test.head()
```

Out[29]:

	vendor_id	passenger_count	trip_duration	Month	DayofMonth	dayofwe
<b>546961</b>	2	1	5.590987	2	10	
<b>708968</b>	2	6	6.440947	5	24	
<b>688012</b>	1	1	6.717805	3	9	
<b>673321</b>	1	1	5.937536	2	19	
<b>556530</b>	2	2	6.302619	4	8	



In [30]:

```
1 test.shape
```

Out[30]:

(158566, 8)

In [31]:

1

In [32]:

```
1 #calculating Error for Simple Mean
2
3 simple_mean_error = MSE(test['trip_duration'] , test['simple_mean'])
4 simple_mean_error
```

Out[32]:

0.6414769712789729

In [33]:

```
1 hour= pd.pivot_table(train, values='trip_duration', index = ['Hour'], aggfun
2 hour
```

Out[33]:

trip_duration	
Hour	
0	6.419300
1	6.375759
2	6.325873
3	6.316918
4	6.319188
5	6.248041
6	6.185225
7	6.346175
8	6.462101
9	6.480601
10	6.486776
11	6.517891
12	6.517866
13	6.531802
14	6.570180
15	6.566781
16	6.543516
17	6.532548
18	6.504891
19	6.440187
20	6.416309
21	6.432378
22	6.470023
23	6.464092

In [34]:

```

1 test['hour_mean'] = 0
2
3 # For every unique entry in Hour
4 for i in train['Hour'].unique():
5     # Assign the mean value corresponding to unique entry
6     test['hour_mean'][test['Hour'] == i] = train['trip_duration'][train['Hour'] == i].mean()

```

In [35]:

```

1 #calculating mean squared error
2 hour_error = MSE(test['trip_duration'] , test['hour_mean'] )
3 hour_error

```

Out[35]:

0.6348134684473961

In [36]:

```

1 weekday= pd.pivot_table(train, values='trip_duration', index = ['dayofweek', 'weekday'])
2 weekday

```

Out[36]:

	trip_duration
dayofweek	
0	6.425561
1	6.489393
2	6.516227
3	6.523728
4	6.500830
5	6.417368
6	6.371678

In [37]:

```

1 test['weekday_mean'] = 0
2
3 # For every unique entry in dayofweek
4 for i in train['dayofweek'].unique():
5     # Assign the mean value corresponding to unique entry
6     test['weekday_mean'][test['dayofweek'] == i] = train['trip_duration'][tr

```

In [38]:

```

1 #calculating mean squared error
2 weekday_error = MSE(test['trip_duration'] , test['weekday_mean'] )
3 weekday_error

```

Out[38]:

0.6386881642731397

In [39]:

```

1 #seperating independent and dependent variables
2 x = data.drop(['trip_duration'], axis=1)
3 y = data['trip_duration']
4 x.shape, y.shape

```

Out[39]:

((729283, 7), (729283,))

In [40]:

```

1 # Importing MinMax Scaler
2 from sklearn.preprocessing import MinMaxScaler
3 scaler = MinMaxScaler()
4 x_scaled = scaler.fit_transform(x)

```

In [41]:

```

1 x = pd.DataFrame(x_scaled)

```

In [42]:

```

1 # Importing Train test split
2 from sklearn.model_selection import train_test_split
3 train_x, test_x, train_y, test_y = train_test_split(x, y, random_state = 56)

```

In [43]:

```
1 #importing KNN regressor
2 from sklearn.neighbors import KNeighborsRegressor as KNN
```

In [44]:

```
1 reg = KNN(n_neighbors = 5)
2
3 # Fitting the model
4 reg.fit(train_x, train_y)
5
6 # Predicting over the Train Set and calculating MSE
7 test_predict = reg.predict(test_x)
8 k = MSE(test_predict, test_y)
9 print('Test MSE    ', k )
```

Test MSE 0.3201870115350543

In [45]:

```
1 def Elbow(K):
2     #initiating empty list
3     test_mse = []
4
5     #training model for every value of K
6     for i in K:
7         #Instance of KNN
8         reg = KNN(n_neighbors = i)
9         reg.fit(train_x, train_y)
10        #Appending mse value to empty list calculated using the predictions
11        tmp = reg.predict(test_x)
12        tmp = MSE(tmp, test_y)
13        test_mse.append(tmp)
14
15    return test_mse
```

In [46]:

```
1 #Defining K range
2 k = range(1,11,2)
```

In [47]:

```
1 # calling above defined function
2 test = Elbow(k)
```

In [48]:

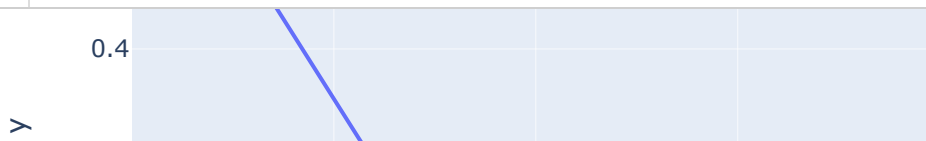
```
1 import plotly.express as px
```

In [49]:

```
1 # plotting the Curves  
2 fig=px.line(x=k,y=test)
```

In [50]:

```
1 fig.show()
```



In [51]:

```
1 reg = KNN(n_neighbors = 5)  
2 # Fitting the model  
3 reg.fit(train_x, train_y)  
4  
5 # Predicting over the Train Set and calculating MSE  
6 train_predict = reg.predict(train_x)  
7 train_k = MSE(train_predict, train_y)
```

In [52]:

```
1 print('Train MSE ', train_k)
```

Train MSE 0.20800164109882655

In [205]:

```

1 reg = KNN(n_neighbors = 5)
2 # Fitting the model
3 reg.fit(train_x, train_y)
4
5 # Predicting over the Test Set and calculating MSE
6 test_predict = reg.predict(test_x)
7 test_k = MSE(test_predict, test_y)
8 print('Test MSE ', test_k )

```

Test MSE 0.26665160784965847

In [54]:

```
1 data.head()
```

Out[54]:

	vendor_id	passenger_count	trip_duration	Month	DayofMonth	dayofwe
<b>267686</b>	1	1	6.588926	6	3	
<b>217626</b>	2	2	6.386879	6	25	
<b>540092</b>	1	1	5.407172	4	19	
<b>488315</b>	2	1	6.614726	5	27	
<b>529358</b>	2	1	5.583496	6	28	

In [181]:

```

1 #seperating independent and dependent variables
2 x = data.drop(['trip_duration'], axis=1)
3 y = data['trip_duration']
4 x.shape, y.shape

```

Out[181]:

((729283, 7), (729283,))

In [182]:

```

1 # Importing the train test split function
2 from sklearn.model_selection import train_test_split
3 train_x, test_x, train_y, test_y = train_test_split(x, y, random_state = 56)

```



In [183]:

```
1 #importing Linear Regression
2 from sklearn.linear_model import LinearRegression as LR
```

In [184]:

```
1 # Creating instance of Linear Regression
2 lr = LR()
3
4 # Fitting the model
5 lr.fit(train_x, train_y)
```

Out[184]:

LinearRegression()

In [185]:

```
1 # Predicting over the Train Set and calculating error
2 train_predict = lr.predict(train_x)
3 train_lr_k = MSE(train_predict, train_y)
4 print('Training Mean Square Error', train_lr_k )
```

Training Mean Square Error 0.38964204246525314

In [186]:

```
1 # Predicting over the Test Set and calculating error
2 test_predict = lr.predict(test_x)
3 test_lr_k = MSE(test_predict, test_y)
4 print('Test Mean Absolute Error ', test_lr_k)
```

Test Mean Absolute Error 0.38458672319254134

In [187]:

```
1 # Importing ridge from sklearn's linear_model module
2 from sklearn.linear_model import Ridge
```

In [188]:

```
1 rr = Ridge(alpha=0.1)
```

In [189]:

```
1 rr.fit(train_x, train_y)
```

Out[189]:

Ridge(alpha=0.1)

In [190]:

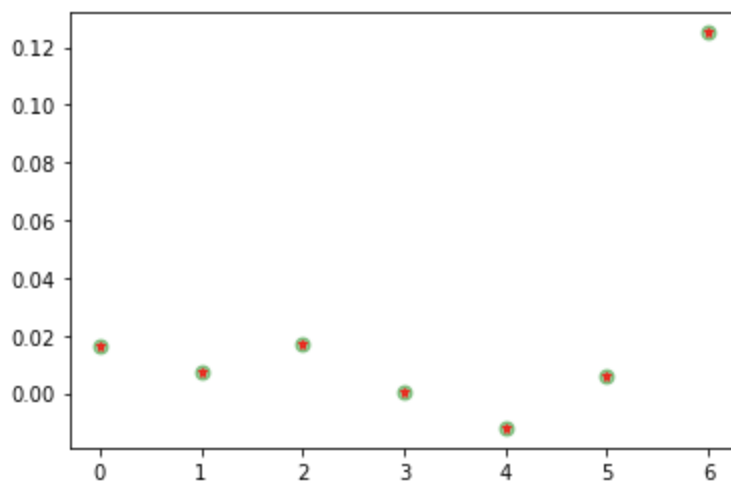
```
1 train_score=lr.score(train_x, train_y)
2 test_score=lr.score(test_x, test_y)
3 Ridge_train_score = rr.score(train_x,train_y)
4 Ridge_test_score = rr.score(test_x, test_y)
5
```

In [191]:

```
1 plt.plot(rr.coef_,alpha=0.7,linestyle='none',marker='*',markersize=5,color='red')
2 plt.plot(lr.coef_,alpha=0.4,linestyle='none',marker='o',markersize=7,color='green')
```

Out[191]:

[&lt;matplotlib.lines.Line2D at 0x160982da2e0&gt;]



In [152]:

```
1 rr.coef_
```

Out[152]:

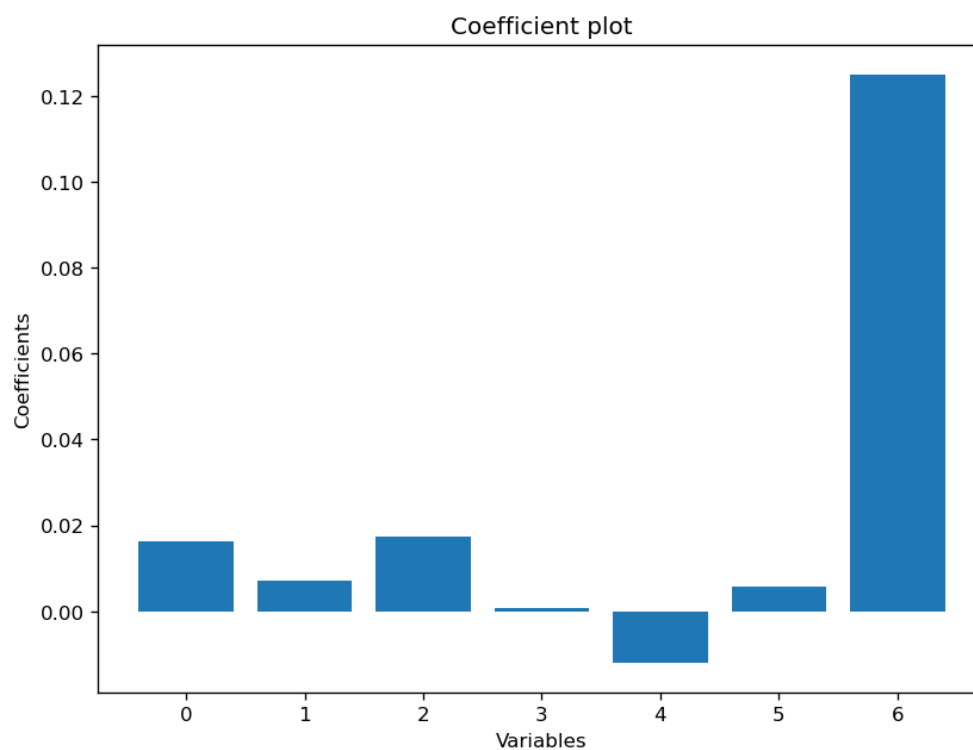
```
array([ 0.016295 ,  0.00716117,  0.01730813,  0.0006261 , -0.01
200293,
        0.00583826,  0.1250478 ])
```

In [153]:

```
1 plt.figure(figsize=(8, 6), dpi=120, facecolor='w', edgecolor='b')
2 x = range(len(train_x.columns))
3 y = lr.coef_
4 plt.bar( x, y )
5 plt.xlabel( "Variables")
6 plt.ylabel('Coefficients')
7 plt.title('Coefficient plot')
```

Out[153]:

Text(0.5, 1.0, 'Coefficient plot')



In [154]:

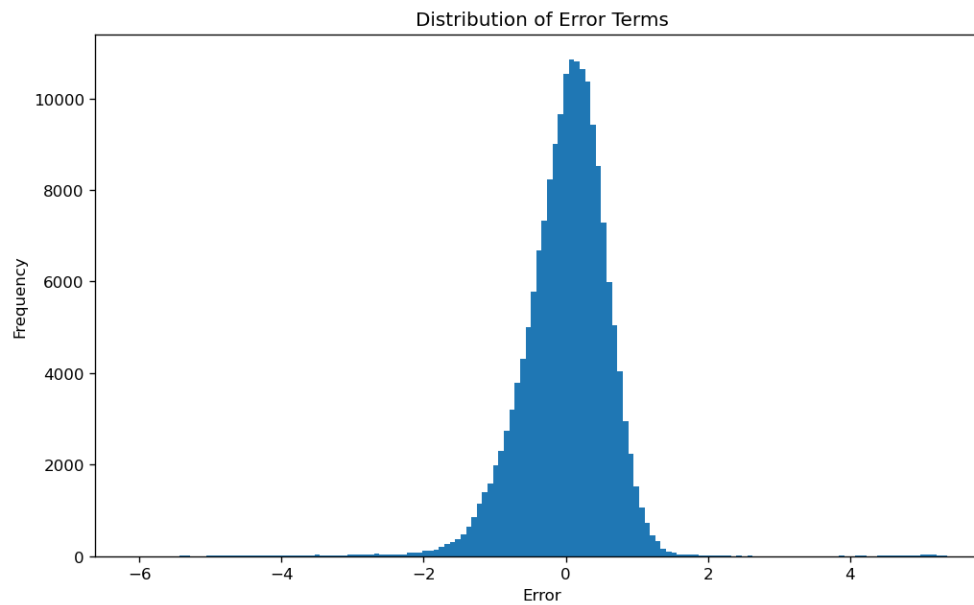
```
1 # Arranging and calculating the Residuals
2 residuals = pd.DataFrame({
3     'fitted values' : test_y,
4     'predicted values' : test_predict,
5 })
6
7 residuals['residuals'] = residuals['fitted values'] - residuals['predicted
8 residuals.head()
```

Out[154]:

	fitted values	predicted values	residuals
<b>347029</b>	6.1	6.2	-0.026
<b>397789</b>	8.2	7.2	0.93
<b>122614</b>	7.5	6.6	0.89
<b>603659</b>	5.6	6.3	-0.72
<b>259950</b>	5.1	6	-0.95

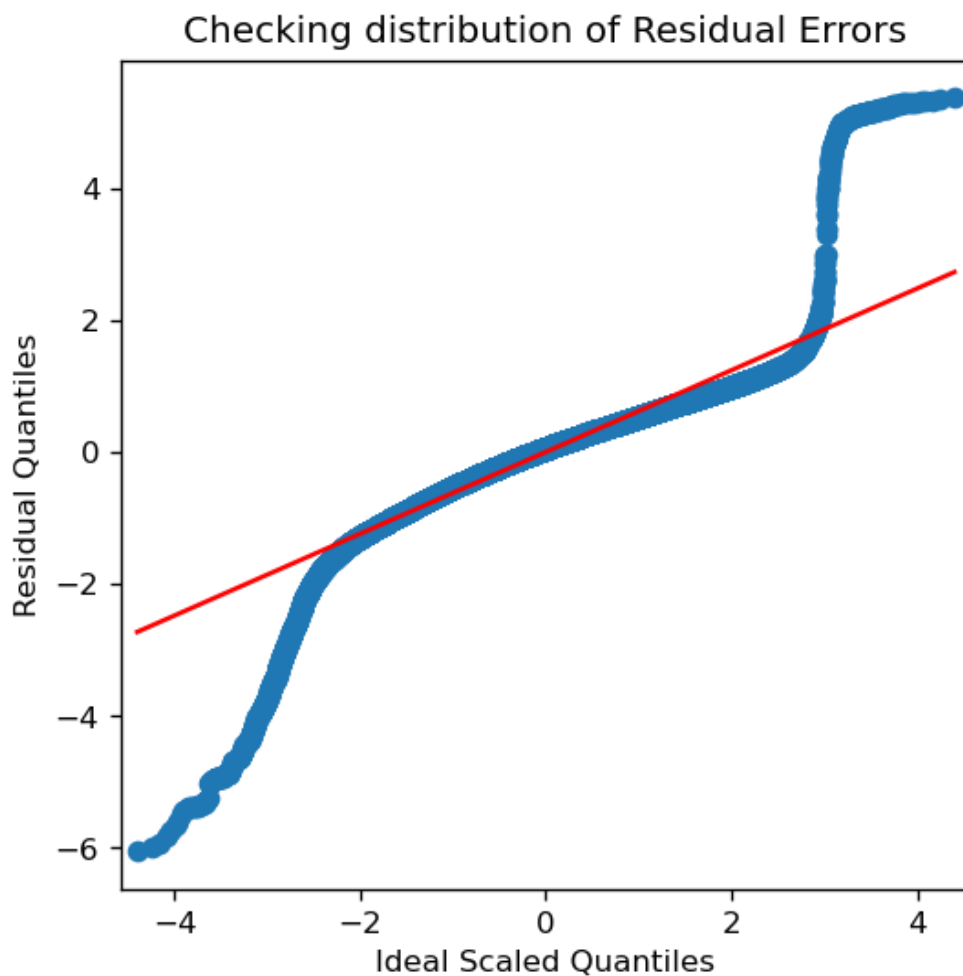
In [173]:

```
1 # Histogram for distribution
2 plt.figure(figsize=(10, 6), dpi=120, facecolor='w', edgecolor='b')
3 plt.hist(residuals.residuals, bins = 150)
4 plt.xlabel('Error')
5 plt.ylabel('Frequency')
6 plt.title('Distribution of Error Terms')
7 plt.show()
```



In [156]:

```
1 # importing the QQ-plot from the statsmodels
2 from statsmodels.graphics.gofplots import qqplot
3
4 ## Plotting the QQ plot
5 fig, ax = plt.subplots(figsize=(5,5) , dpi = 120)
6 qqplot(residuals.residuals, line = 's' , ax = ax)
7 plt.ylabel('Residual Quantiles')
8 plt.xlabel('Ideal Scaled Quantiles')
9 plt.title('Checking distribution of Residual Errors')
10 plt.show()
```



In [157]:

```
1 from statsmodels.stats.outliers_influence import variance_inflation_factor
2 from statsmodels.tools.tools import add_constant
3
4 # Calculating VIF for every column (only works for the not Catagorical)
5 VIF = pd.Series([variance_inflation_factor(data.values, i) for i in range(n)]
6 VIF
```

Out[157]:

```
vendor_id      11
passenger_count 2.8
trip_duration   25
Month           5.3
DayofMonth      4.1
dayofweek       3.3
Hour            5.4
distance        2.2
dtype: float64
```

In [158]:

```
1 lr = LR(normalize = True)
2
3 # Fitting the model
4 lr.fit(train_x, train_y)
```

Out[158]:

```
LinearRegression(normalize=True)
```

In [193]:

```
1 train_predict = lr.predict(train_x)
2 train_k_normal = MSE(train_predict, train_y)
3 print('Training Mean Absolute Error', k )
```

Training Mean Absolute Error 0.38516877731830523

In [194]:

```
1 test_predict = lr.predict(test_x)
2 test_k_normal = MSE(test_predict, test_y)
3 print('Test Mean Absolute Error', k )
```

Test Mean Absolute Error 0.38516877731830523

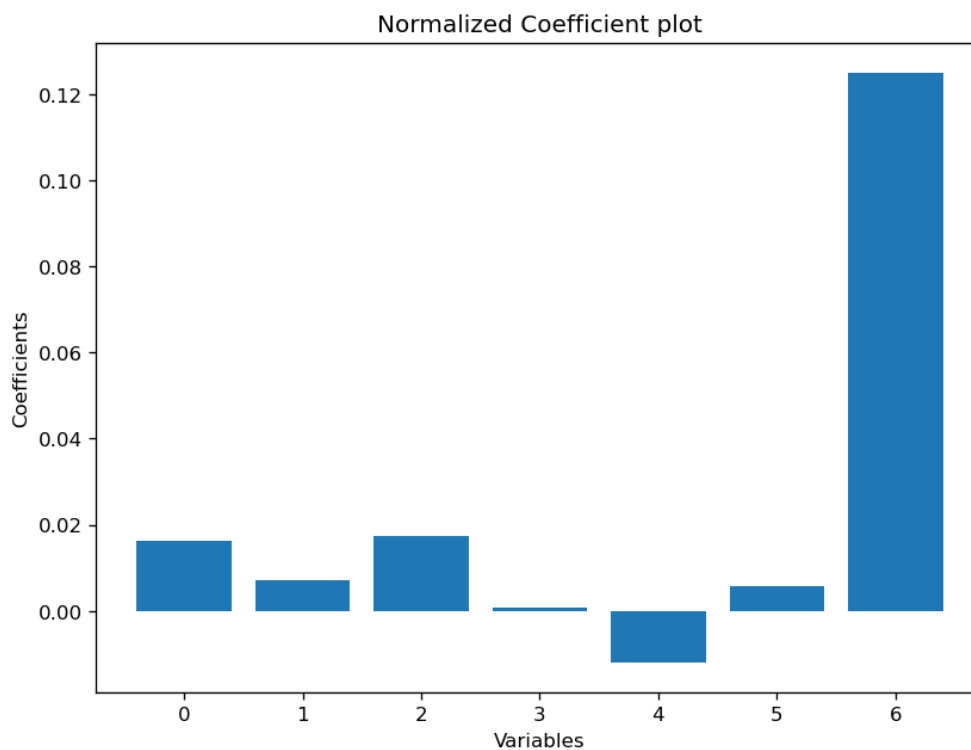


In [161]:

```
1 plt.figure(figsize=(8, 6), dpi=120, facecolor='w', edgecolor='b')
2 x = range(len(train_x.columns))
3 y = lr.coef_
4 plt.bar( x, y )
5 plt.xlabel( "Variables")
6 plt.ylabel('Coefficients')
7 plt.title('Normalized Coefficient plot')
```

Out[161]:

Text(0.5, 1.0, 'Normalized Coefficient plot')



In [162]:

```
1 x = data.drop(['trip_duration'], axis=1)
2 y = data['trip_duration']
3 x.shape, y.shape
```

Out[162]:

((729283, 7), (729283,))

In [163]:

```
1 Coefficients = pd.DataFrame({
2     'Variable' : x.columns,
3     'coefficient' : lr.coef_
4 })
5 Coefficients
```

Out[163]:

	Variable	coefficient
0	vendor_id	0.016
1	passenger_count	0.0072
2	Month	0.017
3	DayofMonth	0.00063
4	dayofweek	-0.012
5	Hour	0.0058
6	distance	0.13

In [164]:

```
1 sig_var = Coefficients[Coefficients.coefficient > 0.001]
```

In [165]:

```
1 subset = data[sig_var['Variable'].values]
2 subset.head()
```

Out[165]:

	vendor_id	passenger_count	Month	Hour	distance
267686	1	1	6	9	0.88
217626	2	2	6	11	3.5
540092	1	1	4	7	0.96
488315	2	1	5	17	1.9
529358	2	1	6	16	1.9

In [166]:

```
1 from sklearn.model_selection import train_test_split
2 train_x, test_x, train_y, test_y = train_test_split(subset, y , random_state :
```

In [167]:

```
1 from sklearn.linear_model import LinearRegression as LR
2
```

In [168]:

```
1 lr = LR(normalize = True)
2
3 # Fitting the model
4 lr.fit(train_x, train_y)
```

Out[168]:

LinearRegression(normalize=True)

In [196]:

```
1 train_predict = lr.predict(train_x)
2 coef_adj_train_k = MSE(train_predict, train_y)
3 print('Training Mean Absolute Error', k )
```

Training Mean Absolute Error 0.38516877731830523

In [197]:

```
1 test_predict = lr.predict(test_x)
2 coef_adj_test_k = MSE(test_predict, test_y)
3 print('Test Mean Absolute Error    ', k )
```

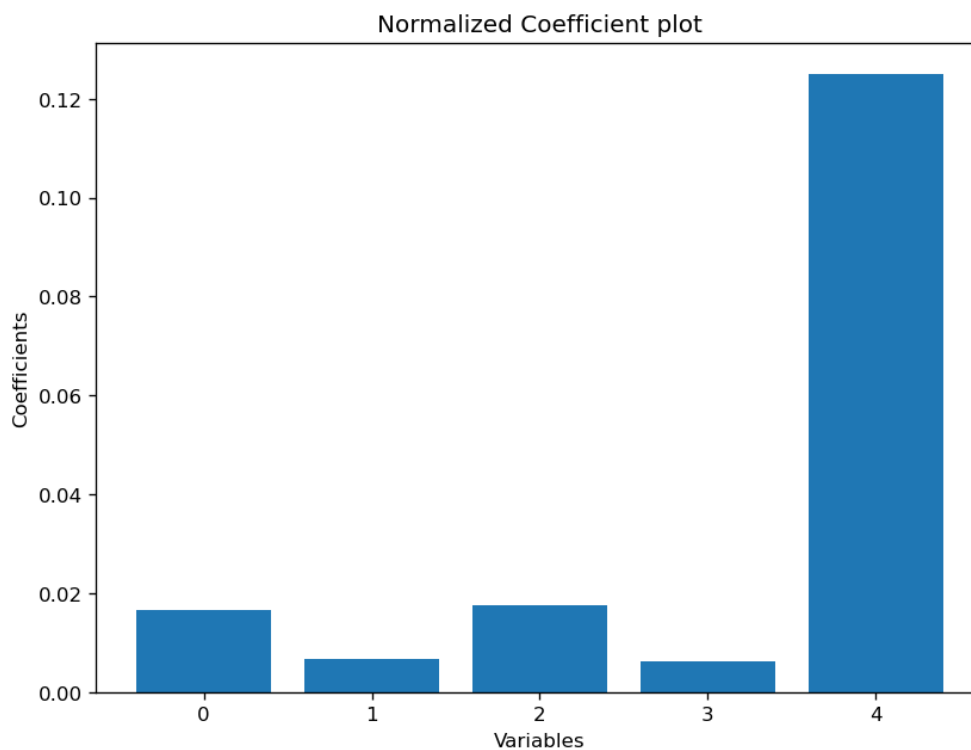
Test Mean Absolute Error 0.38516877731830523

In [171]:

```
1 plt.figure(figsize=(8, 6), dpi=120, facecolor='w', edgecolor='b')
2 x = range(len(train_x.columns))
3 y = lr.coef_
4 plt.bar( x, y )
5 plt.xlabel( "Variables")
6 plt.ylabel('Coefficients')
7 plt.title('Normalized Coefficient plot')
```

Out[171]:

Text(0.5, 1.0, 'Normalized Coefficient plot')



In [202]:

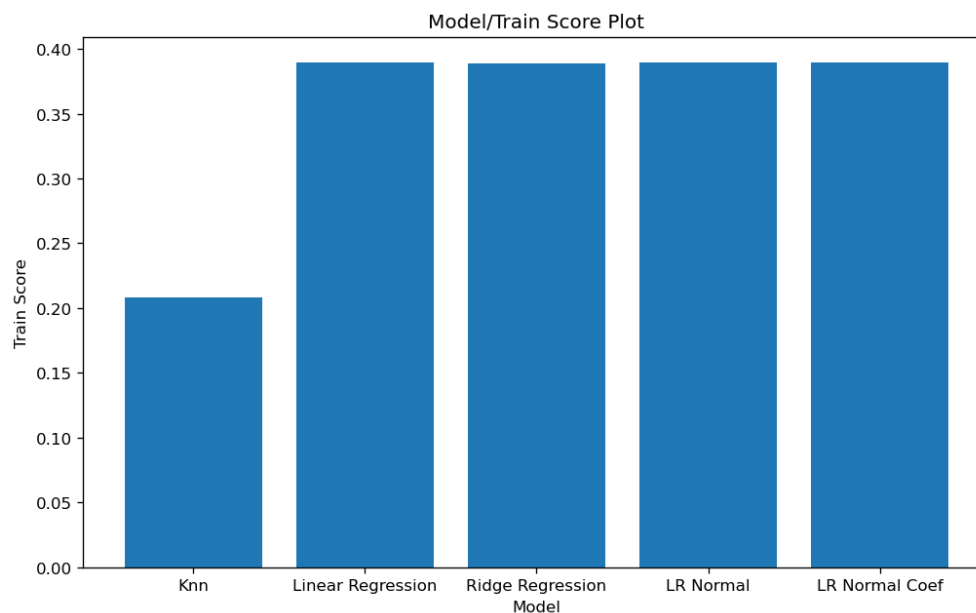
```
1 model=['Knn','Linear Regression','Ridge Regression','LR Normal','LR Normal  
2 train_score=[train_k,train_lr_k,Ridge_train_score,train_k_normal,coef_adj_t
```

In [203]:

```
1 plt.figure(figsize=(10, 6), dpi=120, facecolor='w', edgecolor='b')  
2 plt.bar( model,train_score )  
3 plt.xlabel( "Model")  
4 plt.ylabel('Train Score')  
5 plt.title('Model/Train Score Plot')
```

Out[203]:

Text(0.5, 1.0, 'Model/Train Score Plot')

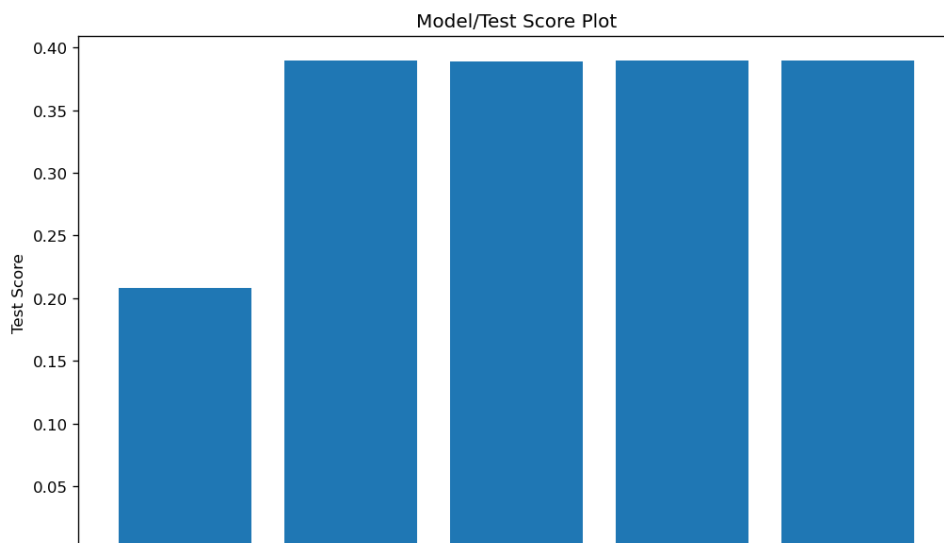


In [206]:

```
1 model=['Knn','Linear Regression','Ridge Regression','LR Normal','LR Normal
2 test_score=[test_k,test_lr_k,Ridge_test_score,test_k_normal,coef_adj_test_l
3 plt.figure(figsize=(10, 6), dpi=120, facecolor='w', edgecolor='b')
4 plt.bar( model,train_score )
5 plt.xlabel( "Model")
6 plt.ylabel('Test Score')
7 plt.title('Model/Test Score Plot')
```

Out[206]:

Text(0.5, 1.0, 'Model/Test Score Plot')



In [ ]:

1