Evalution Metric is used RMSE I have used Root mean squared error RMSE as the evaluation metric since the target variable is continous variable.RMSE is most commonly used evaluation metrics for regression problems since it brings down the unit of the differnce between predicted and actual values to the same unit as the target variable

4 id0232939

Out[9]:

In [1]: #importing the required libraries import pandas as pd import numpy as np

Build a BenchMark model for the given dataset

import matplotlib.pyplot as plt import warnings warnings.filterwarnings('ignore') df = pd.read_csv("nyc_taxi_trip_duration.csv")

df.head()

id vendor_id dropoff_datetime passenger_count pickup_longitude pickup_latitude dropoff_longitude dropoff_latitude store_and_fwd_flag trip_duration Out[2]: pickup_datetime

0 id1080784 2 2016-02-29 16:40:21 2016-02-29 16:47:01 1

-73.953918 **1** id0889885 1 2016-03-11 23:35:37 2016-03-11 23:53:57 **2** id0857912 2 2 2016-02-21 17:59:33 2016-02-21 18:26:48

40.778873

40.731743

40.721458

40.759720

40.708469

40.762035

40.747677

40.762791

40.792503

40.732590

40.768059

40.766018

40.793781

40.798557

-73.972267

-73.982704

-73.973518

-73.988609

-74.000526

-73.967606

-74.004219

-73.976006

-73.979736

40.781265

40.741161

40.762909

40.758369

40.742283

40.763073

40.742523

40.792339

40.777878

0.182500

0.054167

0.110556

0.189444

0.096111

0.066667

0.186389

0.011944

0.090278

Ν

Ν

Ν

Ν

Ν

Ν

-74.017120

-73.963875

-73.994751

-73.948029

-73.956779

-73.988182

40.771164

40.694931

40.774918

40.780628

40.740631

Ν

Ν

Ν

Ν

Ν

400

1100

1635

1141

848

-73.988312 -73.997314 **3** id3744273 2 2016-01-05 09:44:31 2016-01-05 10:03:32 6 -73.961670

1 2016-02-17 06:42:23 2016-02-17 06:56:31 1 df['trip_duration_hour'] = df['trip_duration'].apply(lambda x: x/3600) df.drop(columns=['trip_duration'], inplace=True)

In [4]: | df['pickup_datetime'] = pd.to_datetime(df['pickup_datetime'], format = '%Y-%m-%d %H:%M:%S') df['dropoff_datetime'] = pd.to_datetime(df['dropoff_datetime'], format = '%Y-%m-%d %H:%M:%S')

df['pickup_hour'] = df['pickup_datetime'].dt.hour df['pickup_weekday'] = df['pickup_datetime'].dt.weekday df['pickup_day'] = df['pickup_datetime'].dt.day

df['pickup_month'] = df['pickup_datetime'].dt.month df['pickup_year'] = df['pickup_datetime'].dt.year df['pickup_date'] = df['pickup_datetime'].dt.date df['dropoff_hour'] = df['dropoff_datetime'].dt.hour

df['dropoff_weekday'] = df['dropoff_datetime'].dt.weekday df['dropoff_day'] = df['dropoff_datetime'].dt.day df['dropoff_month'] = df['dropoff_datetime'].dt.month df['dropoff_year'] = df['dropoff_datetime'].dt.year df['dropoff_date'] = df['dropoff_datetime'].dt.date In [6]: def time_of_day(x): # to calculate what time of it is now

if x in range(6,12): return 'Morning'

elif x in range(12,16):

return 'Afternoon' elif x in range(16,22): return 'Evening'

else: return 'Late night'

df['pickup_time_of_day'] = df['pickup_hour'].apply(time_of_day) df['dropoff_time_of_day'] = df['dropoff_hour'].apply(time_of_day) df.drop(columns=['pickup_hour','pickup_weekday','pickup_day','pickup_month','pickup_year','dropoff_hour','dropoff_weekday','dropoff_day','dropoff_month In [9]: #simple predictive model df["trip_duration_hour_mean"]=df["trip_duration_hour"].mean() df["trip_duration_hour_mean"].head()

0.264508 1 0.264508 2 0.264508 3 0.264508 0.264508 Name: trip_duration_hour_mean, dtype: float64 plt.figure(dpi=100) k = range(0, len(df))

plt.scatter(k,df["trip_duration_hour"].sort_values(),color="blue",label="actual trip duration hour") plt.plot(k,df["trip_duration_hour_mean"].sort_values(),color ="purple",label="predicted trip duration hour") plt.xlabel("fitted points (ascending)") plt.ylabel("trip duration hour") plt.title("overall mean") plt.legend() <matplotlib.legend.Legend at 0x1bc1dc98dc0> Out[17]: overall mean actual trip duration hour predicted trip duration hour 500 400 300

trip duration hour 200 100 0 100000 200000 300000 400000 500000 600000 700000 fitted points (ascending) From the above simple predictive model we can draw an insight that the values of actual trip duration and predicted trip duration are almost same. #importing the shuffle library from sklearn.utils import shuffle

Shuffling the Dataset df = shuffle(df, random_state = 42) #creating 4 divisions div = int(df.shape[0]/4)

In [10]: # 3 parts to train set and 1 part to test set train = df.loc[:3*div+1,:]test = df.loc[3*div+1:]train.head() In [11]: Out[11]: 2016-05-21 2016-05-21

id vendor_id pickup_datetime dropoff_datetime passenger_count pickup_longitude pickup_latitude dropoff_longitude dropoff_latitude store_and_fwd_flag trip_duration_hour 469114 id2380741 -73.981796 10:40:14 10:51:11 2016-01-08 2016-01-08 694852 id3946961 -73.980965 18:49:27 18:52:42 2016-05-22 2016-05-22 696324 id0833913 01:08:10 2016-06-11 2016-06-11 356496 id1336849 -73.987625 10:32:12 10:38:50 2016-04-03 2016-04-03 645318 id1610858 3 -73.964333 10:45:51 10:57:13 test.head() In [12]: id vendor_id pickup_datetime dropoff_datetime passenger_count pickup_longitude pickup_latitude dropoff_longitude dropoff_latitude store_and_fwd_flag trip_duration_hour Out[12]: 2016-05-25 2016-05-25 **546991** id2240736 1 -73.991364 07:59:16 08:05:02

2016-01-18 2016-01-18 **43126** id1423404 2 -73.966225 12:21:13 12:17:13 2016-03-02 2016-03-02 641450 id1317268 1 -73.994926 18:39:01 18:50:12 2016-04-06 2016-04-06 611380 id3335546 -73.974388 19:17:20 19:18:03 2016-06-21 2016-06-21 62690 id2174190 2 3 -73.963440 18:40:56 18:35:31

storing simple mean in a new column in the test set as "simple_mean" df["trip_duration_hour_mean"]=df["trip_duration_hour"].mean() df["trip_duration_hour_mean"] 469114 0.264508 Out[13]: 694852 0.264508 696324 0.264508 356496 0.264508 645318 0.264508 0.264508 259178 365838 0.264508 131932 0.264508

Name: trip_duration_hour_mean, Length: 729322, dtype: float64

#trip duration mean with respect to the mean of pickup time of the day

trip_mean_error = sqrt((mse(test['trip_duration_hour'] , test['trip_duration_hour_mean'])))

pickup = pd.pivot_table(train, values='trip_duration_hour', index = ['pickup_time_of_day'], aggfunc=np.mean)

test['pickup'][test['pickup_time_of_day'] == str(i)] = train['trip_duration_hour'][train['pickup_time_of_day'] == str(i)].mean()

test['dropoff'][test['dropoff_time_of_day'] == str(i)] = train['trip_duration_hour'][train['dropoff_time_of_day'] == str(i)].mean()

test['pass_count'][test['passenger_count'] == str(i)] = train['trip_duration_hour'][train['passenger_count'] == str(i)].mean()

test['store_and_fwd'][test['store_and_fwd_flag'] == str(i)] = train['trip_duration_hour'][train['store_and_fwd_flag'] == str(i)].mean()

combo = pd.pivot_table(train, values = 'trip_duration_hour', index = ['passenger_count', 'pickup_time_of_day', 'dropoff_time_of_day'], aggfunc = np.mean)

• THe error of simple mean of trip duration hour is 0.89040676 which is also equal to rmse of dropp time and pickup time of the day where as the str_fwd_error is 0.89040207 there is

In [31]: store_and_fwd = pd.pivot_table(train, values='trip_duration_hour', index = ["store_and_fwd_flag"], aggfunc=np.mean)

trip_duration_hour

0.305417

0.054352

0.106944

0.023856

0.432222

0.282054

0.313738 0.376369

16.731944 0.267228

The RMSE of the trip duration hour mean with respect to the mean of pickup time of the day is higher than the simple mean of the trip duration hour

dropoff =pd.pivot_table(train, values='trip_duration_hour', index = ['dropoff_time_of_day'], aggfunc=np.mean)

pass_count = pd.pivot_table(train, values='trip_duration_hour', index = ["passenger_count"], aggfunc=np.mean)

from sklearn.metrics import mean_squared_error as mse

trip_duration_hour

For every unique entry in pickup longitude for i in train['pickup_time_of_day'].unique():

0.291531

0.264078 0.255589

0.250610

Assign the mean value corresponding to unique entry

pickup_error = sqrt(mse(test['trip_duration_hour'] , test['pickup']))

In [21]: #trip duration mean with respect to the mean of dropoff time of the day

trip_duration_hour

0.284258

0.269666 0.256343

0.247905

Assign the mean value corresponding to unique entry

In [24]: ##trip duration mean with respect to the mean of passenger count

dropoff_error = sqrt(mse(test['trip_duration_hour'] , test['dropoff']))

0.264508

0.264508

#calculating RMSE

trip_mean_error

pickup_time_of_day

Afternoon

Evening

Late night Morning

test['pickup'] = 0

In [18]: #calculating RMSE

dropoff

dropoff_time_of_day

Afternoon

Evening

Late night

Morning

test['dropoff'] = 0

In [23]: #calculating mean absolute error

dropoff_error

pass_count

passenger_count

0

2

6

In [29]: # initializing new column to zero test['pass_count'] = 0

pass_count_error

store_and_fwd

store_and_fwd_flag

0.9290781075032716

The pass count error is more than both pickup and dropoff error

initializing new column to zero

test['store_and_fwd'] = 0

0

6

75 rows × 1 columns

conclusions

only slightly deferrence and

str_and_fwd_error 0.8904020727484228

combo

0.8903865469170373

initializing new column to zero

For every unique entry in pickup latitude for i in train['dropoff_time_of_day'].unique():

The pick up time error is same as droppoff time error

trip_duration_hour

0.092981 0.255343

0.277822 0.287332 0.285759 0.299641

0.300193

Assign the mean value corresponding to unique entry

pass_count_error = sqrt(mse(test['trip_duration_hour'] , test['pass_count']))

For every unique entry in passenger count for i in train['passenger_count'].unique():

trip_duration_hour

For every unique entry in pickup latitude for i in train['store_and_fwd_flag'].unique():

passenger_count pickup_time_of_day dropoff_time_of_day

Afternoon

Evening

Late night

Morning

Late night

Morning

• The rmse of passanger count is little high 0.92907810

The pickup time error, dropoff time error and pass_count error 0.89040676

0.264109 0.304058

Assign the mean value corresponding to unique entry

In [33]: str_and_fwd_error = sqrt(mse(test['store_and_fwd'] , test['trip_duration_hour']))

Afternoon

Evening

Late night Late night

Afternoon

Late night

Morning

Afternoon Late night

Morning

pickup_error

0.8903816016040094

initializing new column to zero

0.8904067655425832

from math import sqrt

671155 121958

pickup

In [14]:

Out[14]:

Out[15]:

In [16]:

Out[18]:

Out[21]:

Out[23]:

Out[24]:

Out[30]:

Out[31]:

Out[33]:

Out[34]: