```
1 from google.colab import drive
2 drive.mount('/content/drive')
```

Go to this URL in a browser: https://accounts.google.com/o/oauth2/auth?client_id=9473189

Enter your authorization code:
..........
Mounted at /content/drive

```
1 import os
2 os.chdir("/content/drive/My Drive/NYTP")
3 !ls -l
```

```
total 7327210
-rw------- 1 root root   23895892 Jul 30 17:49  df_test.pkl
-rw------- 1 root root   55749012 Jul 30 17:49  df_train.pkl
-rw------- 1 root root      13680 Aug  1 23:59 'kmeans.cluster_centers_ .pkl'
-rw------- 1 root root     150151 Jul 26 15:33  mydask.png
-rw------- 1 root root   15181864 Aug  1 23:59  regions_cum.pkl
-rw------- 1 root root    4707106 Jul 30 17:49  tsne_test_output.pkl
-rw------- 1 root root   10498836 Jul 30 17:49  tsne_train_output.pkl
-r-------- 1 root root 1985964692 Mar  1  2018  yellow_tripdata_2015-01.csv
-r-------- 1 root root 1708674492 Mar  1  2018  yellow_tripdata_2016-01.csv
-r-------- 1 root root 1783554554 Mar  1  2018  yellow_tripdata_2016-02.csv
-r-------- 1 root root 1914669757 Mar  1  2018  yellow_tripdata_2016-03.csv
```

# Taxi demand prediction in New York City



```
1 !pip install gpxpy
```

Collecting gpxpy
  Downloading https://files.pythonhosted.org/packages/dd/23/a1c04fb3ea8d57d4b46cf2956c99
     |████████████████████████████████| 112kB 2.8MB/s
Building wheels for collected packages: gpxpy
  Building wheel for gpxpy (setup.py) ... done
  Created wheel for gpxpy: filename=gpxpy-1.4.2-cp36-none-any.whl size=42546 sha256=3b5e
  Stored in directory: /root/.cache/pip/wheels/d9/df/ed/b52985999b3967fa0ef8de22b3dc8ad3
Successfully built gpxpy
Installing collected packages: gpxpy
Successfully installed gpxpy-1.4.2

```
1 #Importing Libraries
2 # pip3 install graphviz
3 #pip3 install dask
```

```
 4 #pip3 install toolz
 5 #pip3 install cloudpickle
 6 # https://www.youtube.com/watch?v=ieW3G7ZzRZ0
 7 # https://github.com/dask/dask-tutorial
 8 # please do go through this python notebook: https://github.com/dask/dask-tutorial/blob/ma
 9 import dask.dataframe as dd#similar to pandas
10
11 import pandas as pd#pandas to create small dataframes
12
13 # pip3 install foliun
14 # if this doesnt work refere install_folium.JPG in drive
15 import folium #open street map
16
17 # unix time: https://www.unixtimestamp.com/
18 import datetime #Convert to unix time
19
20 import time #Convert to unix time
21
22 # if numpy is not installed already : pip3 install numpy
23 import numpy as np#Do aritmetic operations on arrays
24
25 # matplotlib: used to plot graphs
26 import matplotlib
27 # matplotlib.use('nbagg') : matplotlib uses this protocall which makes plots more user int
28 matplotlib.use('nbagg')
29 import matplotlib.pylab as plt
30 import seaborn as sns#Plots
31 from matplotlib import rcParams#Size of plots
32
33 # this lib is used while we calculate the stight line distance between two (lat,lon) pairs
34 import gpxpy.geo #Get the haversine distance
35
36 from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
37 import math
38 import pickle
39 import os
40
41 # download migwin: https://mingw-w64.org/doku.php/download/mingw-builds
42 # install it in your system and keep the path, migw_path ='installed path'
43 mingw_path = 'C:\\Program Files\\mingw-w64\\x86_64-5.3.0-posix-seh-rt_v4-rev0\\mingw64\\bi
44 os.environ['PATH'] = mingw_path + ';' + os.environ['PATH']
45
46 # to install xgboost: pip3 install xgboost
47 # if it didnt happen check install_xgboost.JPG
48 import xgboost as xgb
49
50 # to install sklearn: pip install -U scikit-learn
51 from sklearn.ensemble import RandomForestRegressor
52 from sklearn.metrics import mean_squared_error
53 from sklearn.metrics import mean_absolute_error
54 import warnings
55
```

```
55 warnings.filterwarnings("ignore")
```

## ▾ Data Information

Ge the data from : http://[www.nyc.gov/html/tlc/html/about/trip_record_data.shtml](http://www.nyc.gov/html/tlc/html/about/trip_record_data.shtml) (2016 data) The data used in the attached datasets were collected and provided to the NYC Taxi and Limousine Commission (TLC)

## Information on taxis:

Yellow Taxi: Yellow Medallion Taxicabs

These are the famous NYC yellow taxis that provide transportation exclusively through street-hails. The number of taxicabs is limited by a finite number of medallions issued by the TLC. You access this mode of transportation by standing in the street and hailing an available taxi with your hand. The pickups are not pre-arranged.

For Hire Vehicles (FHVs)

FHV transportation is accessed by a pre-arrangement with a dispatcher or limo company. These FHVs are not permitted to pick up passengers via street hails, as those rides are not considered pre-arranged.

Green Taxi: Street Hail Livery (SHL)

The SHL program will allow livery vehicle owners to license and outfit their vehicles with green borough taxi branding, meters, credit card machines, and ultimately the right to accept street hails in addition to pre-arranged rides.

Credits: Quora

Footnote:

In the given notebook we are considering only the yellow taxis for the time period between Jan - Mar 2015 & Jan - Mar 2016

## ▾ Data Collection

We Have collected all yellow taxi trips data from jan-2015 to dec-2016(Will be using only 2015 data)

| file name | file name size | number of records | number of features |
|---|---|---|---|
| yellow_tripdata_2016-01 | 1. 59G | 10906858 | 19 |
| yellow_tripdata_2016-02 | 1. 66G | 11382049 | 19 |
| yellow_tripdata_2016-03 | 1. 78G | 12210952 | 19 |

| | | | |
|---|---|---|---|
| yellow_tripdata_2016-04 | 1. 74G | 11934338 | 19 |
| yellow_tripdata_2016-05 | 1. 73G | 11836853 | 19 |
| yellow_tripdata_2016-06 | 1. 62G | 11135470 | 19 |
| yellow_tripdata_2016-07 | 884Mb | 10294080 | 17 |
| yellow_tripdata_2016-08 | 854Mb | 9942263 | 17 |
| yellow_tripdata_2016-09 | 870Mb | 10116018 | 17 |
| yellow_tripdata_2016-10 | 933Mb | 10854626 | 17 |
| yellow_tripdata_2016-11 | 868Mb | 10102128 | 17 |
| yellow_tripdata_2016-12 | 897Mb | 10449408 | 17 |
| yellow_tripdata_2015-01 | 1.84Gb | 12748986 | 19 |
| yellow_tripdata_2015-02 | 1.81Gb | 12450521 | 19 |
| yellow_tripdata_2015-03 | 1.94Gb | 13351609 | 19 |
| yellow_tripdata_2015-04 | 1.90Gb | 13071789 | 19 |
| yellow_tripdata_2015-05 | 1.91Gb | 13158262 | 19 |
| yellow_tripdata_2015-06 | 1.79Gb | 12324935 | 19 |
| yellow_tripdata_2015-07 | 1.68Gb | 11562783 | 19 |
| yellow_tripdata_2015-08 | 1.62Gb | 11130304 | 19 |
| yellow_tripdata_2015-09 | 1.63Gb | 11225063 | 19 |
| yellow_tripdata_2015-10 | 1.79Gb | 12315488 | 19 |
| yellow_tripdata_2015-11 | 1.65Gb | 11312676 | 19 |
| yellow_tripdata_2015-12 | 1.67Gb | 11460573 | 19 |

```
1 #Looking at the features
2 # dask dataframe  : # https://github.com/dask/dask-tutorial/blob/master/07_dataframe.ipynb
3 month = dd.read_csv('yellow_tripdata_2015-01.csv')
4 print(month.columns)
```

```
Index(['VendorID', 'tpep_pickup_datetime', 'tpep_dropoff_datetime',
       'passenger_count', 'trip_distance', 'pickup_longitude',
       'pickup_latitude', 'RateCodeID', 'store_and_fwd_flag',
       'dropoff_longitude', 'dropoff_latitude', 'payment_type', 'fare_amount',
       'extra', 'mta_tax', 'tip_amount', 'tolls_amount',
       'improvement_surcharge', 'total_amount'],
      dtype='object')
```

```
1 # However unlike Pandas, operations on dask.dataframes don't trigger immediate computation
2 # instead they add key-value pairs to an underlying Dask graph. Recall that in the diagram
3 # circles are operations and rectangles are results.
4
5 # to see the visulaization you need to install graphviz
6 # pip3 install graphviz if this doesnt work please check the install_graphviz.jpg in the d
7 month.visualize()
```
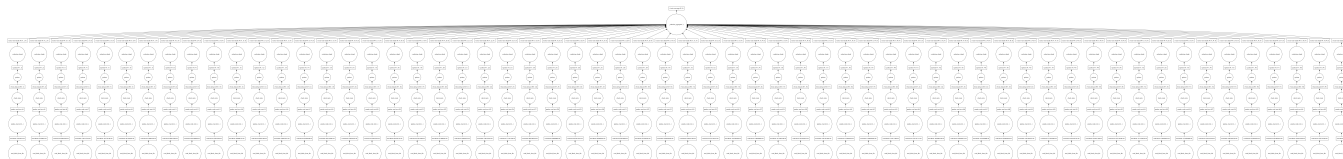
```
1 month.fare_amount.sum().visualize()
```



# Features in the dataset:

```
<tr>
    <td>Dropoff_longitude</td>
    <td>Longitude where the meter was disengaged.</td>
</tr>
<tr>
    <td>Dropoff_ latitude</td>
    <td>Latitude where the meter was disengaged.</td>
</tr>
<tr>
    <td>Payment_type</td>
    <td>A numeric code signifying how the passenger paid for the trip.
    <ol>
        <li> Credit card </li>
        <li> Cash </li>
        <li> No charge </li>
        <li> Dispute</li>
        <li> Unknown </li>
        <li> Voided trip</li>
    </ol>
    </td>
</tr>
<tr>
    <td>Fare_amount</td>
    <td>The time-and-distance fare calculated by the meter.</td>
```

```
    </tr>
    <tr>
        <td>Extra</td>
        <td>Miscellaneous extras and surcharges. Currently, this only includes. the $0.50 and $1 rush hou
    </tr>
    <tr>
        <td>MTA_tax</td>
        <td>0.50 MTA tax that is automatically triggered based on the metered rate in use.</td>
    </tr>
    <tr>
        <td>Improvement_surcharge</td>
        <td>0.30 improvement surcharge assessed trips at the flag drop. the improvement surcharge began b
    </tr>
    <tr>
        <td>Tip_amount</td>
        <td>Tip amount – This field is automatically populated for credit card tips.Cash tips are not inc
    </tr>
    <tr>
        <td>Tolls_amount</td>
        <td>Total amount of all tolls paid in trip.</td>
    </tr>
    <tr>
        <td>Total_amount</td>
        <td>The total amount charged to passengers. Does not include cash tips.</td>
    </tr>
```

| Field Name | Description |
| --- | --- |
| VendorID | A code indicating the TPEP provider that provided the record.<br>1. Creative Mobile Technologies<br>2. VeriFone Inc. |
| tpep_pickup_datetime | The date and time when the meter was engaged. |
| tpep_dropoff_datetime | The date and time when the meter was disengaged. |
| Passenger_count | The number of passengers in the vehicle. This is a driver-entered value. |
| Trip_distance | The elapsed trip distance in miles reported by the taximeter. |
| Pickup_longitude | Longitude where the meter was engaged. |
| Pickup_latitude | Latitude where the meter was engaged. |
| RateCodeID | The final rate code in effect at the end of the trip.<br>1. Standard rate<br>2. JFK<br>3. Newark<br>4. Nassau or Westchester<br>5. Negotiated fare<br>6. Group ride |

| Store_and_fwd_flag | This flag indicates whether the trip record was held in vehicle memory before sending to the vendor, aka "store and forward," because the vehicle did not have a connection to the server. Y= store and forward trip N= not a store and forward trip |
|---|---|

# ML Problem Formulation

**Time-series forecasting and Regression**

*- To find number of pickups, given location cordinates(latitude and longitude) and time, in the query reigion and surrounding regions.*

To solve the above we would be using data collected in Jan - Mar 2015 to predict the pickups in Jan - Mar 2016.

## Performance metrics

1. Mean Absolute percentage error.
2. Mean Squared error.

## Data Cleaning

In this section we will be doing univariate analysis and removing outlier/illegitimate values which may be caused due to some error

```
1 #table below shows few datapoints along with all our features
2 month.head(5)
```
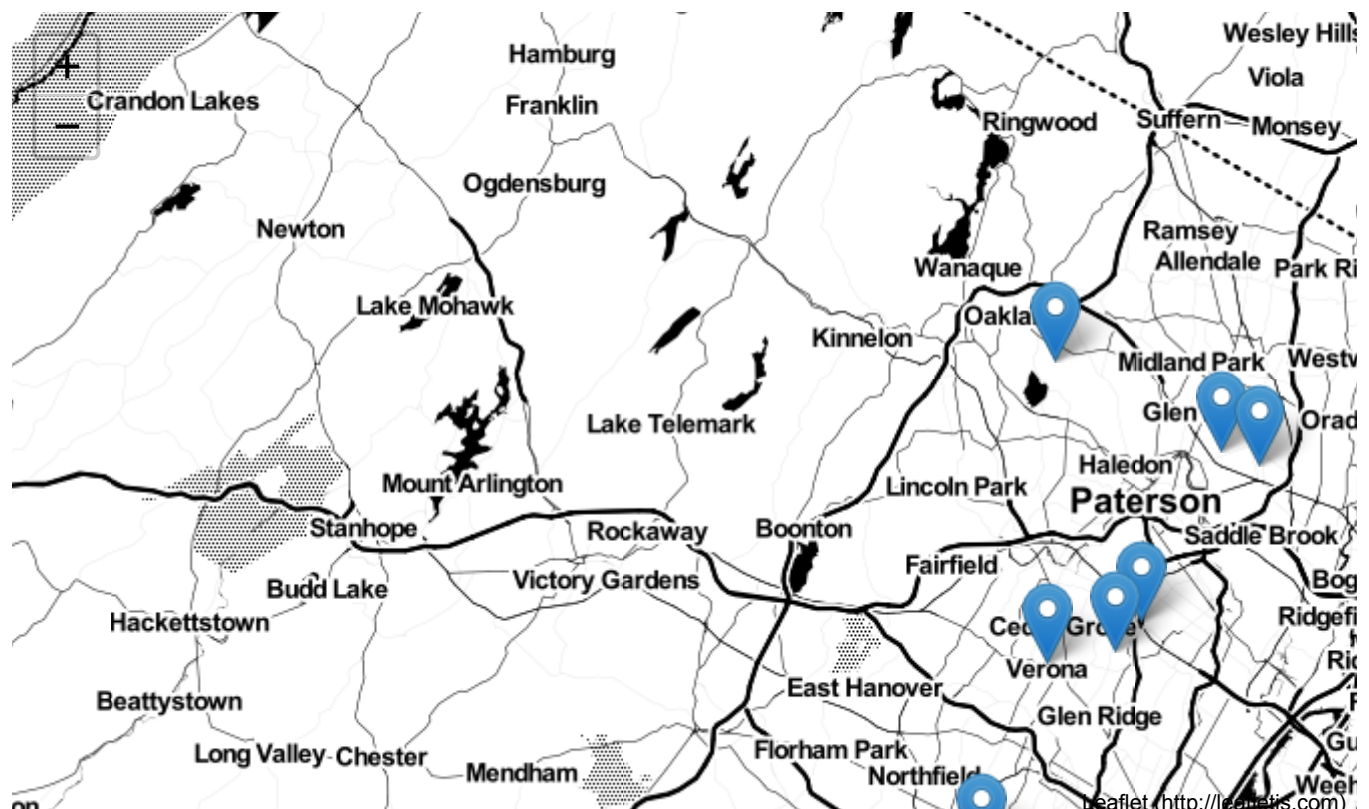
|   | VendorID | tpep_pickup_datetime | tpep_dropoff_datetime | passenger_count | trip_distanc |
|---|---|---|---|---|---|
| **0** | 2 | 2015-01-15 19:05:39 | 2015-01-15 19:23:42 | 1 | 1.5 |
| **1** | 1 | 2015-01-10 20:33:38 | 2015-01-10 20:53:28 | 1 | 3.3 |
| **2** | 1 | 2015-01-10 20:33:38 | 2015-01-10 20:43:41 | 1 | 1.8 |
| **3** | 1 | 2015-01-10 20:33:39 | 2015-01-10 20:35:31 | 1 | 0.5 |
| **4** | 1 | 2015-01-10 20:33:39 | 2015-01-10 20:52:58 | 1 | 3.0 |

## 1. Pickup Latitude and Pickup Longitude

It is inferred from the source https://www.flickr.com/places/info/2459115 that New York is bounded by the location cordinates(lat,long) - (40.5774, -74.15) & (40.9176,-73.7004) so hence any cordinates not within these cordinates are not considered by us as we are only concerned with pickups which originate within New York.

```
1  # Plotting pickup cordinates which are outside the bounding box of New-York
2  # we will collect all the points outside the bounding box of newyork city to outlier_locat
3  outlier_locations = month[((month.pickup_longitude <= -74.15) | (month.pickup_latitude <=
4                            (month.pickup_longitude >= -73.7004) | (month.pickup_latitude >= 40.917
5
6  # creating a map with the a base location
7  # read more about the folium here: http://folium.readthedocs.io/en/latest/quickstart.html
8
9  # note: you dont need to remember any of these, you dont need indeepth knowledge on these
10
11 map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')
12
13 # we will spot only first 100 outliers on the map, plotting all the outliers will take mor
14 sample_locations = outlier_locations.head(10000)
15 for i,j in sample_locations.iterrows():
16     if int(j['pickup_latitude']) != 0:
17         folium.Marker(list((j['pickup_latitude'],j['pickup_longitude']))).add_to(map_osm)
18 map_osm
```

**Observation:-** As you can see above that there are some points just outside the boundary but there are a few that are in either South america, Mexico or Canada


## ▼ 2. Dropoff Latitude & Dropoff Longitude


It is inferred from the source https://www.flickr.com/places/info/2459115 that New York is bounded by the location cordinates(lat,long) - (40.5774, -74.15) & (40.9176,-73.7004) so hence any cordinates not within these cordinates are not considered by us as we are only concerned with dropoffs which are within New York.


```
1 # Plotting dropoff cordinates which are outside the bounding box of New-York
2 # we will collect all the points outside the bounding box of newyork city to outlier_locat
3 outlier_locations = month[((month.dropoff_longitude <= -74.15) | (month.dropoff_latitude <
4                     (month.dropoff_longitude >= -73.7004) | (month.dropoff_latitude >= 40.9
5
6 # creating a map with the a base location
7 # read more about the folium here: http://folium.readthedocs.io/en/latest/quickstart.html
8
9 # note: you dont need to remember any of these, you dont need indeepth knowledge on these
10
11 map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')
12
13 # we will spot only first 100 outliers on the map, plotting all the outliers will take mor
14 sample_locations = outlier_locations.head(10000)
15 for i,j in sample_locations.iterrows():
```

**Observation:-** The observations here are similar to those obtained while analysing pickup latitude

## 3. Trip Durations:

According to NYC Taxi & Limousine Commision Regulations **the maximum allowed trip duration in a 24 hour interval is 12 hours.**

```python
1 #The timestamps are converted to unix so as to get duration(trip-time) & speed also pickup
2
3 # in out data we have time in the formate "YYYY-MM-DD HH:MM:SS" we convert thiss sting to
4 # https://stackoverflow.com/a/27914405
5 def convert_to_unix(s):
6     return time.mktime(datetime.datetime.strptime(s, "%Y-%m-%d %H:%M:%S").timetuple())
7
8
9
10 # we return a data frame which contains the columns
11 # 1.'passenger_count' : self explanatory
12 # 2.'trip_distance' : self explanatory
13 # 3.'pickup_longitude' : self explanatory
14 # 4.'pickup_latitude' : self explanatory
15 # 5.'dropoff_longitude' : self explanatory
16 # 6.'dropoff_latitude' : self explanatory
17 # 7.'total_amount' : total fair that was paid
18 # 8.'trip_times' : duration of each trip
19 # 9.'pickup_times : pickup time converted into unix time
20 # 10.'Speed' : velocity of each trip
21 def return_with_trip_times(month):
22     duration = month[['tpep_pickup_datetime','tpep_dropoff_datetime']].compute()
23     #pickups and dropoffs to unix time
24     duration_pickup = [convert_to_unix(x) for x in duration['tpep_pickup_datetime'].values
25     duration_drop = [convert_to_unix(x) for x in duration['tpep_dropoff_datetime'].values]
26     #calculate duration of trips
27     durations = (np.array(duration_drop) - np.array(duration_pickup))/float(60)
28
29     #append durations of trips and speed in miles/hr to a new dataframe
30     new_frame = month[['passenger_count','trip_distance','pickup_longitude','pickup_latitu
31
32     new_frame['trip_times'] = durations
33     new_frame['pickup_times'] = duration_pickup
34     new_frame['Speed'] = 60*(new_frame['trip_distance']/new_frame['trip_times'])
35
36     return new_frame
37
38 # print(frame_with_durations.head())
39 #  passenger_count  trip_distance pickup_longitude  pickup_latitude dropoff_longitude drop
40 #  1                 1.59          -73.993896        40.750111       -73.974785         40.7
41 #  1                 3.30          -74.001648        40.724243       -73.994415         40.75910
42 #  1                 1.80          -73.963341        40.802788       -73.951820         40.824
```
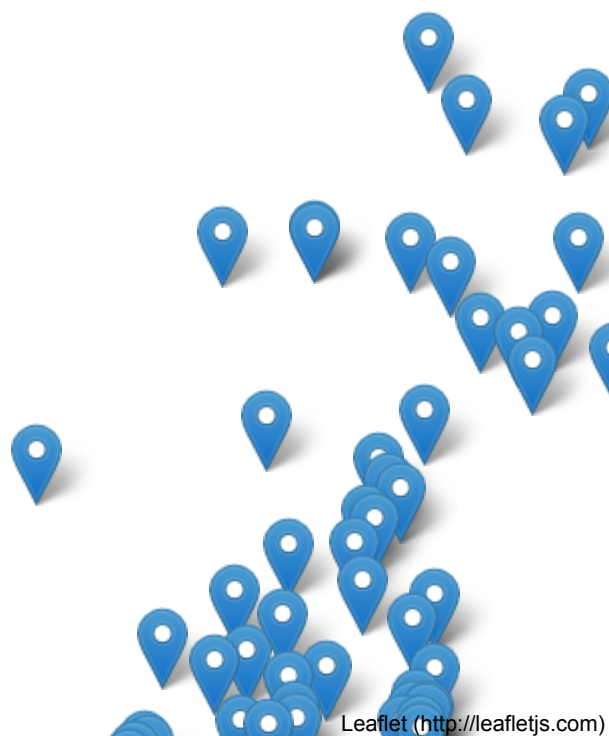
```
16     if int(j['pickup_latitude']) != 0:
17         folium.Marker(list((j['dropoff_latitude'],j['dropoff_longitude']))).add_to(map_osm
18 map_osm
```

Make this Notebook Trusted to load map: File -> Trust Notebook



Leaflet (http://leafletjs.com)

```
43 #    1                      0.50      -74.009087        40.713818     -74.004326        40.71998
44 #    1                      3.00      -73.971176        40.762428     -74.004181        40.74265
45 frame_with_durations = return_with_trip_times(month)
```
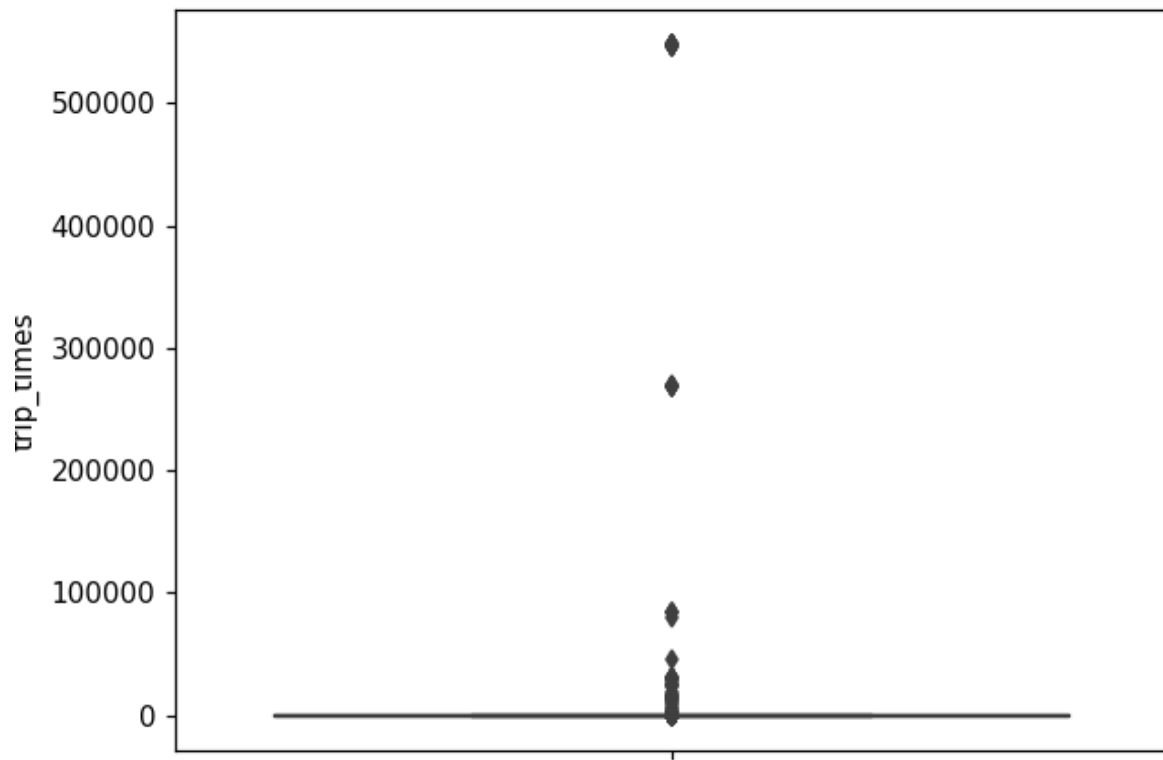
```
1 # the skewed box plot shows us the presence of outliers
2 sns.boxplot(y="trip_times", data =frame_with_durations)
3 plt.show()
```



```
1 #calculating 0-100th percentile to find a the correct percentile value for removal of outl
2 for i in range(0,100,10):
3     var =frame_with_durations["trip_times"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print ("100 percentile value is ",var[-1])
```

```
    0 percentile value is -1211.0166666666667
    10 percentile value is 3.8333333333333335
    20 percentile value is 5.38333333333334
```

```
1 #looking further from the 99th percecntile
2 for i in range(90,100):
3     var =frame_with_durations["trip_times"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print ("100 percentile value is ",var[-1])
```

```
    90 percentile value is 23.45
    91 percentile value is 24.35
    92 percentile value is 25.383333333333333
    93 percentile value is 26.55
    94 percentile value is 27.933333333333334
    95 percentile value is 29.583333333333332
    96 percentile value is 31.683333333333334
    97 percentile value is 34.46666666666667
    98 percentile value is 38.71666666666667
    99 percentile value is 46.75
    100 percentile value is  548555.633333
```

```
1 #removing data based on our analysis and TLC regulations
2 frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_times>1) & (
```

```
1 #box-plot after removal of outliers
2 sns.boxplot(y="trip_times", data =frame_with_durations_modified)
3 plt.show()
```
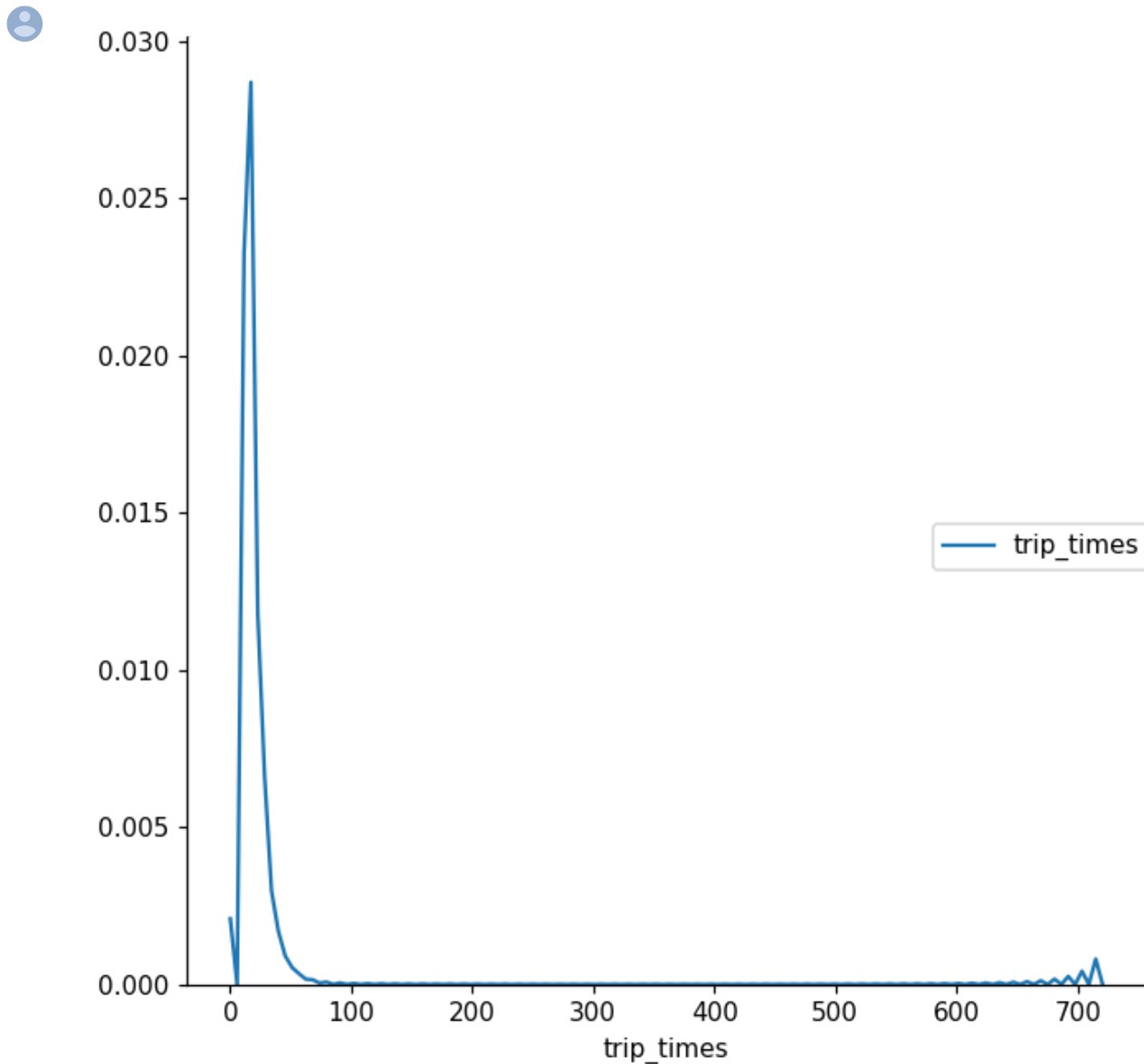
700

```
1 #pdf of trip-times after removing the outliers
2 sns.FacetGrid(frame_with_durations_modified,size=6) \
3       .map(sns.kdeplot,"trip_times") \
4       .add_legend();
5 plt.show();
```



```
1 #converting the values to log-values to chec for log-normal
2 import math
3 frame_with_durations_modified['log_times']=[math.log(i) for i in frame_with_durations_modi
```

```
1 #pdf of log-values
2 sns.FacetGrid(frame_with_durations_modified,size=6) \
```

```
3        .map(sns.kdeplot,"log_times") \
4        .add_legend();
5 plt.show();
```
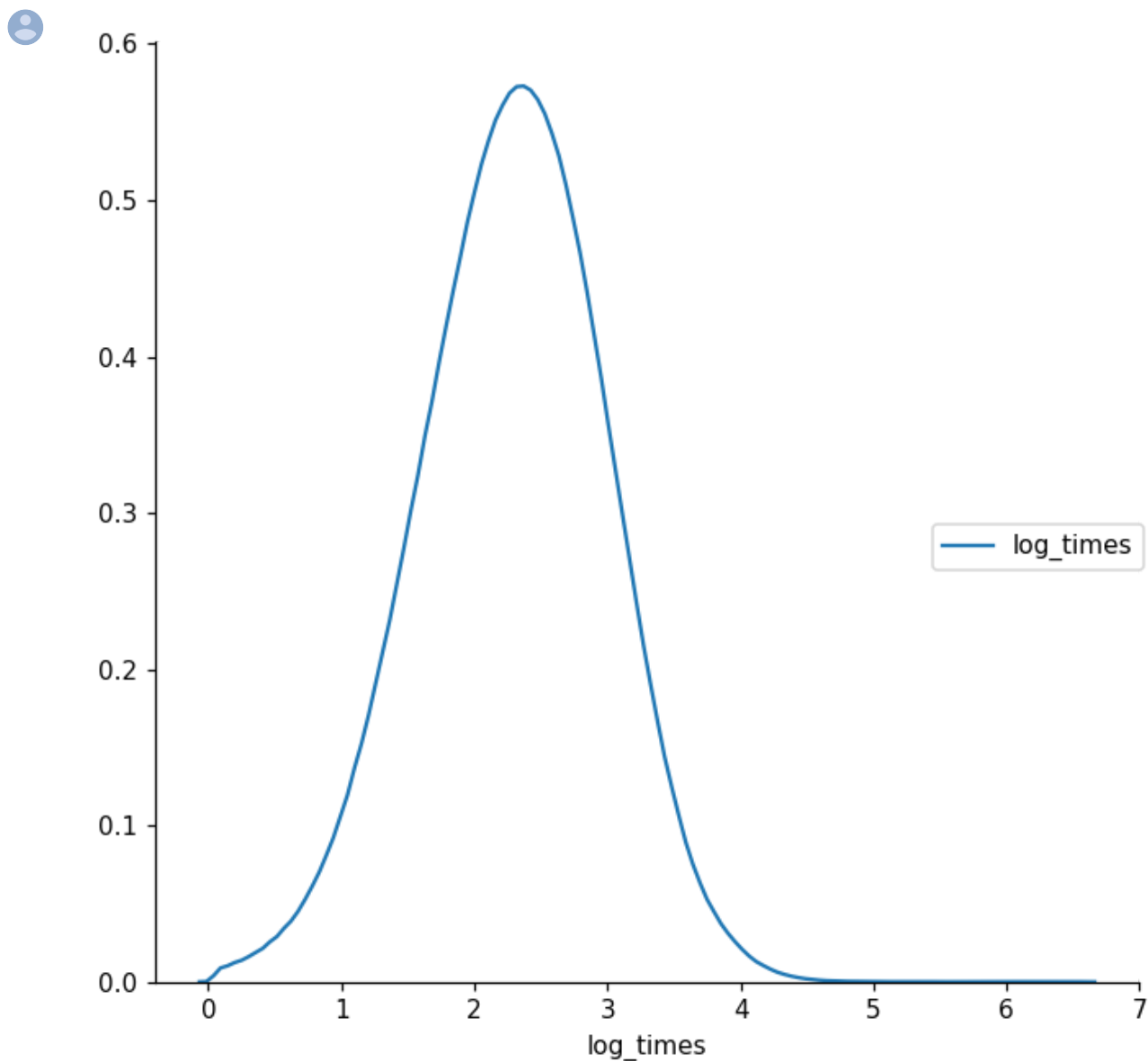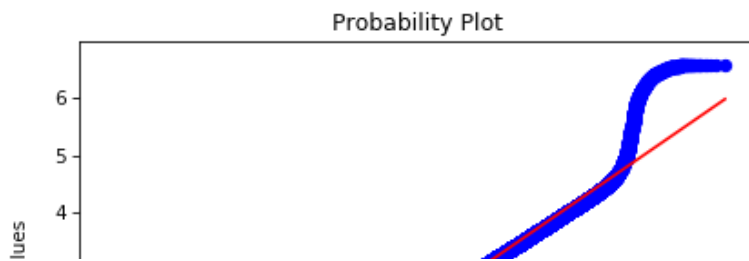


```
1 #Q-Q plot for checking if trip-times is log-normal
2 scipy.stats.probplot(frame_with_durations_modified['log_times'].values, plot=plt)
3 plt.show()
```

**Probability Plot**

## 4. Speed

```
1 # check for any outliers in the data after trip duration outliers removed
2 # box-plot for speeds with outliers
3 frame_with_durations_modified['Speed'] = 60*(frame_with_durations_modified['trip_distance'
4 sns.boxplot(y="Speed", data =frame_with_durations_modified)
5 plt.show()
```

```
1 #calculating speed values at each percntile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var =frame_with_durations_modified["Speed"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
   0 percentile value is 0.0
   10 percentile value is 6.409495548961425
   20 percentile value is 7.80952380952381
   30 percentile value is 8.929133858267717
```

```
1 #calculating speed values at each percntile 90,91,92,93,94,95,96,97,98,99,100
2 for i in range(90,100):
3     var =frame_with_durations_modified["Speed"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
   90 percentile value is 20.186915887850468
   91 percentile value is 20.91645569620253
   92 percentile value is 21.752988047808763
   93 percentile value is 22.721893491124263
   94 percentile value is 23.844155844155843
   95 percentile value is 25.182552504038775
   96 percentile value is 26.80851063829787
   97 percentile value is 28.84304932735426
   98 percentile value is 31.591128254580514
   99 percentile value is 35.7513566847558
   100 percentile value is  192857142.857
```

```
1 #calculating speed values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,9
2 for i in np.arange(0.0, 1.0, 0.1):
3     var =frame_with_durations_modified["Speed"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
   99.0 percentile value is 35.7513566847558
   99.1 percentile value is 36.31084727468969
   99.2 percentile value is 36.91470054446461
   99.3 percentile value is 37.588235294117645
   99.4 percentile value is 38.33035714285714
   99.5 percentile value is 39.17580340264651
   99.6 percentile value is 40.15384615384615
   99.7 percentile value is 41.338301043219076
   99.8 percentile value is 42.86631016042781
   99.9 percentile value is 45.3107822410148
   100 percentile value is  192857142.857
```

```
1 #removing further outliers based on the 99.9th percentile value
2 frame_with_durations_modified=frame_with_durations[(frame_with_durations.Speed>0) & (frame_
```

```
1 #avg.speed of cabs in New-York
2 sum(frame_with_durations_modified["Speed"]) / float(len(frame_with_durations_modified["Spe
```

```
   12.450173996027528
```

**The avg speed in Newyork speed is 12.45miles/hr, so a cab driver can travel 2 miles per 10min on avg.**

## 4. Trip Distance

```
1 # up to now we have removed the outliers based on trip durations and cab speeds
2 # lets try if there are any outliers in trip distances
3 # box-plot showing outliers in trip-distance values
4 sns.boxplot(y="trip_distance", data =frame_with_durations_modified)
5 plt.show()
```



```
1 #calculating trip distance values at each percntile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var =frame_with_durations_modified["trip_distance"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
     0 percentile value is 0.01
```

```
1 #calculating trip distance values at each percntile 90,91,92,93,94,95,96,97,98,99,100
2 for i in range(90,100):
3     var =frame_with_durations_modified["trip_distance"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```
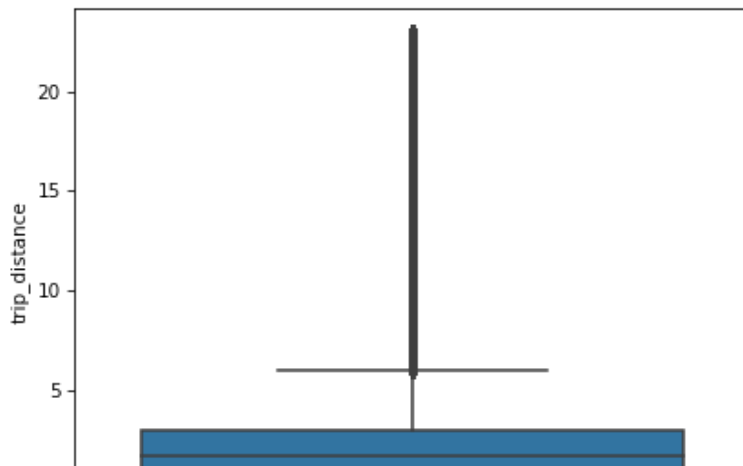
```
    90 percentile value is 5.97
    91 percentile value is 6.45
    92 percentile value is 7.07
    93 percentile value is 7.85
    94 percentile value is 8.72
    95 percentile value is 9.6
    96 percentile value is 10.6
    97 percentile value is 12.1
    98 percentile value is 16.03
    99 percentile value is 18.17
    100 percentile value is  258.9
```

```
1 #calculating trip distance values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.
2 for i in np.arange(0.0, 1.0, 0.1):
3     var =frame_with_durations_modified["trip_distance"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
    99.0 percentile value is 18.17
    99.1 percentile value is 18.37
    99.2 percentile value is 18.6
    99.3 percentile value is 18.83
    99.4 percentile value is 19.13
    99.5 percentile value is 19.5
    99.6 percentile value is 19.96
    99.7 percentile value is 20.5
    99.8 percentile value is 21.22
    99.9 percentile value is 22.57
    100 percentile value is  258.9
```

```
1 #removing further outliers based on the 99.9th percentile value
2 frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_distance>0)
```

```
1 #box-plot after removal of outliers
2 sns.boxplot(y="trip_distance", data = frame_with_durations_modified)
3 plt.show()
```

## 5. Total Fare

```
1 # up to now we have removed the outliers based on trip durations, cab speeds, and trip dis
2 # lets try if there are any outliers in based on the total_amount
3 # box-plot showing outliers in fare
4 sns.boxplot(y="total_amount", data =frame_with_durations_modified)
5 plt.show()
```



```
1 #calculating total fare amount values at each percntile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var = frame_with_durations_modified["total_amount"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
0 percentile value is -242.55
10 percentile value is 6.3
20 percentile value is 7.8
30 percentile value is 8.8
40 percentile value is 9.8
50 percentile value is 11.16
60 percentile value is 12.8
70 percentile value is 14.8
80 percentile value is 18.3
90 percentile value is 25.8
100 percentile value is  3950611.6
```
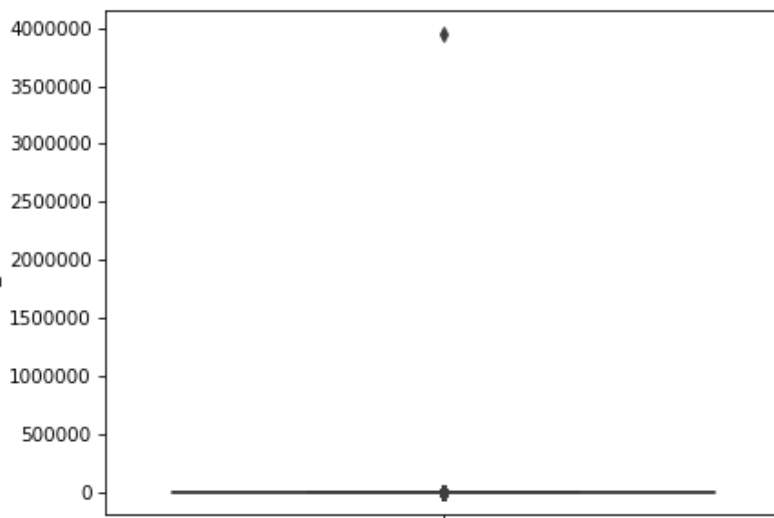
```
1 #calculating total fare amount values at each percntile 90,91,92,93,94,95,96,97,98,99,100
2 for i in range(90,100):
3     var = frame_with_durations_modified["total_amount"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
90 percentile value is 25.8
91 percentile value is 27.3
92 percentile value is 29.3
93 percentile value is 31.8
94 percentile value is 34.8
95 percentile value is 38.53
96 percentile value is 42.6
97 percentile value is 48.13
98 percentile value is 58.13
99 percentile value is 66.13
100 percentile value is  3950611.6
```

```
1 #calculating total fare amount values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.6
2 for i in np.arange(0.0, 1.0, 0.1):
3     var = frame_with_durations_modified["total_amount"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
99.0 percentile value is 68.13
99.1 percentile value is 69.13
99.2 percentile value is 69.6
99.3 percentile value is 69.73
99.4 percentile value is 69.73
99.5 percentile value is 69.76
99.6 percentile value is 72.46
99.7 percentile value is 72.73
99.8 percentile value is 80.05
99.9 percentile value is 95.55
100 percentile value is  3950611.6
```

**Observation:-** As even the 99.9th percentile value doesnt look like an outlier,as there is not much difference between the 99.8th percentile and 99.9th percentile, we move on to do graphical analyis

```
1 #below plot shows us the fare values(sorted) to find a sharp increase to remove those valu
2 # plot the fare amount excluding last two values in sorted data
3 plt.plot(var[:-2])
4 plt.show()
```



```
1 # a very sharp increase in fare values can be seen
2 # plotting last three total fare values, and we can observe there is share increase in the
3 plt.plot(var[-3:])
4 plt.show()
```



```
1 #now looking at values not including the last two points we again find a drastic increase
2 # we plot last 50 values excluding last two values
```

```
3 plt.plot(var[-50:-2])
4 plt.show()
```



## Remove all outliers/erronous points.

```
 1 #removing all outliers based on our univariate analysis above
 2 def remove_outliers(new_frame):
 3
 4
 5    a = new_frame.shape[0]
 6    print ("Number of pickup records = ",a)
 7    temp_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_frame.dropoff_l
 8                    (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitu
 9                    ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_latitud
10                    (new_frame.pickup_longitude <= -73.7004) & (new_frame.pickup_latitu
11    b = temp_frame.shape[0]
12    print ("Number of outlier coordinates lying outside NY boundaries:",(a-b))
13
14
15    temp_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times < 720)]
16    c = temp_frame.shape[0]
17    print ("Number of outliers from trip times analysis:",(a-c))
18
19
20    temp_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_distance < 23)]
21    d = temp_frame.shape[0]
22    print ("Number of outliers from trip distance analysis:",(a-d))
23
24    temp_frame = new_frame[(new_frame.Speed <= 65) & (new_frame.Speed >= 0)]
25    e = temp_frame.shape[0]
26    print ("Number of outliers from speed analysis:",(a-e))
```
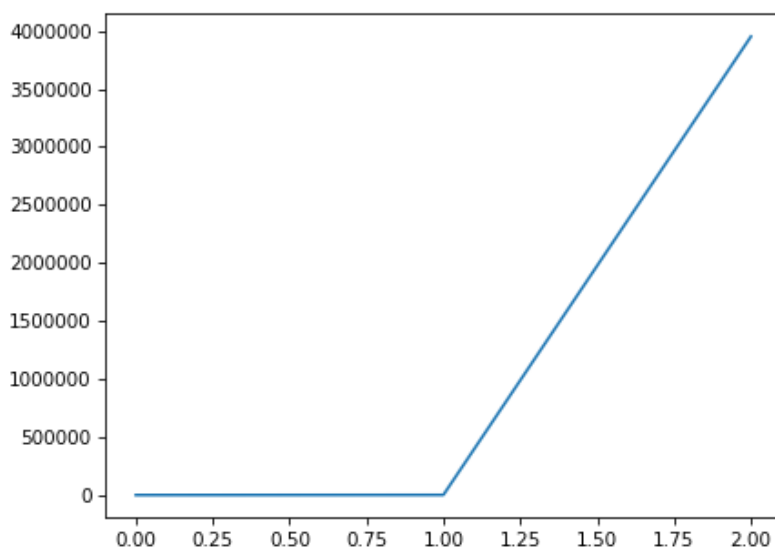
```
26     print ( Number of outliers from speed analysis: ,(a-c))
27
28     temp_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_amount >0)]
29     f = temp_frame.shape[0]
30     print ("Number of outliers from fare analysis:",(a-f))
31
32
33     new_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_frame.dropoff_lo
34                     (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitu
35                     ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_latitud
36                     (new_frame.pickup_longitude <= -73.7004) & (new_frame.pickup_latitu
37
38     new_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times < 720)]
39     new_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_distance < 23)]
40     new_frame = new_frame[(new_frame.Speed < 45.31) & (new_frame.Speed > 0)]
41     new_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_amount >0)]
42
43     print ("Total outliers removed",a - new_frame.shape[0])
44     print ("---")
45     return new_frame
```

```
1 print ("Removing outliers in the month of Jan-2015")
2 print ("----")
3 frame_with_durations_outliers_removed = remove_outliers(frame_with_durations)
4 print("fraction of data points that remain after removing outliers", float(len(frame_with_
```

```
Removing outliers in the month of Jan-2015
----
Number of pickup records =  12748986
Number of outlier coordinates lying outside NY boundaries: 293919
Number of outliers from trip times analysis: 23889
Number of outliers from trip distance analysis: 92597
Number of outliers from speed analysis: 24473
Number of outliers from fare analysis: 5275
Total outliers removed 377910
---
fraction of data points that remain after removing outliers 0.9703576425607495
```

# Data-preperation

## Clustering/Segmentation

```
1 #trying different cluster sizes to choose the right K in K-means
2 coords = frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']].va
3 neighbours=[]
4
5 def find_min_distance(cluster_centers, cluster_len):
6     nice_points = 0
7     wrong_points = 0
```

```python
 8      less2 = []
 9      more2 = []
10      min_dist=1000
11      for i in range(0, cluster_len):
12          nice_points = 0
13          wrong_points = 0
14          for j in range(0, cluster_len):
15              if j!=i:
16                  distance = gpxpy.geo.haversine_distance(cluster_centers[i][0], cluster_cen
17                  min_dist = min(min_dist,distance/(1.60934*1000))
18                  if (distance/(1.60934*1000)) <= 2:
19                      nice_points +=1
20                  else:
21                      wrong_points += 1
22          less2.append(nice_points)
23          more2.append(wrong_points)
24      neighbours.append(less2)
25      print ("On choosing a cluster size of ",cluster_len,"\nAvg. Number of Clusters within
26
27 def find_clusters(increment):
28      kmeans = MiniBatchKMeans(n_clusters=increment, batch_size=10000,random_state=42).fit(c
29      frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_du
30      cluster_centers = kmeans.cluster_centers_
31      cluster_len = len(cluster_centers)
32      return cluster_centers, cluster_len
33
34 # we need to choose number of clusters so that, there are more number of cluster regions
35 #that are close to any cluster center
36 # and make sure that the minimum inter cluster should not be very less
37 for increment in range(10, 100, 10):
38      cluster_centers, cluster_len = find_clusters(increment)
39      find_min_distance(cluster_centers, cluster_len)
```

```
On choosing a cluster size of  10
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 2.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 8.0
Min inter-cluster distance =  1.0945442325142662
---
On choosing a cluster size of  20
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 4.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 16.0
Min inter-cluster distance =  0.7131298007388065
---
On choosing a cluster size of  30
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 22.0
Min inter-cluster distance =  0.5185088176172186
---
On choosing a cluster size of  40
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 32.0
Min inter-cluster distance =  0.5069768450365043
---
On choosing a cluster size of  50
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 12.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 38.0
Min inter-cluster distance =  0.36536302598358383
---
On choosing a cluster size of  60
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 14.0
```

## Inference:

- The main objective was to find a optimal min. distance(Which roughly estimates to the radius of a cluster) between the clusters which we got was 40

```
---
1 # if check for the 50 clusters you can observe that there are two clusters with only 0.3 m
2 # so we choose 40 clusters for solve the further problem
3
4 # Getting 40 clusters using the kmeans
5
6 coords = frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']].va
7 kmeans = MiniBatchKMeans(n_clusters=40, batch_size=10000,random_state=0).fit(coords)
8 frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_durati
```

```
1 len(kmeans.cluster_centers_[0])
```

  2

```
1 len(set(kmeans.labels_))
```

  40

## Plotting the cluster centers:

```
1 # Plotting the cluster centers on OSM
2 cluster_centers = kmeans.cluster_centers_
3 cluster_len = len(cluster_centers)
4 map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')
5 for i in range(cluster_len):
6     folium.Marker(list((cluster_centers[i][0],cluster_centers[i][1])), popup=(str(cluster_
7 map_osm
```

### Plotting the clusters:



```
1  #Visualising the clusters on a map
2  def plot_clusters(frame):
3      city_long_border = (-74.03, -73.75)
4      city_lat_border = (40.63, 40.85)
5      fig, ax = plt.subplots(ncols=1, nrows=1)
6      ax.scatter(frame.pickup_longitude.values[:100000], frame.pickup_latitude.values[:10000
7                  c=frame.pickup_cluster.values[:100000], cmap='tab20', alpha=0.2)
8      ax.set_xlim(city_long_border)
9      ax.set_ylim(city_lat_border)
10     ax.set_xlabel('Longitude')
11     ax.set_ylabel('Latitude')
12     plt.show()
13
14 plot_clusters(frame_with_durations_outliers_removed)
```

# ▾ Time-binning

```
1  #Refer:https://www.unixtimestamp.com/
2  # 1420070400 : 2015-01-01 00:00:00
3  # 1422748800 : 2015-02-01 00:00:00
4  # 1425168000 : 2015-03-01 00:00:00
5  # 1427846400 : 2015-04-01 00:00:00
6  # 1430438400 : 2015-05-01 00:00:00
7  # 1433116800 : 2015-06-01 00:00:00
8
9  # 1451606400 : 2016-01-01 00:00:00
10 # 1454284800 : 2016-02-01 00:00:00
11 # 1456790400 : 2016-03-01 00:00:00
12 # 1459468800 : 2016-04-01 00:00:00
13 # 1462060800 : 2016-05-01 00:00:00
14 # 1464739200 : 2016-06-01 00:00:00
15
16 def add_pickup_bins(frame,month,year):
17     unix_pickup_times=[i for i in frame['pickup_times'].values]
18     unix_times = [[1420070400,1422748800,1425168000,1427846400,1430438400,1433116800],\
19                   [1451606400,1454284800,1456790400,1459468800,1462060800,1464739200]]
20
21     start_pickup_unix=unix_times[year-2015][month-1]
22     # https://www.timeanddate.com/time/zones/est
23     # (int((i-start_pickup_unix)/600)+33) : our unix time is in gmt to we are converting i
24     tenminutewise_binned_unix_pickup_times=[(int((i-start_pickup_unix)/600)+33) for i in u
25     frame['pickup_bins'] = np.array(tenminutewise_binned_unix_pickup_times)
26     return frame
```

```
1 # clustering, making pickup bins and grouping by pickup cluster and pickup bins
2 frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_durati
3 jan_2015_frame = add_pickup_bins(frame_with_durations_outliers_removed,1,2015)
4 jan_2015_groupby = jan_2015_frame[['pickup_cluster','pickup_bins','trip_distance']].groupb
```

```
1 # we add two more columns 'pickup_cluster'(to which cluster it belogns to)
2 # and 'pickup_bins' (to which 10min intravel the trip belongs to)
3 jan_2015_frame.head()
```

| | passenger_count | trip_distance | pickup_longitude | pickup_latitude | dropoff_longitude |
|---|---|---|---|---|---|
| 0 | 1 | 1.59 | -73.993896 | 40.750111 | -73.974785 |
| 1 | 1 | 3.30 | -74.001648 | 40.724243 | -73.994415 |
| 2 | 1 | 1.80 | -73.963341 | 40.802788 | -73.951820 |
| 3 | 1 | 0.50 | -74.009087 | 40.713818 | -74.004326 |
| 4 | 1 | 3.00 | -73.971176 | 40.762428 | -74.004181 |

```
1 # hear the trip_distance represents the number of pickups that are happend in that particu
2 # this data frame has two indices
3 # primary index: pickup_cluster (cluster number)
4 # secondary index : pickup_bins (we devid whole months time into 10min intravels 24*31*60/
5 jan_2015_groupby.head()
```

|                | trip_distance |     |
|----------------|---------------|-----|
| pickup_cluster | pickup_bins   |     |
| 0              | 33            | 104 |
|                | 34            | 200 |
|                | 35            | 208 |
|                | 36            | 141 |
|                | 37            | 155 |

```
 1 # upto now we cleaned data and prepared data for the month 2015,
 2
 3 # now do the same operations for months Jan, Feb, March of 2016
 4 # 1. get the dataframe which inlcudes only required colums
 5 # 2. adding trip times, speed, unix time stamp of pickup_time
 6 # 4. remove the outliers based on trip_times, speed, trip_duration, total_amount
 7 # 5. add pickup_cluster to each data point
 8 # 6. add pickup_bin (index of 10min intravel to which that trip belongs to)
 9 # 7. group by data, based on 'pickup_cluster' and 'pickuo_bin'
10
11 # Data Preparation for the months of Jan,Feb and March 2016
12 def datapreparation(month,kmeans,month_no,year_no):
13
14     print ("Return with trip times..")
15
16     frame_with_durations = return_with_trip_times(month)
17
18     print ("Remove outliers..")
19     frame_with_durations_outliers_removed = remove_outliers(frame_with_durations)
20
21     print ("Estimating clusters..")
22     frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_du
23     #frame_with_durations_outliers_removed_2016['pickup_cluster'] = kmeans.predict(frame_w
24
25     print ("Final groupbying..")
26     final_updated_frame = add_pickup_bins(frame_with_durations_outliers_removed,month_no,y
27     final_groupby_frame = final_updated_frame[['pickup_cluster','pickup_bins','trip_distan
28
29     return final_updated_frame,final_groupby_frame
30
31 month_jan_2016 = dd.read_csv('yellow_tripdata_2016-01.csv')
32 month_feb_2016 = dd.read_csv('yellow_tripdata_2016-02.csv')
```

```
33 month_mar_2016 = dd.read_csv('yellow_tripdata_2016-03.csv')
34
35 jan_2016_frame,jan_2016_groupby = datapreparation(month_jan_2016,kmeans,1,2016)
36 feb_2016_frame,feb_2016_groupby = datapreparation(month_feb_2016,kmeans,2,2016)
37 mar_2016_frame,mar_2016_groupby = datapreparation(month_mar_2016,kmeans,3,2016)
```

Return with trip times..
Remove outliers..
Number of pickup records =  10906858
Number of outlier coordinates lying outside NY boundaries: 214677
Number of outliers from trip times analysis: 27190
Number of outliers from trip distance analysis: 79742
Number of outliers from speed analysis: 21047
Number of outliers from fare analysis: 4991
Total outliers removed 297784
---
Estimating clusters..
Final groupbying..
Return with trip times..
Remove outliers..
Number of pickup records =  11382049
Number of outlier coordinates lying outside NY boundaries: 223161
Number of outliers from trip times analysis: 27670
Number of outliers from trip distance analysis: 81902
Number of outliers from speed analysis: 22437
Number of outliers from fare analysis: 5476
Total outliers removed 308177
---
Estimating clusters..
Final groupbying..
Return with trip times..
Remove outliers..
Number of pickup records =  12210952
Number of outlier coordinates lying outside NY boundaries: 232444
Number of outliers from trip times analysis: 30868
Number of outliers from trip distance analysis: 87318
Number of outliers from speed analysis: 23889
Number of outliers from fare analysis: 5859
Total outliers removed 324635
---
Estimating clusters..
Final groupbying..

## Smoothing

```
1 # Gets the unique bins where pickup values are present for each each region
2
3 # for each cluster region we will collect all the indices of 10min intravels in which the
4 # we got an observation that there are some pickpbins that doesnt have any pickups
5 def return_unq_pickup_bins(frame):
6     values = []
7     for i in range(0 40):
```

```
 7      for i in range(0,40):
 8          new = frame[frame['pickup_cluster'] == i]
 9          list_unq = list(set(new['pickup_bins']))
10          list_unq.sort()
11          values.append(list_unq)
12      return values
```

```
 1 # for every month we get all indices of 10min intravels in which atleast one pickup got ha
 2
 3 #jan
 4 jan_2015_unique = return_unq_pickup_bins(jan_2015_frame)
 5 jan_2016_unique = return_unq_pickup_bins(jan_2016_frame)
 6
 7 #feb
 8 feb_2016_unique = return_unq_pickup_bins(feb_2016_frame)
 9
10 #march
11 mar_2016_unique = return_unq_pickup_bins(mar_2016_frame)
```

```
 1 # for each cluster number of 10min intravels with 0 pickups
 2 for i in range(40):
 3      print("for the ",i,"th cluster number of 10min intavels with zero pickups: ",4464 - le
 4      print('-'*60)
```

```
for the  0 th cluster number of 10min intavels with zero pickups:  40
-------------------------------------------------------------
for the  1 th cluster number of 10min intavels with zero pickups:  1985
-------------------------------------------------------------
for the  2 th cluster number of 10min intavels with zero pickups:  29
-------------------------------------------------------------
for the  3 th cluster number of 10min intavels with zero pickups:  354
-------------------------------------------------------------
for the  4 th cluster number of 10min intavels with zero pickups:  37
-------------------------------------------------------------
for the  5 th cluster number of 10min intavels with zero pickups:  153
-------------------------------------------------------------
for the  6 th cluster number of 10min intavels with zero pickups:  34
-------------------------------------------------------------
for the  7 th cluster number of 10min intavels with zero pickups:  34
-------------------------------------------------------------
for the  8 th cluster number of 10min intavels with zero pickups:  117
-------------------------------------------------------------
for the  9 th cluster number of 10min intavels with zero pickups:  40
-------------------------------------------------------------
for the  10 th cluster number of 10min intavels with zero pickups:  25
-------------------------------------------------------------
for the  11 th cluster number of 10min intavels with zero pickups:  44
-------------------------------------------------------------
for the  12 th cluster number of 10min intavels with zero pickups:  42
-------------------------------------------------------------
for the  13 th cluster number of 10min intavels with zero pickups:  28
-------------------------------------------------------------
for the  14 th cluster number of 10min intavels with zero pickups:  26
-------------------------------------------------------------
for the  15 th cluster number of 10min intavels with zero pickups:  31
-------------------------------------------------------------
for the  16 th cluster number of 10min intavels with zero pickups:  40
-------------------------------------------------------------
for the  17 th cluster number of 10min intavels with zero pickups:  58
-------------------------------------------------------------
for the  18 th cluster number of 10min intavels with zero pickups:  1190
-------------------------------------------------------------
for the  19 th cluster number of 10min intavels with zero pickups:  1357
-------------------------------------------------------------
for the  20 th cluster number of 10min intavels with zero pickups:  53
-------------------------------------------------------------
for the  21 th cluster number of 10min intavels with zero pickups:  29
-------------------------------------------------------------
for the  22 th cluster number of 10min intavels with zero pickups:  29
-------------------------------------------------------------
for the  23 th cluster number of 10min intavels with zero pickups:  163
-------------------------------------------------------------
for the  24 th cluster number of 10min intavels with zero pickups:  35
-------------------------------------------------------------
for the  25 th cluster number of 10min intavels with zero pickups:  41
-------------------------------------------------------------
for the  26 th cluster number of 10min intavels with zero pickups:  31
-------------------------------------------------------------
for the  27 th cluster number of 10min intavels with zero pickups:  214
```

there are two ways to fill up these values

- Fill the missing value with 0's
- Fill the missing values with the avg values

  - Case 1:(values missing at the start)

    Ex1: \_ \_ \_ x =>ceil(x/4), ceil(x/4), ceil(x/4), ceil(x/4)

    Ex2: \_ \_ x => ceil(x/3), ceil(x/3), ceil(x/3)

  - Case 2:(values missing in middle)

    Ex1: x \_ \_ y => ceil((x+y)/4), ceil((x+y)/4), ceil((x+y)/4), ceil((x+y)/4)

    Ex2: x \_ \_ \_ y => ceil((x+y)/5), ceil((x+y)/5), ceil((x+y)/5), ceil((x+y)/5), ceil((x+y)/5)

  - Case 3:(values missing at the end)

    Ex1: x \_ \_ \_ => ceil(x/4), ceil(x/4), ceil(x/4), ceil(x/4)

    Ex2: x \_ => ceil(x/2), ceil(x/2)

```
for the  37 th cluster number of 10min intavels with zero pickups:  321
```

```python
1 # Fills a value of zero for every bin where no pickup data is present
2 # the count_values: number pickps that are happened in each region for each 10min intravel
3 # there wont be any value if there are no picksups.
4 # values: number of unique bins
5
6 # for every 10min intravel(pickup_bin) we will check it is there in our unique bin,
7 # if it is there we will add the count_values[index] to smoothed data
8 # if not we add 0 to the smoothed data
9 # we finally return smoothed data
10 def fill_missing(count_values,values):
11     smoothed_regions=[]
12     ind=0
13     for r in range(0,40):
14         smoothed_bins=[]
15         for i in range(4464):
16             if i in values[r]:
17                 smoothed_bins.append(count_values[ind])
18                 ind+=1
19             else:
20                 smoothed_bins.append(0)
21         smoothed_regions.extend(smoothed_bins)
22     return smoothed_regions
```

```python
1 # Fills a value of zero for every bin where no pickup data is present
2 # the count_values: number pickps that are happened in each region for each 10min intravel
3 # there wont be any value if there are no picksups.
4 # values: number of unique bins
5
6 # for every 10min intravel(pickup_bin) we will check it is there in our unique bin,
7 # if it is there we will add the count_values[index] to smoothed data
8 # if not we add smoothed data (which is calculated based on the methods that are discussed
9 # we finally return smoothed data
10 def smoothing(count_values,values):
11     smoothed_regions=[] # stores list of final smoothed values of each reigion
12     ind=0
```

```
12      ind=0
13      repeat=0
14      smoothed_value=0
15      for r in range(0,40):
16          smoothed_bins=[] #stores the final smoothed values
17          repeat=0
18          for i in range(4464):
19              if repeat!=0: # prevents iteration for a value which is already visited/resolv
20                  repeat-=1
21                  continue
22              if i in values[r]: #checks if the pickup-bin exists
23                  smoothed_bins.append(count_values[ind]) # appends the value of the pickup
24              else:
25                  if i!=0:
26                      right_hand_limit=0
27                      for j in range(i,4464):
28                          if  j not in values[r]: #searches for the left-limit or the pickup
29                              continue
30                          else:
31                              right_hand_limit=j
32                              break
33                      if right_hand_limit==0:
34                      #Case 1: When we have the last/last few values are found to be missing
35                          smoothed_value=count_values[ind-1]*1.0/((4463-i)+2)*1.0
36                          for j in range(i,4464):
37                              smoothed_bins.append(math.ceil(smoothed_value))
38                          smoothed_bins[i-1] = math.ceil(smoothed_value)
39                          repeat=(4463-i)
40                          ind-=1
41                      else:
42                      #Case 2: When we have the missing values between two known values
43                          smoothed_value=(count_values[ind-1]+count_values[ind])*1.0/((right
44                          for j in range(i,right_hand_limit+1):
45                              smoothed_bins.append(math.ceil(smoothed_value))
46                          smoothed_bins[i-1] = math.ceil(smoothed_value)
47                          repeat=(right_hand_limit-i)
48                  else:
49                      #Case 3: When we have the first/first few values are found to be missi
50                      right_hand_limit=0
51                      for j in range(i,4464):
52                          if  j not in values[r]:
53                              continue
54                          else:
55                              right_hand_limit=j
56                              break
57                      smoothed_value=count_values[ind]*1.0/((right_hand_limit-i)+1)*1.0
58                      for j in range(i,right_hand_limit+1):
59                          smoothed_bins.append(math.ceil(smoothed_value))
60                      repeat=(right_hand_limit-i)
61              ind+=1
62          smoothed_regions.extend(smoothed_bins)
63      return smoothed_regions
```

64

```python
1  #Filling Missing values of Jan-2015 with 0
2  # here in jan_2015_groupby dataframe the trip_distance represents the number of pickups th
3  jan_2015_fill = fill_missing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)
4
5  #Smoothing Missing values of Jan-2015
6  jan_2015_smooth = smoothing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)
```

```python
1  # number of 10min indices for jan 2015= 24*31*60/10 = 4464
2  # number of 10min indices for jan 2016 = 24*31*60/10 = 4464
3  # number of 10min indices for feb 2016 = 24*29*60/10 = 4176
4  # number of 10min indices for march 2016 = 24*30*60/10 = 4320
5  # for each cluster we will have 4464 values, therefore 40*4464 = 178560 (length of the jan
6  print("number of 10min intravels among all the clusters ",len(jan_2015_fill))
```

number of 10min intravels among all the clusters  178560

```python
1  # Smoothing vs Filling
2  # sample plot that shows two variations of filling missing values
3  # we have taken the number of pickups for cluster region 2
4  plt.figure(figsize=(10,5))
5  plt.plot(jan_2015_fill[4464:8920], label="zero filled values")
6  plt.plot(jan_2015_smooth[4464:8920], label="filled with avg values")
7  plt.legend()
8  plt.show()
```

```python
1  # why we choose, these methods and which method is used for which data?
2
3  # Ans: consider we have data of some month in 2015 jan 1st, 10 _ _ _ 20, i.e there are 10
4  # 10st 10min intravel, 0 pickups happened in 2nd 10mins intravel, 0 pickups happened in 3r
5  # and 20 pickups happened in 4th 10min intravel.
6  # in fill_missing method we replace these values like 10, 0, 0, 20
7  # where as in smoothing method we replace these values as 6,6,6,6,6, if you can check the
8  # that are happened in the first 40min are same in both cases, but if you can observe that
9  # wheen you are using smoothing we are looking at the future number of pickups which might
10
11 # so we use smoothing for jan 2015th data since it acts as our training data
12 # and we use simple fill_misssing method for 2016th data.
```

```python
1  # Jan-2015 data is smoothed, Jan,Feb & March 2016 data missing values are filled with zero
2  jan_2015_smooth = smoothing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)
3  jan_2016_smooth = fill_missing(jan_2016_groupby['trip_distance'].values,jan_2016_unique)
4  feb_2016_smooth = fill_missing(feb_2016_groupby['trip_distance'].values,feb_2016_unique)
5  mar_2016_smooth = fill_missing(mar_2016_groupby['trip_distance'].values,mar_2016_unique)
6
7  # Making list of all the values of pickup data in every bin for a period of 3 months and s
```

```
 8 regions_cum = []
 9
10 # a =[1,2,3]
11 # b = [2,3,4]
12 # a+b = [1, 2, 3, 2, 3, 4]
13
14 # number of 10min indices for jan 2015= 24*31*60/10 = 4464
15 # number of 10min indices for jan 2016 = 24*31*60/10 = 4464
16 # number of 10min indices for feb 2016 = 24*29*60/10 = 4176
17 # number of 10min indices for march 2016 = 24*31*60/10 = 4464
18 # regions_cum: it will contain 40 lists, each list will contain 4464+4176+4464 values whic
19 # that are happened for three months in 2016 data
20
21 for i in range(0,40):
22     regions_cum.append(jan_2016_smooth[4464*i:4464*(i+1)]+feb_2016_smooth[4176*i:4176*(i+1
23
24 # print(len(regions_cum))
25 # 40
26 # print(len(regions_cum[0]))
27 # 13104
```

## Time series and Fourier Transforms

```
 1 def uniqueish_color():
 2     """There're better ways to generate unique colors, but this isn't awful."""
 3     return plt.cm.gist_ncar(np.random.random())
 4 first_x = list(range(0,4464))
 5 second_x = list(range(4464,8640))
 6 third_x = list(range(8640,13104))
 7 for i in range(40):
 8     plt.figure(figsize=(10,4))
 9     plt.plot(first_x,regions_cum[i][:4464], color=uniqueish_color(), label='2016 Jan month
10     plt.plot(second_x,regions_cum[i][4464:8640], color=uniqueish_color(), label='2016 feb
11     plt.plot(third_x,regions_cum[i][8640:], color=uniqueish_color(), label='2016 march mon
12     plt.legend()
13     plt.show()
```

```
 1 # getting peaks: https://blog.ytotech.com/2015/11/01/findpeaks-in-python/
 2 # read more about fft function : https://docs.scipy.org/doc/numpy/reference/generated/nump
 3 Y    = np.fft.fft(np.array(jan_2016_smooth)[0:4464])
 4 # read more about the fftfreq: https://docs.scipy.org/doc/numpy/reference/generated/numpy.
 5 freq = np.fft.fftfreq(4464)
 6 n = len(freq)
```

```
 1 plt.plot( freq[:int(n/2)], np.abs(Y)[:int(n/2)] )
 2 plt.xlabel("Frequency")
```

```
 2 plt.xlabel("Frequency")
 3 plt.ylabel("Amplitude")
 4 plt.show()
```

```
 1 #Preparing the Dataframe only with x(i) values as jan-2015 data and y(i) values as jan-201
 2 ratios_jan = pd.DataFrame()
 3 ratios_jan['Given']=jan_2015_smooth
 4 ratios_jan['Prediction']=jan_2016_smooth
 5 ratios_jan['Ratios']=ratios_jan['Prediction']*1.0/ratios_jan['Given']*1.0
```

## Modelling: Baseline Models

Now we get into modelling in order to forecast the pickup densities for the months of Jan, Feb and March of 2016 for which we are using multiple models with two variations

1. Using Ratios of the 2016 data to the 2015 data i.e $R_t = P_t^{2016}/P_t^{2015}$
2. Using Previous known values of the 2016 data itself to predict the future values

## Simple Moving Averages

The First Model used is the Moving Averages Model which uses the previous n values in order to predict the next value

Using Ratio Values - $R_t = (R_{t-1} + R_{t-2} + R_{t-3}\ldots R_{t-n})/n$

```
 1 def MA_R_Predictions(ratios,month):
 2     predicted_ratio=(ratios['Ratios'].values)[0]
 3     error=[]
 4     predicted_values=[]
 5     window_size=3
 6     predicted_ratio_values=[]
 7     for i in range(0,4464*40):
 8         if i%4464==0:
 9             predicted_ratio_values.append(0)
10             predicted_values.append(0)
11             error.append(0)
12             continue
13         predicted_ratio_values.append(predicted_ratio)
14         predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
15         error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(rat
16         if i+1>=window_size:
17             predicted_ratio=sum((ratios['Ratios'].values)[(i+1)-window_size:(i+1)])/window_
18         else:
19             predicted_ratio=sum((ratios['Ratios'].values)[0:(i+1)])/(i+1)
20
```

```
20
21
22      ratios['MA_R_Predicted'] = predicted_values
23      ratios['MA_R_Error'] = error
24      mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Predi
25      mse_err = sum([e**2 for e in error])/len(error)
26      return ratios,mape_err,mse_err
```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 3 is optimal for getting the best results using Moving Averages using previous Ratio values therefore we get $R_t = (R_{t-1} + R_{t-2} + R_{t-3})/3$

Next we use the Moving averages of the 2016 values itself to predict the future value using $P_t = (P_{t-1} + P_{t-2} + P_{t-3} \ldots . P_{t-n})/n$

```
 1 def MA_P_Predictions(ratios,month):
 2      predicted_value=(ratios['Prediction'].values)[0]
 3      error=[]
 4      predicted_values=[]
 5      window_size=1
 6      predicted_ratio_values=[]
 7      for i in range(0,4464*40):
 8          predicted_values.append(predicted_value)
 9          error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
10          if i+1>=window_size:
11              predicted_value=int(sum((ratios['Prediction'].values)[(i+1)-window_size:(i+1)]
12          else:
13              predicted_value=int(sum((ratios['Prediction'].values)[0:(i+1)])/(i+1))
14
15      ratios['MA_P_Predicted'] = predicted_values
16      ratios['MA_P_Error'] = error
17      mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Predi
18      mse_err = sum([e**2 for e in error])/len(error)
19      return ratios,mape_err,mse_err
```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 1 is optimal for getting the best results using Moving Averages using previous 2016 values therefore we get $P_t = P_{t-1}$

## Weighted Moving Averages

The Moving Avergaes Model used gave equal importance to all the values in the window used, but we know intuitively that the future is more likely to be similar to the latest values and less similar to the older values. Weighted Averages converts this analogy into a mathematical relationship giving

the highest weight while computing the averages to the latest previous value and decreasing weights to the subsequent older ones

Weighted Moving Averages using Ratio Values -

$$R_t = (N * R_{t-1} + (N-1) * R_{t-2} + (N-2) * R_{t-3} \ldots 1 * R_{t-n})/(N * (N+1)/2)$$

```
1 def WA_R_Predictions(ratios,month):
2     predicted_ratio=(ratios['Ratios'].values)[0]
3     alpha=0.5
4     error=[]
5     predicted_values=[]
6     window_size=5
7     predicted_ratio_values=[]
8     for i in range(0,4464*40):
9         if i%4464==0:
10            predicted_ratio_values.append(0)
11            predicted_values.append(0)
12            error.append(0)
13            continue
14        predicted_ratio_values.append(predicted_ratio)
15        predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
16        error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(rat
17        if i+1>=window_size:
18            sum_values=0
19            sum_of_coeff=0
20            for j in range(window_size,0,-1):
21                sum_values += j*(ratios['Ratios'].values)[i-window_size+j]
22                sum_of_coeff+=j
23            predicted_ratio=sum_values/sum_of_coeff
24        else:
25            sum_values=0
26            sum_of_coeff=0
27            for j in range(i+1,0,-1):
28                sum_values += j*(ratios['Ratios'].values)[j-1]
29                sum_of_coeff+=j
30            predicted_ratio=sum_values/sum_of_coeff
31
32    ratios['WA_R_Predicted'] = predicted_values
33    ratios['WA_R_Error'] = error
34    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Predi
35    mse_err = sum([e**2 for e in error])/len(error)
36    return ratios,mape_err,mse_err
```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 5 is optimal for getting the best results using Weighted Moving Averages using previous Ratio values therefore we get

$$R_t = (5 * R_{t-1} + 4 * R_{t-2} + 3 * R_{t-3} + 2 * R_{t-4} + R_{t-5})/15$$

Weighted Moving Averages using Previous 2016 Values -

$$P_t = (N * P_{t-1} + (N-1) * P_{t-2} + (N-2) * P_{t-3} \dots 1 * P_{t-n})/(N * (N+1)/2)$$

```
1  def WA_P_Predictions(ratios,month):
2      predicted_value=(ratios['Prediction'].values)[0]
3      error=[]
4      predicted_values=[]
5      window_size=2
6      for i in range(0,4464*40):
7          predicted_values.append(predicted_value)
8          error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
9          if i+1>=window_size:
10             sum_values=0
11             sum_of_coeff=0
12             for j in range(window_size,0,-1):
13                 sum_values += j*(ratios['Prediction'].values)[i-window_size+j]
14                 sum_of_coeff+=j
15             predicted_value=int(sum_values/sum_of_coeff)
16
17         else:
18             sum_values=0
19             sum_of_coeff=0
20             for j in range(i+1,0,-1):
21                 sum_values += j*(ratios['Prediction'].values)[j-1]
22                 sum_of_coeff+=j
23             predicted_value=int(sum_values/sum_of_coeff)
24
25      ratios['WA_P_Predicted'] = predicted_values
26      ratios['WA_P_Error'] = error
27      mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Predi
28      mse_err = sum([e**2 for e in error])/len(error)
29      return ratios,mape_err,mse_err
```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 2 is optimal for getting the best results using Weighted Moving Averages using previous 2016 values therefore we get $P_t = (2 * P_{t-1} + P_{t-2})/3$

## Exponential Weighted Moving Averages

https://en.wikipedia.org/wiki/Moving_average#Exponential_moving_average Through weighted averaged we have satisfied the analogy of giving higher weights to the latest value and decreasing weights to the subsequent ones but we still do not know which is the correct weighting scheme as there are infinetly many possibilities in which we can assign weights in a non-increasing order and tune the the hyperparameter window-size. To simplify this process we use Exponential Moving

Averages which is a more logical way towards assigning weights and at the same time also using an optimal window-size.

In exponential moving averages we use a single hyperparameter alpha $(\alpha)$ which is a value between 0 & 1 and based on the value of the hyperparameter alpha the weights and the window sizes are configured.

For eg. If $\alpha = 0.9$ then the number of days on which the value of the current iteration is based is~ $1/(1-\alpha) = 10$ i.e. we consider values 10 days prior before we predict the value for the current iteration. Also the weights are assigned using $2/(N+1) = 0.18$ ,where N = number of prior values being considered, hence from this it is implied that the first or latest value is assigned a

$$R'_t = \alpha * R_{t-1} + (1-\alpha) * R'_{t-1}$$

```
1 def EA_R1_Predictions(ratios,month):
2     predicted_ratio=(ratios['Ratios'].values)[0]
3     alpha=0.6
4     error=[]
5     predicted_values=[]
6     predicted_ratio_values=[]
7     for i in range(0,4464*40):
8         if i%4464==0:
9             predicted_ratio_values.append(0)
10            predicted_values.append(0)
11            error.append(0)
12            continue
13        predicted_ratio_values.append(predicted_ratio)
14        predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
15        error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(rat
16        predicted_ratio = (alpha*predicted_ratio) + (1-alpha)*((ratios['Ratios'].values)[i
17
18    ratios['EA_R1_Predicted'] = predicted_values
19    ratios['EA_R1_Error'] = error
20    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Predi
21    mse_err = sum([e**2 for e in error])/len(error)
22    return ratios,mape_err,mse_err
```

$$P'_t = \alpha * P_{t-1} + (1-\alpha) * P'_{t-1}$$

```
1 def EA_P1_Predictions(ratios,month):
2     predicted_value= (ratios['Prediction'].values)[0]
3     alpha=0.3
4     error=[]
5     predicted_values=[]
6     for i in range(0,4464*40):
7         if i%4464==0:
8             predicted_values.append(0)
9             error.append(0)
```

```
10                continue
11          predicted_values.append(predicted_value)
12          error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
13          predicted_value =int((alpha*predicted_value) + (1-alpha)*((ratios['Prediction'].va
14
15      ratios['EA_P1_Predicted'] = predicted_values
16      ratios['EA_P1_Error'] = error
17      mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Predi
18      mse_err = sum([e**2 for e in error])/len(error)
19      return ratios,mape_err,mse_err
```

```
1 mean_err=[0]*10
2 median_err=[0]*10
3 ratios_jan,mean_err[0],median_err[0]=MA_R_Predictions(ratios_jan,'jan')
4 ratios_jan,mean_err[1],median_err[1]=MA_P_Predictions(ratios_jan,'jan')
5 ratios_jan,mean_err[2],median_err[2]=WA_R_Predictions(ratios_jan,'jan')
6 ratios_jan,mean_err[3],median_err[3]=WA_P_Predictions(ratios_jan,'jan')
7 ratios_jan,mean_err[4],median_err[4]=EA_R1_Predictions(ratios_jan,'jan')
8 ratios_jan,mean_err[5],median_err[5]=EA_P1_Predictions(ratios_jan,'jan')
```

## ▾ Comparison between baseline models

We have chosen our error metric for comparison between models as **MAPE (Mean Absolute Percentage Error)** so that we can know that on an average how good is our model with predictions and **MSE (Mean Squared Error)** is also used so that we have a clearer understanding as to how well our forecasting model performs with outliers so that we make sure that there is not much of a error margin between our prediction and the actual value

```
1 print ("Error Metric Matrix (Forecasting Methods) - MAPE & MSE")
2 print ("----------------------------------------------------------------------------------
3 print ("Moving Averages (Ratios) -                        MAPE: ",mean_err[0],"
4 print ("Moving Averages (2016 Values) -                   MAPE: ",mean_err[1],"
5 print ("----------------------------------------------------------------------------------
6 print ("Weighted Moving Averages (Ratios) -               MAPE: ",mean_err[2],"
7 print ("Weighted Moving Averages (2016 Values) -          MAPE: ",mean_err[3],"
8 print ("----------------------------------------------------------------------------------
9 print ("Exponential Moving Averages (Ratios) -            MAPE: ",mean_err[4],"      MSE
10 print ("Exponential Moving Averages (2016 Values) -       MAPE: ",mean_err[5],"      MSE
```

```
Error Metric Matrix (Forecasting Methods) - MAPE & MSE
--------------------------------------------------------------------------------
Moving Averages (Ratios) -                          MAPE:   0.182115517339      MSE:
```

**Plese Note:-** The above comparisons are made using Jan 2015 and Jan 2016 only

```
weighted Moving Averages (Ratios) -                 MAPE.   0.178488923438      MSE.
```

From the above matrix it is inferred that the best forecasting model for our prediction would be:-

$P'_t = \alpha * P_{t-1} + (1 - \alpha) * P'_{t-1}$ i.e Exponential Moving Averages using 2016 Values

# Regression Models

## Train-Test Split

Before we start predictions using the tree based regression models we take 3 months of 2016 pickup data and split it such that for every region we have 70% data in train and 30% in test, ordered date-wise for every region

```
1 with open('/content/drive/My Drive/NYTP/kmeans.cluster_centers_ .pkl', 'ab') as fr:
2   for i in kmeans.cluster_centers_:
3     pickle.dump(i,fr)
```

```
1 with open('/content/drive/My Drive/NYTP/regions_cum.pkl', 'ab') as fr:
2   for i in regions_cum:
3     pickle.dump(i,fr)
```

```
1 import os
2 import pickle
```

```
1 kmeans_cluster_center_ =[]
2
3 with open('/content/drive/My Drive/NYTP/kmeans.cluster_centers_ .pkl', 'rb') as fr:
4   try:
5     while True:
6       kmeans_cluster_center_ .append(pickle.load(fr))
7   except EOFError:
8     pass
```

```
1 regions_cum =[]
2
3 with open('/content/drive/My Drive/NYTP/regions_cum.pkl', 'rb') as fr:
4   try:
5     while True:
6       regions_cum.append(pickle.load(fr))
7     t FOFF
```

```
 7    except EOFError:
 8      pass
```

```
 1 len((list(kmeans.cluster_center_)))
```

> 80

```
 1 len(regions_cum)
```

> 40

```
 1 # Preparing data to be split into train and test, The below prepares data in cumulative fo
 2 # number of 10min indices for jan 2015= 24*31*60/10 = 4464
 3 # number of 10min indices for jan 2016 = 24*31*60/10 = 4464
 4 # number of 10min indices for feb 2016 = 24*29*60/10 = 4176
 5 # number of 10min indices for march 2016 = 24*31*60/10 = 4464
 6 # regions_cum: it will contain 40 lists, each list will contain 4464+4176+4464 values whic
 7 # that are happened for three months in 2016 data
 8 # print(len(regions_cum))
 9 # 40
10 # print(len(regions_cum[0]))
11 # 12960
12
13 # we take number of pickups that are happened in last 5 10min intravels
14 number_of_time_stamps = 5
15
16 # output varaible
17 # it is list of lists
18 # it will contain number of pickups 13099 for each cluster
19 output = []
20
21
22 # tsne_lat will contain 13104-5=13099 times lattitude of cluster center for every cluster
23 # Ex: [[cent_lat 13099times],[cent_lat 13099times], [cent_lat 13099times].... 40 lists]
24 # it is list of lists
25 tsne_lat = []
26
27
28 # tsne_lon will contain 13104-5=13099 times logitude of cluster center for every cluster
29 # Ex: [[cent_long 13099times],[cent_long 13099times], [cent_long 13099times].... 40 lists]
30 # it is list of lists
31 tsne_lon = []
32
33 # we will code each day
34 # sunday = 0, monday=1, tue = 2, wed=3, thur=4, fri=5,sat=6
35 # for every cluster we will be adding 13099 values, each value represent to which day of t
36 # it is list of lists
37 tsne_weekday = []
38
39 # its an numbpy array, of shape (523960, 5)
```

```
40 # each row corresponds to an entry in out data
41 # for the first row we will have [f0,f1,f2,f3,f4] fi=number of pickups happened in i+1th 1
42 # the second row will have [f1,f2,f3,f4,f5]
43 # the third row will have [f2,f3,f4,f5,f6]
44 # and so on...
45 tsne_feature = []
46
47
48 tsne_feature = [0]*number_of_time_stamps
49 for i in range(0,40):
50     tsne_lat.append([kmeans_cluster_center_[i][0]]*13099)
51     tsne_lon.append([kmeans_cluster_center_[i][1]]*13099)
52     # jan 1st 2016 is thursday, so we start our day from 4: "(int(k/144))%7+4"
53     # our prediction start from 5th 10min intravel since we need to have number of pickups
54     tsne_weekday.append([int(((int(k/144))%7+4)%7) for k in range(5,4464+4176+4464)])
55     # regions_cum is a list of lists [[x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3..x
56     tsne_feature = np.vstack((tsne_feature, [regions_cum[i][r:r+number_of_time_stamps] for
57     output.append(regions_cum[i][5:])
58 tsne_feature = tsne_feature[1:]
```

```
1 len(tsne_lat[0])*len(tsne_lat) == tsne_feature.shape[0] == len(tsne_weekday)*len(tsne_week
```

> True

```
 1 # Getting the predictions of exponential moving averages to be used as a feature in cumula
 2
 3 # upto now we computed 8 features for every data point that starts from 50th min of the da
 4 # 1. cluster center lattitude
 5 # 2. cluster center longitude
 6 # 3. day of the week
 7 # 4. f_t_1: number of pickups that are happened previous t-1th 10min intravel
 8 # 5. f_t_2: number of pickups that are happened previous t-2th 10min intravel
 9 # 6. f_t_3: number of pickups that are happened previous t-3th 10min intravel
10 # 7. f_t_4: number of pickups that are happened previous t-4th 10min intravel
11 # 8. f_t_5: number of pickups that are happened previous t-5th 10min intravel
12
13 # from the baseline models we said the exponential weighted moving avarage gives us the be
14 # we will try to add the same exponential weighted moving avarage at t as a feature to our
15 # exponential weighted moving avarage => p'(t) = alpha*p'(t-1) + (1-alpha)*P(t-1)
16 alpha=0.3
17
18 # it is a temporary array that store exponential weighted moving avarage for each 10min in
19 # for each cluster it will get reset
20 # for every cluster it contains 13104 values
21 predicted_values=[]
22
23 # it is similar like tsne_lat
24 # it is list of lists
25 # predict_list is a list of lists [[x5,x6,x7..x13104], [x5,x6,x7..x13104], [x5,x6,x7..x131
26 predict_list = []
27 tsne_flat_exp_avg = []
```

```
27 tshc_fiat_exp_avg = []
28 for r in range(0,40):
29     for i in range(0,13104):
30         if i==0:
31             predicted_value= regions_cum[r][0]
32             predicted_values.append(0)
33             continue
34         predicted_values.append(predicted_value)
35         predicted_value =int((alpha*predicted_value) + (1-alpha)*(regions_cum[r][i]))
36     predict_list.append(predicted_values[5:])
37     predicted_values=[]
```

```
1 # train, test split : 70% 30% split
2 # Before we start predictions using the tree based regression models we take 3 months of 2
3 # and split it such that for every region we have 70% data in train and 30% in test,
4 # ordered date-wise for every region
5 print("size of train data :", int(13099*0.7))
6 print("size of test data :", int(13099*0.3))
```

```
size of train data : 9169
size of test data : 3929
```

```
1 from statsmodels.tsa.api import ExponentialSmoothing,SimpleExpSmoothing, Holt
```

```
1 def double_exponential_smoothing(series, alpha, beta):
2     result = [series[0]]
3     for n in range(1, len(series)+1):
4         if n == 1:
5             level, trend = series[0], series[1] - series[0]
6         if n >= len(series): # we are forecasting
7             value = result[-1]
8         else:
9             value = series[n]
10        last_level, level = level, alpha*value + (1-alpha)*(level+trend)
11        trend = beta*(level-last_level) + (1-beta)*trend
12        result.append(level+trend)
13    return result
```

```
1 exp_smt=[]
2 for i in range(40):
3   exp_smts=double_exponential_smoothing(regions_cum[i],0.3,0.716)
4   exp_smt.append(exp_smts[6:])
```

```
1 def initial_trend(series, slen):
2     sum = 0.0
3     for i in range(slen):
4         sum += float(series[i+slen] - series[i]) / slen
5     return sum / slen
6
```

```
 7
 8 def initial_seasonal_components(series, slen):
 9     seasonals = {}
10     season_averages = []
11     n_seasons = int(len(series)/slen)
12     # compute season averages
13     for j in range(n_seasons):
14         season_averages.append(sum(series[slen*j:slen*j+slen])/float(slen))
15     # compute initial values
16     for i in range(slen):
17         sum_of_vals_over_avg = 0.0
18         for j in range(n_seasons):
19             sum_of_vals_over_avg += series[slen*j+i]-season_averages[j]
20         seasonals[i] = sum_of_vals_over_avg/n_seasons
21     return seasonals
22
23 def triple_exponential_smoothing(series, slen, n_preds,gamma=0.993, alpha=0.716, beta=0.3)
24     result = []
25     seasonals = initial_seasonal_components(series, slen)
26     for i in range(len(series)+n_preds):
27         if i == 0: # initial values
28             smooth = series[0]
29             trend = initial_trend(series, slen)
30             result.append(series[0])
31             continue
32         if i >= len(series): # we are forecasting
33             m = i - len(series) + 1
34             result.append((smooth + m*trend) + seasonals[i%slen])
35         else:
36             val = series[i]
37             last_smooth, smooth = smooth, alpha*(val-seasonals[i%slen]) + (1-alpha)*(smoot
38             trend = beta * (smooth-last_smooth) + (1-beta)*trend
39             seasonals[i%slen] = gamma*(val-smooth) + (1-gamma)*seasonals[i%slen]
40             result.append(smooth+trend+seasonals[i%slen])
41     return result
```

```
1 tri_exp=[]
2 for i in range(40):
3   tri_exps=triple_exponential_smoothing(regions_cum[i],144,5)
4   tri_exp.append(tri_exps[10:])
```

```
1 amplitude_lists = []
2 frequency_lists = []
3 for i in range(40):
4     ampli  = np.abs(np.fft.fft(regions_cum[i][:13099]))
5     freq = np.abs(np.fft.fftfreq(13099, 1))
6     ampli_indices = np.argsort(-ampli)[1:]        #it will return an array of indices for
7     amplitude_values = []
8     frequency_values = []
9     for j in range(0, 9, 2):   #taking top five amplitudes and frequencies
```

```
10          amplitude_values.append(ampli[ampli_indices[j]])
11          frequency_values.append(freq[ampli_indices[j]])
12     for k in range(13104):    #those top 5 frequencies and amplitudes are same for all the
13          amplitude_lists.append(amplitude_values)
14          frequency_lists.append(frequency_values)
```

```
1 train_fft_f = [frequency_lists[i*13099:(13099*i+9169)] for i in range(40)]
2 test_fft_f = [frequency_lists[(i*13099)+9169:(13099*(i+1))] for i in range(40)]
```

```
1 train_amp_f = [amplitude_lists[i*13099:(13099*i+9169)] for i in range(40)]
2 test_amp_f = [amplitude_lists[(i*13099)+9169:(13099*(i+1))] for i in range(40)]
```

```
1 train_freq = []
2 test_freq = []
3 train_amp = []
4 test_amp = []
5 for i in range(40):
6     train_freq.extend(train_fft_f[i])
7     test_freq.extend(test_fft_f[i])
8     train_amp.extend(train_amp_f[i])
9     test_amp.extend(test_amp_f[i])
```

```
1 train_freq_amp = np.hstack((train_freq, train_amp))
2 test_freq_amp = np.hstack((test_freq, test_amp))
```

```
1 # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our train
2 train_features =  [tsne_feature[i*13099:(13099*i+9169)] for i in range(0,40)]
3 # temp = [0]*(12955 - 9068)
4 test_features = [tsne_feature[(13099*(i))+9169:13099*(i+1)] for i in range(0,40)]
```

```
1 print("Number of data clusters",len(train_features), "Number of data points in trian data"
2 print("Number of data clusters",len(train_features), "Number of data points in test data",
```

> Number of data clusters 40 Number of data points in trian data 9169 Each data point con1
> Number of data clusters 40 Number of data points in test data 3930 Each data point conta

```
1 # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our train
2 tsne_train_flat_lat = [i[:9169] for i in tsne_lat]
3 tsne_train_flat_lon = [i[:9169] for i in tsne_lon]
4 tsne_train_flat_weekday = [i[:9169] for i in tsne_weekday]
5 tsne_train_flat_output = [i[:9169] for i in output]
6 tsne_train_flat_exp_avg = [i[:9169] for i in predict_list]
7 tsne_train_flat_tri_exp = [i[:9169] for i in tri_exp]
8 tsne_train_flat_exp_smt = [i[:9169] for i in exp_smt]
```

```
1 # extracting the rest of the timestamp values i.e 30% of 12956 (total timestamps) for our
2 tsne test flat lat = [i[9169:] for i in tsne lat]
```

```
2 tsne_test_flat_lat = [i[9169:] for i in tsne_lat]
3 tsne_test_flat_lon = [i[9169:] for i in tsne_lon]
4 tsne_test_flat_weekday = [i[9169:] for i in tsne_weekday]
5 tsne_test_flat_output = [i[9169:] for i in output]
6 tsne_test_flat_exp_avg = [i[9169:] for i in predict_list]
7 tsne_test_flat_tri_exp = [i[9169:] for i in tri_exp]
8 tsne_test_flat_exp_smt = [i[9169:] for i in exp_smt]
```

```
1 # the above contains values in the form of list of lists (i.e. list of values of each regi
2 train_new_features = []
3 for i in range(0,40):
4     train_new_features.extend(train_features[i])
5 test_new_features = []
6 for i in range(0,40):
7     test_new_features.extend(test_features[i])
```

```
 1 # converting lists of lists into sinle list i.e flatten
 2 # a  = [[1,2,3,4],[4,6,7,8]]
 3 # print(sum(a,[]))
 4 # [1, 2, 3, 4, 4, 6, 7, 8]
 5
 6 tsne_train_lat = sum(tsne_train_flat_lat, [])
 7 tsne_train_lon = sum(tsne_train_flat_lon, [])
 8 tsne_train_weekday = sum(tsne_train_flat_weekday, [])
 9 tsne_train_output = sum(tsne_train_flat_output, [])
10 tsne_train_exp_avg = sum(tsne_train_flat_exp_avg,[])
11 tsne_train_tri_exp = sum(tsne_train_flat_tri_exp,[])
12 tsne_train_exp_smt = sum(tsne_train_flat_exp_smt,[])
```

```
 1 # converting lists of lists into sinle list i.e flatten
 2 # a  = [[1,2,3,4],[4,6,7,8]]
 3 # print(sum(a,[]))
 4 # [1, 2, 3, 4, 4, 6, 7, 8]
 5
 6 tsne_test_lat = sum(tsne_test_flat_lat, [])
 7 tsne_test_lon = sum(tsne_test_flat_lon, [])
 8 tsne_test_weekday = sum(tsne_test_flat_weekday, [])
 9 tsne_test_output = sum(tsne_test_flat_output, [])
10 tsne_test_exp_avg = sum(tsne_test_flat_exp_avg,[])
11 tsne_test_tri_exp = sum(tsne_test_flat_tri_exp,[])
12 tsne_test_exp_smt = sum(tsne_test_flat_exp_smt,[])
```

```
1 len(tsne_train_exp_smt)
```

    366760

```
1 # Preparing the data frame for our train data
2 columns = ['ft_5','ft_4','ft_3','ft_2','ft_1']
3 df_train = pd.DataFrame(data=train_new_features, columns=columns)
```

```
 4 df_train['lat'] = tsne_train_lat
 5 df_train['lon'] = tsne_train_lon
 6 df_train['weekday'] = tsne_train_weekday
 7 df_train['exp_avg'] = tsne_train_exp_avg
 8 df_train['tri_exp'] =tsne_train_tri_exp
 9 df_train['exp_smt'] =tsne_train_exp_smt
10 print(df_train.shape)
```

> (366760, 11)

```
1 #train dataframe
2 columns = ['freq1', 'freq2','freq3','freq4','freq5', 'Amp1', 'Amp2', 'Amp3', 'Amp4', 'Amp5
3 Train_DF = pd.DataFrame(data = train_freq_amp, columns = columns)
```

```
1 Train_DF.shape
```

> (366760, 10)

```
1 df_train=pd.concat((df_train,Train_DF),axis=1)
2 print(df_train.shape)
3 print(df_train.columns)
```

> (366760, 21)
>     Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'lat', 'lon', 'weekday',
>            'exp_avg', 'tri_exp', 'exp_smt', 'freq1', 'freq2', 'freq3', 'freq4',
>            'freq5', 'Amp1', 'Amp2', 'Amp3', 'Amp4', 'Amp5'],
>         dtype='object')

```
 1 # Preparing the data frame for our train data
 2 columns = ['ft_5','ft_4','ft_3','ft_2','ft_1']
 3 df_test = pd.DataFrame(data=test_new_features, columns=columns)
 4 df_test['lat'] = tsne_test_lat
 5 df_test['lon'] = tsne_test_lon
 6 df_test['weekday'] = tsne_test_weekday
 7 df_test['exp_avg'] = tsne_test_exp_avg
 8 df_test['tri_exp'] = tsne_test_tri_exp
 9 df_test['exp_smt'] = tsne_test_exp_smt
10 print(df_test.shape)
```

> (157200, 11)

```
1 #train dataframe
2 columns = ['freq1', 'freq2','freq3','freq4','freq5', 'Amp1', 'Amp2', 'Amp3', 'Amp4', 'Amp5
3 Test_DF = pd.DataFrame(data = test_freq_amp, columns = columns)
```

```
1 Test_DF.shape
```

```
                        (157200  10)
   1 df_test=pd.concat((df_test,Test_DF),axis=1)
   2 print(df_test.shape)
   3 print(df_test.columns)
```

```
   (157200, 21)
   Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'lat', 'lon', 'weekday',
          'exp_avg', 'tri_exp', 'exp_smt', 'freq1', 'freq2', 'freq3', 'freq4',
          'freq5', 'Amp1', 'Amp2', 'Amp3', 'Amp4', 'Amp5'],
         dtype='object')
```

```
   1 df_train.head()
```

|   | ft_5 | ft_4 | ft_3 | ft_2 | ft_1 | lat | lon | weekday | exp_avg | tri_exp | exp_sm |
|---|------|------|------|------|------|-----------|------------|---------|---------|----------|--------|
| 0 | 0 | 0 | 0 | 0 | 0 | 40.776228 | -73.982119 | 4 | 0 | 4.830344 | 0. |
| 1 | 0 | 0 | 0 | 0 | 0 | 40.776228 | -73.982119 | 4 | 0 | 4.443229 | 0. |
| 2 | 0 | 0 | 0 | 0 | 0 | 40.776228 | -73.982119 | 4 | 0 | 3.147928 | 0. |
| 3 | 0 | 0 | 0 | 0 | 0 | 40.776228 | -73.982119 | 4 | 0 | 2.386990 | 0. |
| 4 | 0 | 0 | 0 | 0 | 0 | 40.776228 | -73.982119 | 4 | 0 | 1.412932 | 0. |

## Using SKlearn Linear Regression

```
   1 # find more about LinearRegression function here http://scikit-learn.org/stable/modules/ge
   2 # -------------------------
   3 # default paramters
   4 # sklearn.linear_model.LinearRegression(fit_intercept=True, normalize=False, copy_X=True,
   5
   6 # some of methods of LinearRegression()
   7 # fit(X, y[, sample_weight])  Fit linear model.
   8 # get_params([deep])  Get parameters for this estimator.
   9 # predict(X)  Predict using the linear model
  10 # score(X, y[, sample_weight])  Returns the coefficient of determination R^2 of the predic
  11 # set_params(**params)  Set the parameters of this estimator.
  12 # ----------------------
  13 # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geom
  14 # ----------------------
  15
  16 from sklearn.linear_model import LinearRegression
  17 lr_reg=LinearRegression().fit(df_train, tsne_train_output)
  18
  19 y_pred = lr_reg.predict(df_test)
  20 lr_test_predictions = [round(value) for value in y_pred]
  21 y_pred = lr_reg.predict(df_train)
  22 lr_train_predictions = [round(value) for value in y_pred]
```

```
1 mean_absolute_error(tsne_train_output, lr_train_predictions)/(sum(tsne_train_output)/len(t
```

    0.08589237732000982

```
1 mean_absolute_error(tsne_test_output, lr_test_predictions)/(sum(tsne_test_output)/len(tsne
```

    0.08388462129045682

## using SGD sq-loss implementaiton

```
1 from sklearn.preprocessing import StandardScaler
2 from sklearn.linear_model import SGDRegressor
3 from sklearn.model_selection import GridSearch
4
5
6 clf = SGDRegressor(loss = "squared_loss", penalty = "l2")
7 alpha = [10 ** x for x in range(-10, 5)]
8 hyper_parameter = {"alpha": alpha}
9
10 best_parameter = GridSearchCV(clf, hyper_parameter, verbose =5, scoring = "neg_mean_absolu
11 best_parameter.fit(df_train,tsne_train_output )
12 alpha = best_parameter.best_params_["alpha"]
13
14 #applying linear regression with best hyper-parameter
15 clf = SGDRegressor(loss = "squared_loss", penalty = "l2", alpha = alpha)
16 clf.fit(df_train, tsne_train_output)
17 lr_train_predictions = clf.predict(df_train)
18 lr_test_predictions = clf.predict(df_test)
19
```

    Fitting 3 folds for each of 15 candidates, totalling 45 fits
    [Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
    [Parallel(n_jobs=-1)]: Done  14 tasks      | elapsed: 13.8min
    [Parallel(n_jobs=-1)]: Done  45 out of  45 | elapsed: 42.6min finished

```
1 best_parameter.best_params_
```

    {'alpha': 0.0001}

```
1 from sklearn.linear_model import SGDRegressor
2 clf = SGDRegressor(loss = "squared_loss", penalty = "l2",alpha=0.0001)
3 clf.fit(df_train, tsne_train_output)
4 y_pred= clf.predict(df_train)
5 lr_train_predictions =[round(value) for value in y_pred]
6 y_pred= clf.predict(df_test)
7 lr_test_predictions =[round(value) for value in y_pred]
```

8

```
1 mean_absolute_error(tsne_train_output, lr_train_predictions)/(sum(tsne_train_output)/len(t
```

1.0249292872453195e+17

```
1 mean_absolute_error(tsne_test_output, lr_test_predictions)/(sum(tsne_test_output)/len(tsne
```

9.31162189176855e+16

## Using Random Forest Regressor

```
1 depth=[i for i in np.arange(9,12)]
2 estimators= [50,100,300,800]
3 param = {'max_depth':depth,'n_estimators':estimators}
4
5 from sklearn.model_selection import RandomizedSearchCV
6 from sklearn.ensemble import RandomForestRegressor
7 clf=RandomForestRegressor()
8 temp_gscv= RandomizedSearchCV(clf,param,cv=3,verbose=5,n_jobs=-1,scoring='neg_mean_absolut
9 temp_gscv.fit(df_train,tsne_train_output)
```

```
Fitting 3 folds for each of 10 candidates, totalling 30 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done  14 tasks      | elapsed: 45.5min
[Parallel(n_jobs=-1)]: Done  30 out of  30 | elapsed: 152.8min finished
RandomizedSearchCV(cv=3, error_score=nan,
                   estimator=RandomForestRegressor(bootstrap=True,
                                                   ccp_alpha=0.0,
                                                   criterion='mse',
                                                   max_depth=None,
                                                   max_features='auto',
                                                   max_leaf_nodes=None,
                                                   max_samples=None,
                                                   min_impurity_decrease=0.0,
                                                   min_impurity_split=None,
                                                   min_samples_leaf=1,
                                                   min_samples_split=2,
                                                   min_weight_fraction_leaf=0.0,
                                                   n_estimators=100,
                                                   n_jobs=None, oob_score=False,
                                                   random_state=None, verbose=0,
                                                   warm_start=False),
                   iid='deprecated', n_iter=10, n_jobs=-1,
                   param_distributions={'max_depth': [9, 10, 11],
                                        'n_estimators': [50, 100, 300, 800]},
                   pre_dispatch='2*n_jobs', random_state=None, refit=True,
                   return_train_score=True, scoring='neg_mean_absolute_error',
                   verbose=5)
```

```
1 temp_gscv.best_estimator_
```

```
RandomForestRegressor(bootstrap=True, ccp_alpha=0.0, criterion='mse',
                      max_depth=11, max_features='auto', max_leaf_nodes=None,
                      max_samples=None, min_impurity_decrease=0.0,
                      min_impurity_split=None, min_samples_leaf=1,
                      min_samples_split=2, min_weight_fraction_leaf=0.0,
                      n_estimators=800, n_jobs=None, oob_score=False,
                      random_state=None, verbose=0, warm_start=False)
```

```
1 from sklearn.ensemble import RandomForestRegressor
2 clf=RandomForestRegressor(max_depth=11,n_estimators=800)
3 clf.fit(df_train,tsne_train_output)
4 rndf_train_predictions = clf.predict(df_train)
5 rndf_test_predictions = clf.predict(df_test)
```

```
1 y_pred= clf.predict(df_train)
2 rndf_train_predictions =[round(value) for value in y_pred]
3 y_pred= clf.predict(df_test)
4 rndf_test_predictions =[round(value) for value in y_pred]
5
```

```
1 mean_absolute_error(tsne_train_output, rndf_train_predictions)/(sum(tsne_train_output)/len
```

```
0.08249763487484826
```

```
1 mean_absolute_error(tsne_test_output, rndf_test_predictions)/(sum(tsne_test_output)/len(ts
```

```
0.08531729733882347
```

## Using XgBoost Regressor

```
 1 #hyper-parameter tuning
 2 hyper_parameter = {"max_depth":[2, 3, 4], "n_estimators":[50,100,400,800,1200]}
 3 clf = xgb.XGBRegressor(learning_rate =0.1,
 4  min_child_weight=3,
 5  gamma=0,
 6  subsample=0.8,
 7  reg_alpha=200, reg_lambda=200,
 8  colsample_bytree=0.8,nthread=4)
 9
10 best_parameter = RandomizedSearchCV(clf, hyper_parameter, scoring = "neg_mean_absolute_err
11 best_parameter.fit(df_train, tsne_train_output)
12 estimators = best_parameter.best_params_["n_estimators"]
13 depth = best_parameter.best_params_["max_depth"]
```

```
Fitting 3 folds for each of 10 candidates, totalling 30 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 2 concurrent workers.
[Parallel(n_jobs=-1)]: Done  14 tasks      | elapsed: 18.0min
[Parallel(n_jobs=-1)]: Done  30 out of  30 | elapsed: 30.4min finished
[17:17:16] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now de
```

```
1 print(depth)
2 print(estimators)
```

```
3
800
```

```
1 clf = xgb.XGBRegressor(max_depth = 3, n_estimators = 800)
2 clf.fit(df_train, tsne_train_output)
3 y_pred = clf.predict(df_test)
4 xgb_test_predictions = [round(value) for value in y_pred]
5 y_pred = clf.predict(df_train)
6 xgb_train_predictions = [round(value) for value in y_pred]
```

```
[17:20:23] WARNING: /workspace/src/objective/regression_obj.cu:152: reg:linear is now de
```

```
1 mean_absolute_error(tsne_train_output, xgb_train_predictions)/(sum(tsne_train_output)/len(
```

```
0.0851921237782352
```

```
1 mean_absolute_error(tsne_test_output, xgb_test_predictions)/(sum(tsne_test_output)/len(tsn
```

```
0.08333327793943827
```

## Calculating the error metric values for various models

```
1 train_mape=[]
2 test_mape=[]
3
4 train_mape.append((mean_absolute_error(tsne_train_output,df_train['ft_1'].values))/(sum(ts
5 train_mape.append((mean_absolute_error(tsne_train_output,df_train['exp_avg'].values))/(sum
6 train_mape.append((mean_absolute_error(tsne_train_output,rndf_train_predictions))/(sum(tsn
7 train_mape.append((mean_absolute_error(tsne_train_output, xgb_train_predictions))/(sum(tsn
8 train_mape.append((mean_absolute_error(tsne_train_output, lr_train_predictions))/(sum(tsne
9
10 test_mape.append((mean_absolute_error(tsne_test_output, df_test['ft_1'].values))/(sum(tsne
11 test_mape.append((mean_absolute_error(tsne_test_output, df_test['exp_avg'].values))/(sum(t
12 test_mape.append((mean_absolute_error(tsne_test_output, rndf_test_predictions))/(sum(tsne_
13 test_mape.append((mean_absolute_error(tsne_test_output, xgb_test_predictions))/(sum(tsne_t
14 test_mape.append((mean_absolute_error(tsne_test_output, lr_test_predictions))/(sum(tsne_te
```

```
1 print ("Error Metric Matrix (Tree Based Regression Methods) -  MAPE")
2 print ("-------------------------------------------------------------------
```

```
3 print ("Baseline Model -                                        Train: ",train_ma
4 print ("Exponential Averages Forecasting -                      Train: ",train_ma
5 print ("XGboost Regression -                                    Train: ",train_ma
6 print ("SGD Implementation of Linear Regression -              Train: ",train_map
7 print ("Random Forest Regression -                              Train: ",train_ma
8 print ("Simple Linear Regression -                              Train: ",0.085892
9
```

```
Error Metric Matrix (Tree Based Regression Methods) -  MAPE
--------------------------------------------------------------------------------
Baseline Model -                                        Train:  14.87066699642(
Exponential Averages Forecasting -                      Train:  14.12160356090(
XGboost Regression -                                    Train:  8.519212377823!
SGD Implementation of Linear Regression -              Train:  1.02492928724531
Random Forest Regression -                              Train:  8.249763487484&
Simple Linear Regression -                              Train:  8.589237732000!
```

# Workflow

This Problem statement deals with predicting number of Taxi pickups an area in newyork may have in the upcoming 10 mins.

1. After acquiring the data and cleaning it for cases like lat and lon out side of Newyork city, checking for time consistency, fare, speed, and distance consisitency too, we bin the data in to 10 min time sessions based on the observations made that it would take a driver that much time to switch to nearby areas for pickups. We also segreggate new york into 40 sections.
2. This time binning and sectionization of data under data preparation gives us a clearer picture of data and how to predict the next numberof pickups for the next time bin of 10 min.
3. We simply first make use of moving ratios(avg and exp) forcasting methods then in corparae this features with last five pickups made. These featurizations gave us a MAPE within the range of 10 to 15 percent.
4. further we incorporate fourier feature transformation along with exponential smoothing and holt's winter featurization which further reduced the MAPE to around 8 to 9 percent.

## 5. XGBoost

works the best with this data closely followed by random forest and simple linear regresion without regularization.