# Taxi demand prediction in New York City

In [1]:

```python
#Importing Libraries
# pip3 install graphviz
#pip3 install dask
#pip3 install toolz
#pip3 install cloudpickle
# https://www.youtube.com/watch?v=ieW3G7ZzRZ0
# https://github.com/dask/dask-tutorial
# please do go through this python notebook: https://github.com/dask/dask-tutorial/blob/master/07_dataframe
import dask.dataframe as dd#similar to pandas

import pandas as pd#pandas to create small dataframes

# pip3 install foliun
# if this doesnt work refere install_folium.JPG in drive
import folium #open street map

# unix time: https://www.unixtimestamp.com/
import datetime #Convert to unix time

import time #Convert to unix time

# if numpy is not installed already : pip3 install numpy
import numpy as np#Do aritmetic operations on arrays

# matplotlib: used to plot graphs
import matplotlib
# matplotlib.use('nbagg') : matplotlib uses this protocall which makes plots more user intractive like zoom
matplotlib.use('nbagg')
import matplotlib.pylab as plt
import seaborn as sns#Plots
from matplotlib import rcParams#Size of plots

# this lib is used while we calculate the stight line distance between two (lat,lon) pairs in miles
import gpxpy.geo #Get the haversine distance

from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
import math
import pickle
import os

# 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\\bin'
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  warnings.filterwarnings("ignore")
```

# Data Information

Ge the data from : 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 |

```
In [2]:   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')
```
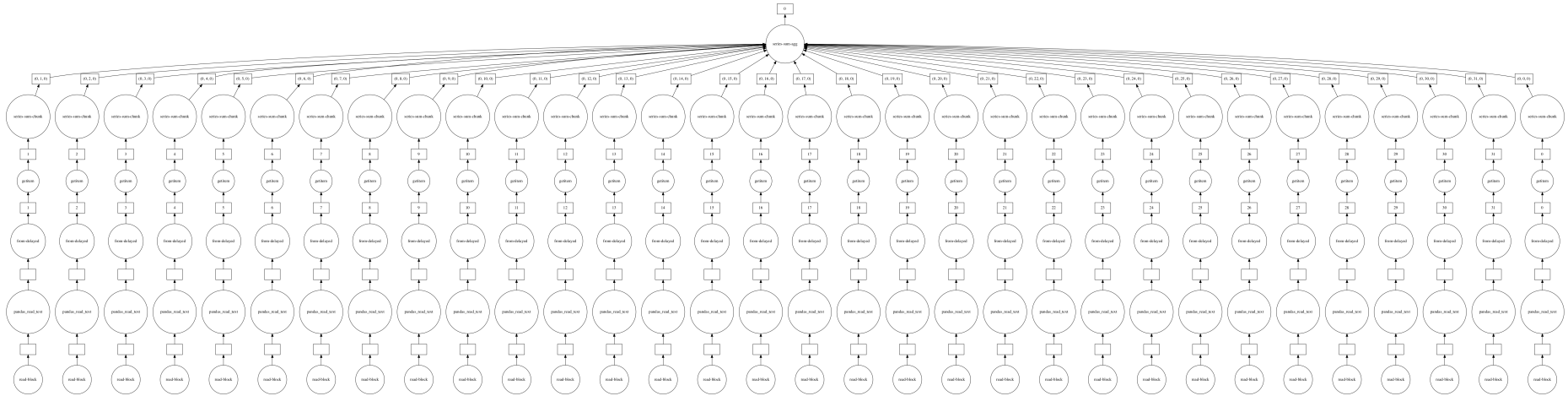
```
In [3]:   1   # However unlike Pandas, operations on dask.dataframes don't trigger immediate computation,
          2   # instead they add key-value pairs to an underlying Dask graph. Recall that in the diagram below,
          3   # circles are operations and rectangles are results.
          4
          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 drive
          7   month.visualize()
```

Out[3]:

In [4]: `1  month.fare_amount.sum().visualize()`

Out[4]:



# 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 hour
 and overnight charges.</td>
</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 bei
ng levied in 2015.</td>
</tr>
<tr>
    <td>Tip_amount</td>
    <td>Tip amount — This field is automatically populated for credit card tips.Cash tips are not inclu
ded.</td>
</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 |

| | 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. |
|---|---|
| Store_and_fwd_flag | 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

```
In [5]:    1  #table below shows few datapoints along with all our features
           2  month.head(5)
```

Out[5]:

| | VendorID | tpep_pickup_datetime | tpep_dropoff_datetime | passenger_count | trip_distance | pickup_longitude | pickup_latitude | RateCodeID | store_and_f |
|---|---|---|---|---|---|---|---|---|---|
| **0** | 2 | 2015-01-15 19:05:39 | 2015-01-15 19:23:42 | 1 | 1.59 | -73.993896 | 40.750111 | 1 | |
| **1** | 1 | 2015-01-10 20:33:38 | 2015-01-10 20:53:28 | 1 | 3.30 | -74.001648 | 40.724243 | 1 | |
| **2** | 1 | 2015-01-10 20:33:38 | 2015-01-10 20:43:41 | 1 | 1.80 | -73.963341 | 40.802788 | 1 | |
| **3** | 1 | 2015-01-10 20:33:39 | 2015-01-10 20:35:31 | 1 | 0.50 | -74.009087 | 40.713818 | 1 | |
| **4** | 1 | 2015-01-10 20:33:39 | 2015-01-10 20:52:58 | 1 | 3.00 | -73.971176 | 40.762428 | 1 | |

## 1. Pickup Latitude and Pickup Longitude

It is inferred from the source https://www.flickr.com/places/info/2459115 (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.

```
In [6]:   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_locations
          3  outlier_locations = month[((month.pickup_longitude <= -74.15) | (month.pickup_latitude <= 40.5774)| \
          4                    (month.pickup_longitude >= -73.7004) | (month.pickup_latitude >= 40.9176))]
          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 maps and plots
         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 more time
         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
```

Out[6]:

**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 (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.

```
In [7]:   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_locations
          3  outlier_locations = month[((month.dropoff_longitude <= -74.15) | (month.dropoff_latitude <= 40.5774)| \
          4                   (month.dropoff_longitude >= -73.7004) | (month.dropoff_latitude >= 40.9176))]
          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 maps and plots
         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 more time
         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['dropoff_latitude'],j['dropoff_longitude']))).add_to(map_osm)
         18  map_osm
```

Out[7]:

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

## 3. Trip Durations:

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

```
In [8]:   1  #The timestamps are converted to unix so as to get duration(trip-time) & speed also pickup-times in unix ar
          2
          3  # in out data we have time in the formate "YYYY-MM-DD HH:MM:SS" we convert thiss sting to python time forma
          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_latitude','dropoff_long
         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    dropoff_latitud
         40  #    1                    1.59         -73.993896          40.750111        -73.974785           40.750618
         41  #    1                    3.30         -74.001648          40.724243        -73.994415           40.759109
```

```
42  #   1                    1.80        -73.963341        40.802788        -73.951820        40.824413
43  #   1                    0.50        -74.009087        40.713818        -74.004326        40.719986
44  #   1                    3.00        -73.971176        40.762428        -74.004181        40.742653
45  frame_with_durations = return_with_trip_times(month)
```

In [9]:
```python
# the skewed box plot shows us the presence of outliers
sns.boxplot(y="trip_times", data =frame_with_durations)
plt.show()
```

```
In [10]:  1  #calculating 0-100th percentile to find a the correct percentile value for removal of outliers
          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 5.383333333333334
30 percentile value is 6.816666666666666
40 percentile value is 8.3
50 percentile value is 9.95
60 percentile value is 11.866666666666667
70 percentile value is 14.283333333333333
80 percentile value is 17.633333333333333
90 percentile value is 23.45
100 percentile value is  548555.6333333333
```

```
In [19]:   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  548555.6333333333
```

```
In [12]:   1  #removing data based on our analysis and TLC regulations
           2  frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_times>1) & (frame_with_durati
```

In [13]:

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

In [14]:
```python
#pdf of trip-times after removing the outliers
sns.FacetGrid(frame_with_durations_modified,size=6) \
      .map(sns.kdeplot,"trip_times") \
      .add_legend();
plt.show();
```
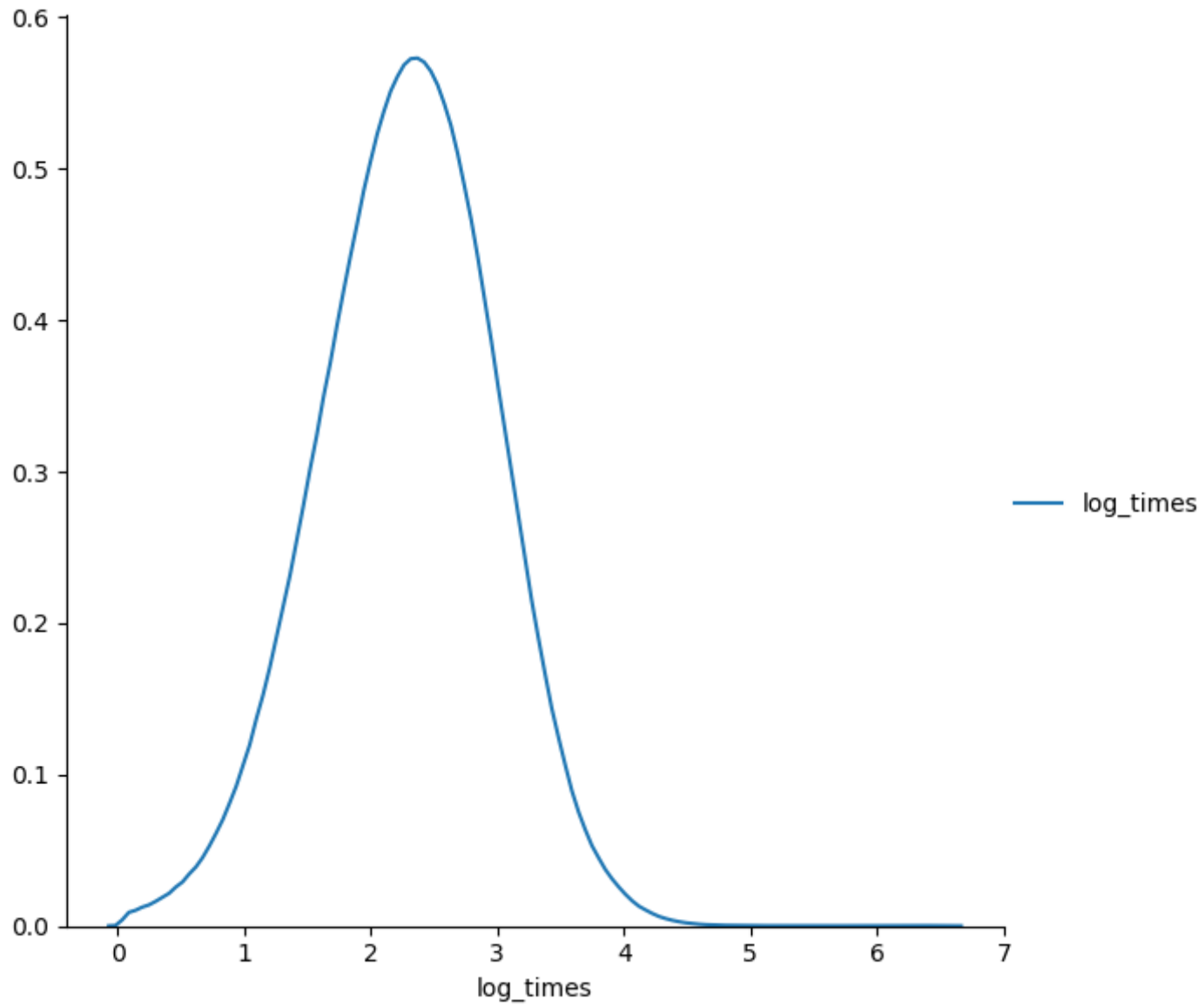
Zoom to rectangle

In [21]:

```
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_modified['trip_times'
```

In [22]:

```python
#pdf of log-values
sns.FacetGrid(frame_with_durations_modified,size=6) \
        .map(sns.kdeplot,"log_times") \
        .add_legend();
plt.show();
```
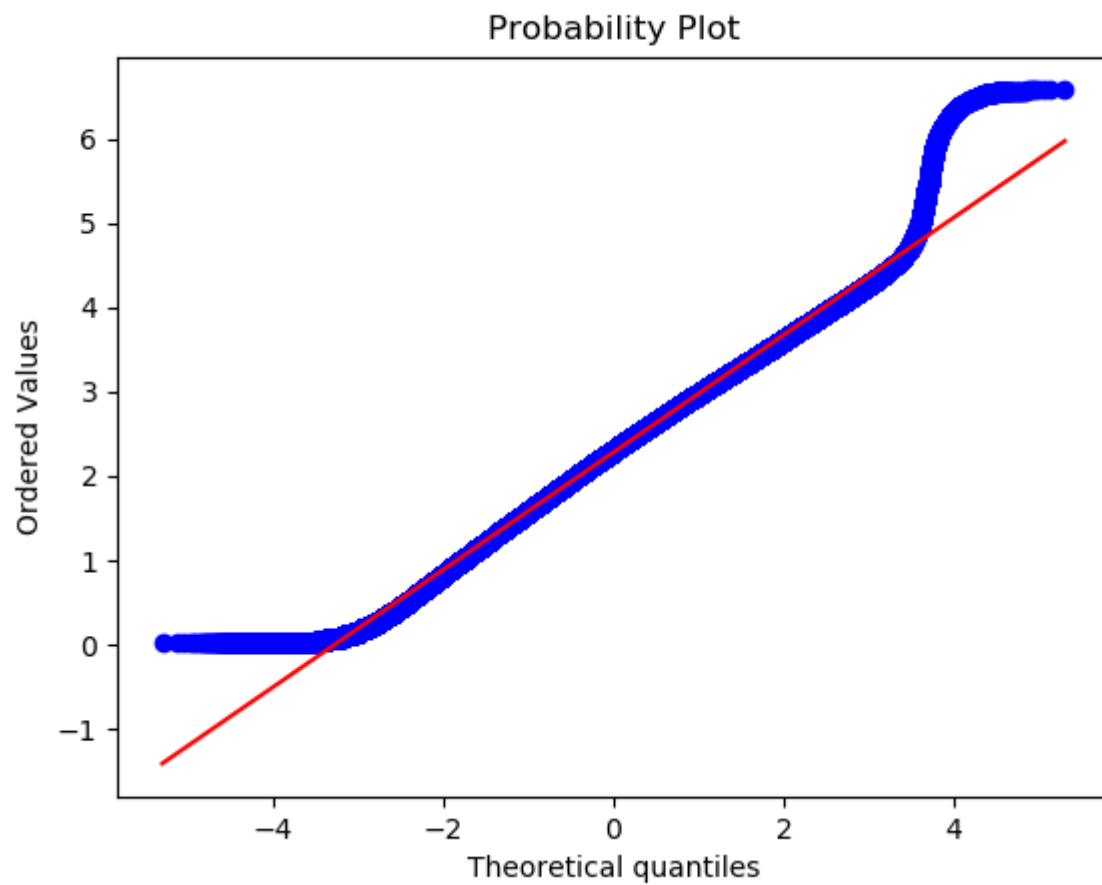
**Figure 5**

In [23]:
```python
#Q-Q plot for checking if trip-times is log-normal
import scipy
scipy.stats.probplot(frame_with_durations_modified['log_times'].values, plot=plt)
plt.show()
```

**Figure 6**



Probability Plot

## 4. Speed

```python
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']/frame_with_dura
4  sns.boxplot(y="Speed", data =frame_with_durations_modified)
5  plt.show()
```

```python
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])
```

```python
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])
```

```python
1  #calculating speed values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100
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])
```

```
In [ ]:   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_with_durations.S
```

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

**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

```
In [ ]:   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()
```

```
In [ ]:   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])
```

```
In [ ]:   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])
```

```
In [ ]:   1   #calculating trip distance values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100
          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])
```

```
In [ ]:   #removing further outliers based on the 99.9th percentile value
          frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_distance>0) & (frame_with_durati
```

```
In [29]:   1  #box-plot after removal of outliers
           2  sns.boxplot(y="trip_distance", data = frame_with_durations_modified)
           3  plt.show()
```

**Figure 7**

## 5. Total Fare

```python
1  # up to now we have removed the outliers based on trip durations, cab speeds, and trip distances
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()
```

```python
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])
```

```python
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])
```

```python
1  #calculating total fare amount values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,1
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])
```

**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

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

```
In [ ]:  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 values
         3  plt.plot(var[-3:])
         4  plt.show()
```

```
In [ ]:  1  #now looking at values not including the last two points we again find a drastic increase at around 1000 fa
         2  # we plot last 50 values excluding last two values
         3  plt.plot(var[-50:-2])
         4  plt.show()
```

# Remove all outliers/erronous points.

In [31]:
```python
#removing all outliers based on our univariate analysis above
def remove_outliers(new_frame):


    a = new_frame.shape[0]
    print ("Number of pickup records = ",a)
    temp_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_frame.dropoff_longitude <= -73.70
                           (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitude <= 40.9176)) &
                           ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_latitude >= 40.5774)& \
                           (new_frame.pickup_longitude <= -73.7004) & (new_frame.pickup_latitude <= 40.9176))]
    b = temp_frame.shape[0]
    print ("Number of outlier coordinates lying outside NY boundaries:",(a-b))


    temp_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times < 720)]
    c = temp_frame.shape[0]
    print ("Number of outliers from trip times analysis:",(a-c))


    temp_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_distance < 23)]
    d = temp_frame.shape[0]
    print ("Number of outliers from trip distance analysis:",(a-d))

    temp_frame = new_frame[(new_frame.Speed <= 65) & (new_frame.Speed >= 0)]
    e = temp_frame.shape[0]
    print ("Number of outliers from speed analysis:",(a-e))

    temp_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_amount >0)]
    f = temp_frame.shape[0]
    print ("Number of outliers from fare analysis:",(a-f))


    new_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_frame.dropoff_longitude <= -73.700
                          (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitude <= 40.9176)) &
                          ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_latitude >= 40.5774)& \
                          (new_frame.pickup_longitude <= -73.7004) & (new_frame.pickup_latitude <= 40.9176))]

    new_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times < 720)]
    new_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_distance < 23)]
    new_frame = new_frame[(new_frame.Speed < 45.31) & (new_frame.Speed > 0)]
    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
```

In [32]:
```
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_durations_outlier
```

```
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

In [33]:

```python
#trying different cluster sizes to choose the right K in K-means
coords = frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']].values
neighbours=[]

def find_min_distance(cluster_centers, cluster_len):
    nice_points = 0
    wrong_points = 0
    less2 = []
    more2 = []
    min_dist=1000
    for i in range(0, cluster_len):
        nice_points = 0
        wrong_points = 0
        for j in range(0, cluster_len):
            if j!=i:
                distance = gpxpy.geo.haversine_distance(cluster_centers[i][0], cluster_centers[i][1],cluste
                min_dist = min(min_dist,distance/(1.60934*1000))
                if (distance/(1.60934*1000)) <= 2:
                    nice_points +=1
                else:
                    wrong_points += 1
        less2.append(nice_points)
        more2.append(wrong_points)
    neighbours.append(less2)
    print ("On choosing a cluster size of ",cluster_len,"\nAvg. Number of Clusters within the vicinity (i.e

def find_clusters(increment):
    kmeans = MiniBatchKMeans(n_clusters=increment, batch_size=10000,random_state=42).fit(coords)
    frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_durations_outliers_
    cluster_centers = kmeans.cluster_centers_
    cluster_len = len(cluster_centers)
    return cluster_centers, cluster_len

# we need to choose number of clusters so that, there are more number of cluster regions
#that are close to any cluster center
# and make sure that the minimum inter cluster should not be very less
for increment in range(10, 100, 10):
    cluster_centers, cluster_len = find_clusters(increment)
    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 =   1.0945442325142543
---
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 =   0.7131298007387813
---
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 =   0.5185088176172206
---
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 =   0.5069768450363972
---
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 =   0.365363025983595
---
On choosing a cluster size of   60
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 14.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 46.0
```

```
Min inter-cluster distance =  0.34704283494187155
---
On choosing a cluster size of  70
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 16.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 54.0
Min inter-cluster distance =  0.30502203163244707
---
On choosing a cluster size of  80
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 18.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 62.0
Min inter-cluster distance =  0.29220324531738534
---
On choosing a cluster size of  90
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 21.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 69.0
Min inter-cluster distance =  0.18257992857034985
---
```

## 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

```
In [34]:  1  # if check for the 50 clusters you can observe that there are two clusters with only 0.3 miles apart from e
          2  # so we choose 40 clusters for solve the further problem
          3
          4  # Getting 40 clusters using the kmeans
          5  kmeans = MiniBatchKMeans(n_clusters=40, batch_size=10000,random_state=0).fit(coords)
          6  frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_durations_outliers_remo
```

## Plotting the cluster centers:

In [35]:

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

Out[35]:

"Leaflet (https://leafletjs.com) | Map tiles by Stamen Design
(http://openstreetmap.org), under ODbL (http://www.openstree

## Plotting the clusters:

```
In [36]:   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[:100000], s=10, lw=0,
           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)
```

**Figure 8**



# Time-binning

```
In [42]:    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 it to est
           24       tenminutewise_binned_unix_pickup_times=[(int((i-start_pickup_unix)/600)+33) for i in unix_pickup_times]
           25       frame['pickup_bins'] = np.array(tenminutewise_binned_unix_pickup_times)
           26       return frame
```

```
In [43]:    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_durations_outliers_remov
            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']].groupby(['pickup_cluste
```

```
In [44]:    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()
```

Out[44]:

| | passenger_count | trip_distance | pickup_longitude | pickup_latitude | dropoff_longitude | dropoff_latitude | total_amount | trip_times | pickup_times | Sp |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 1.59 | -73.993896 | 40.750111 | -73.974785 | 40.750618 | 17.05 | 18.050000 | 1.421329e+09 | 5.28! |
| 1 | 1 | 3.30 | -74.001648 | 40.724243 | -73.994415 | 40.759109 | 17.80 | 19.833333 | 1.420902e+09 | 9.98: |
| 2 | 1 | 1.80 | -73.963341 | 40.802788 | -73.951820 | 40.824413 | 10.80 | 10.050000 | 1.420902e+09 | 10.74( |
| 3 | 1 | 0.50 | -74.009087 | 40.713818 | -74.004326 | 40.719986 | 4.80 | 1.866667 | 1.420902e+09 | 16.07: |
| 4 | 1 | 3.00 | -73.971176 | 40.762428 | -74.004181 | 40.742653 | 16.30 | 19.316667 | 1.420902e+09 | 9.31; |

In [45]:
```
1  # hear the trip_distance represents the number of pickups that are happend in that particular 10min intravel
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*60/10 =4464bins)
5  jan_2015_groupby.head()
```

Out[45]:

| | | trip_distance |
|---|---|---|
| pickup_cluster | pickup_bins | |
| 0 | 1 | 105 |
| | 2 | 199 |
| | 3 | 208 |
| | 4 | 141 |
| | 5 | 155 |

In [41]:

```python
# upto now we cleaned data and prepared data for the month 2015,

# now do the same operations for months Jan, Feb, March of 2016
# 1. get the dataframe which inlcudes only required colums
# 2. adding trip times, speed, unix time stamp of pickup_time
# 4. remove the outliers based on trip_times, speed, trip_duration, total_amount
# 5. add pickup_cluster to each data point
# 6. add pickup_bin (index of 10min intravel to which that trip belongs to)
# 7. group by data, based on 'pickup_cluster' and 'pickuo_bin'

# Data Preparation for the months of Jan,Feb and March 2016
def datapreparation(month,kmeans,month_no,year_no):

    print ("Return with trip times..")

    frame_with_durations = return_with_trip_times(month)

    print ("Remove outliers..")
    frame_with_durations_outliers_removed = remove_outliers(frame_with_durations)

    print ("Estimating clusters..")
    frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_durations_outliers_
    #frame_with_durations_outliers_removed_2016['pickup_cluster'] = kmeans.predict(frame_with_durations_out

    print ("Final groupbying..")
    final_updated_frame = add_pickup_bins(frame_with_durations_outliers_removed,month_no,year_no)
    final_groupby_frame = final_updated_frame[['pickup_cluster','pickup_bins','trip_distance']].groupby(['p

    return final_updated_frame,final_groupby_frame

month_jan_2016 = dd.read_csv('yellow_tripdata_2016-01.csv')
month_feb_2016 = dd.read_csv('yellow_tripdata_2016-02.csv')
month_mar_2016 = dd.read_csv('yellow_tripdata_2016-03.csv')

jan_2016_frame,jan_2016_groupby = datapreparation(month_jan_2016,kmeans,1,2016)
feb_2016_frame,feb_2016_groupby = datapreparation(month_feb_2016,kmeans,2,2016)
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
```

# Smoothing

In [46]:
```python
# Gets the unique bins where pickup values are present for each each reigion

# for each cluster region we will collect all the indices of 10min intravels in which the pickups are happen
# we got an observation that there are some pickpbins that doesnt have any pickups
def return_unq_pickup_bins(frame):
    values = []
    for i in range(0,40):
        new = frame[frame['pickup_cluster'] == i]
        list_unq = list(set(new['pickup_bins']))
        list_unq.sort()
        values.append(list_unq)
    return values
```

In [47]:
```python
# for every month we get all indices of 10min intravels in which atleast one pickup got happened

#jan
jan_2015_unique = return_unq_pickup_bins(jan_2015_frame)
jan_2016_unique = return_unq_pickup_bins(jan_2016_frame)

#feb
feb_2016_unique = return_unq_pickup_bins(feb_2016_frame)

#march
mar_2016_unique = return_unq_pickup_bins(mar_2016_frame)
```

In [48]:
```python
# for each cluster number of 10min intravels with 0 pickups
for i in range(40):
    print("for the ",i,"th cluster number of 10min intavels with zero pickups: ",4464 - len(set(jan_2015_un
    print('-'*60)
```

```
for the  0 th cluster number of 10min intavels with zero pickups:  41
------------------------------------------------------------
for the  1 th cluster number of 10min intavels with zero pickups:  1986
------------------------------------------------------------
for the  2 th cluster number of 10min intavels with zero pickups:  30
------------------------------------------------------------
for the  3 th cluster number of 10min intavels with zero pickups:  355
------------------------------------------------------------
for the  4 th cluster number of 10min intavels with zero pickups:  38
------------------------------------------------------------
for the  5 th cluster number of 10min intavels with zero pickups:  154
------------------------------------------------------------
for the  6 th cluster number of 10min intavels with zero pickups:  35
------------------------------------------------------------
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:  118
------------------------------------------------------------
for the  9 th cluster number of 10min intavels with zero pickups:  41
------------------------------------------------------------
for the  10 th cluster number of 10min intavels with zero pickups:  26
------------------------------------------------------------
for the  11 th cluster number of 10min intavels with zero pickups:  45
------------------------------------------------------------
for the  12 th cluster number of 10min intavels with zero pickups:  43
------------------------------------------------------------
for the  13 th cluster number of 10min intavels with zero pickups:  29
------------------------------------------------------------
for the  14 th cluster number of 10min intavels with zero pickups:  27
------------------------------------------------------------
for the  15 th cluster number of 10min intavels with zero pickups:  32
------------------------------------------------------------
for the  16 th cluster number of 10min intavels with zero pickups:  41
------------------------------------------------------------
for the  17 th cluster number of 10min intavels with zero pickups:  59
```

```
--------------------------------------------------------------
for the  18 th cluster number of 10min intavels with zero pickups:  1191
--------------------------------------------------------------
for the  19 th cluster number of 10min intavels with zero pickups:  1358
--------------------------------------------------------------
for the  20 th cluster number of 10min intavels with zero pickups:  54
--------------------------------------------------------------
for the  21 th cluster number of 10min intavels with zero pickups:  30
--------------------------------------------------------------
for the  22 th cluster number of 10min intavels with zero pickups:  30
--------------------------------------------------------------
for the  23 th cluster number of 10min intavels with zero pickups:  164
--------------------------------------------------------------
for the  24 th cluster number of 10min intavels with zero pickups:  36
--------------------------------------------------------------
for the  25 th cluster number of 10min intavels with zero pickups:  42
--------------------------------------------------------------
for the  26 th cluster number of 10min intavels with zero pickups:  32
--------------------------------------------------------------
for the  27 th cluster number of 10min intavels with zero pickups:  215
--------------------------------------------------------------
for the  28 th cluster number of 10min intavels with zero pickups:  37
--------------------------------------------------------------
for the  29 th cluster number of 10min intavels with zero pickups:  42
--------------------------------------------------------------
for the  30 th cluster number of 10min intavels with zero pickups:  1181
--------------------------------------------------------------
for the  31 th cluster number of 10min intavels with zero pickups:  43
--------------------------------------------------------------
for the  32 th cluster number of 10min intavels with zero pickups:  45
--------------------------------------------------------------
for the  33 th cluster number of 10min intavels with zero pickups:  44
--------------------------------------------------------------
for the  34 th cluster number of 10min intavels with zero pickups:  40
--------------------------------------------------------------
for the  35 th cluster number of 10min intavels with zero pickups:  43
--------------------------------------------------------------
for the  36 th cluster number of 10min intavels with zero pickups:  37
--------------------------------------------------------------
for the  37 th cluster number of 10min intavels with zero pickups:  322
--------------------------------------------------------------
for the  38 th cluster number of 10min intavels with zero pickups:  37
```

```
-------------------------------------------------------------
for the  39 th cluster number of 10min intavels with zero pickups:  44
-------------------------------------------------------------
```

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)

```
In [49]:  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
```

```
In [50]:    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 in the above mar
            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
           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/resolved
           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 bin if it exists
           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-bin value which
           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,hence we have no
           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:
```

```python
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_hand_limit-i)+2)
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 missing,hence we have
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
```
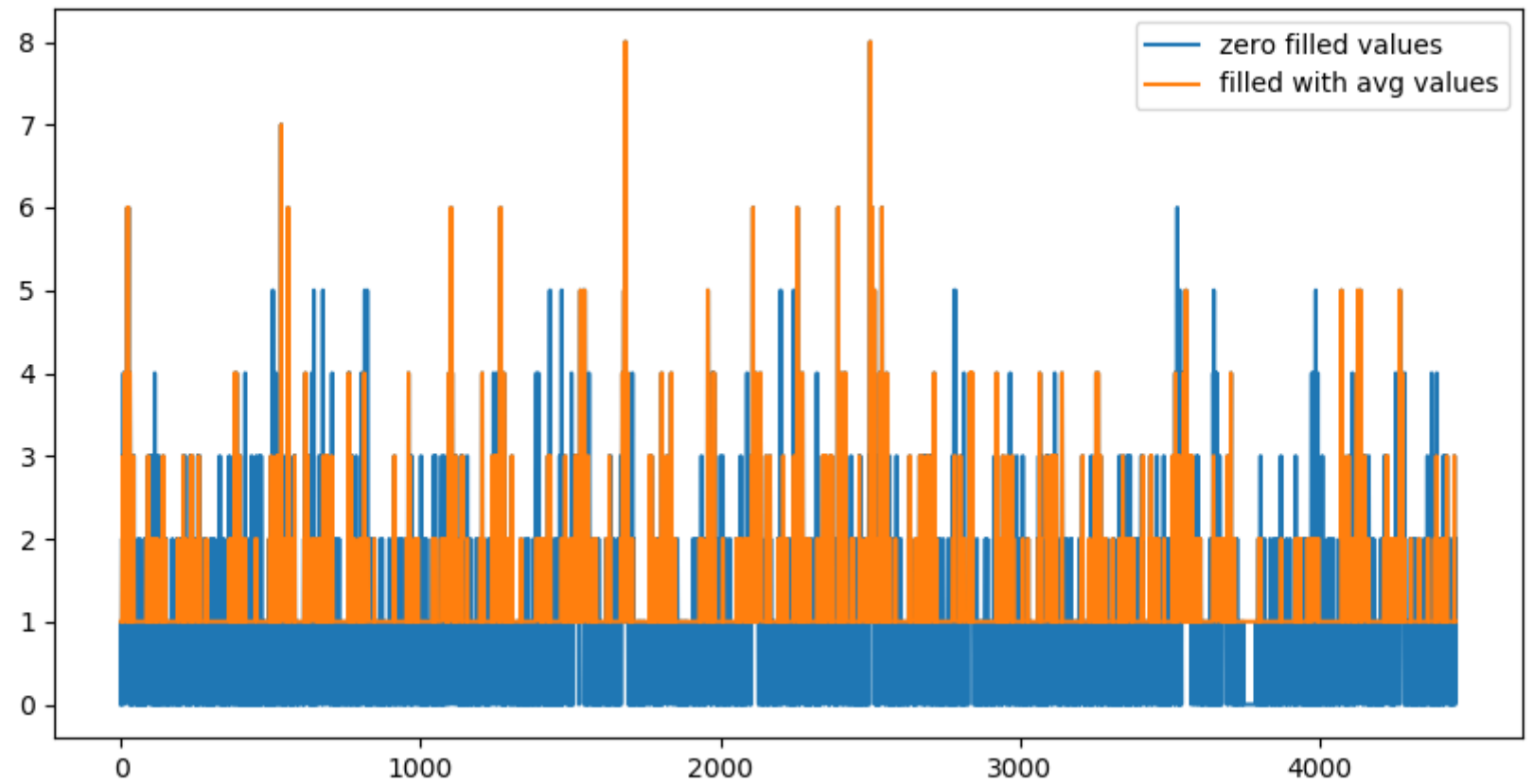
In [51]:
```python
#Filling Missing values of Jan-2015 with 0
# here in jan_2015_groupby dataframe the trip_distance represents the number of pickups that are happened
jan_2015_fill = fill_missing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)

#Smoothing Missing values of Jan-2015
jan_2015_smooth = smoothing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)
```

In [52]:
```python
# number of 10min indices for jan 2015= 24*31*60/10 = 4464
# number of 10min indices for jan 2016 = 24*31*60/10 = 4464
# number of 10min indices for feb 2016 = 24*29*60/10 = 4176
# number of 10min indices for march 2016 = 24*30*60/10 = 4320
# for each cluster we will have 4464 values, therefore 40*4464 = 178560 (length of the jan_2015_fill)
print("number of 10min intravels among all the clusters ",len(jan_2015_fill))
```

```
number of 10min intravels among all the clusters  178560
```

```
In [53]:   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()
```

**Figure 9**



x=541.175

```
In [ ]:    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 pickups that are
           4  # 10st 10min intravel, 0 pickups happened in 2nd 10mins intravel, 0 pickups happened in 3rd 10min intravel
           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 number of pickups
           8  # that are happened in the first 40min are same in both cases, but if you can observe that we looking at th
           9  # wheen you are using smoothing we are looking at the future number of pickups which might cause a data lea
          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.
```

```
In [65]:    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 storing them regio
            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 which represents the
           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)]+mar_2016_smoot
           23
           24  # print(len(regions_cum))
           25  # 40
           26  # print(len(regions_cum[0]))
           27  # 13104
```

```
In [68]:    1  regions_cum[0]
```

```
Out[68]:  [0,
           63,
           217,
           189,
           137,
           135,
           129,
           150,
           164,
           152,
           131,
           138,
           147,
           127,
           138,
           147,
           147,
           124,
           98,
```
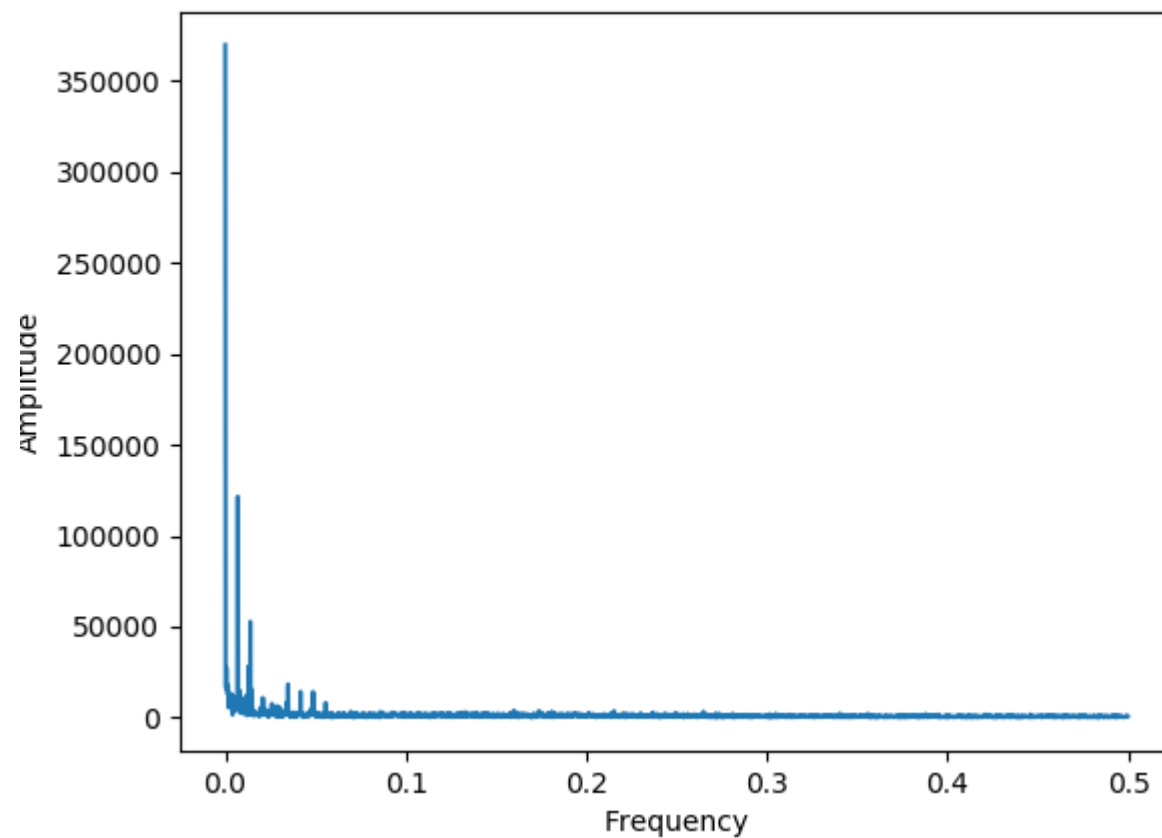
# Time series and Fourier Transforms

```
In [ ]:   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 data')
         10      plt.plot(second_x,regions_cum[i][4464:8640], color=uniqueish_color(), label='2016 feb month data')
         11      plt.plot(third_x,regions_cum[i][8640:], color=uniqueish_color(), label='2016 march month data')
         12      plt.legend()
         13      plt.show()
```

```
In [69]:   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/numpy.fft.fft.html
           3  Y     = np.fft.fft(np.array(jan_2016_smooth)[0:4460])
           4  # read more about the fftfreq: https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fftfreq.html
           5  freq = np.fft.fftfreq(4460, 1)
           6  n = len(freq)
           7  plt.figure()
           8  plt.plot( freq[:int(n/2)], np.abs(Y)[:int(n/2)] )
           9  plt.xlabel("Frequency")
          10  plt.ylabel("Amplitude")
          11  plt.show()
```

**Figure 10**



Zoom to rectangle

```
In [70]:  1  #Preparing the Dataframe only with x(i) values as jan-2015 data and y(i) values as jan-2016
          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
```

```
In [71]:    1  Y[10]
```

Out[71]:   (8717.640047632445+10569.201257952609j)

# 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 - $R_t = (R_{t-1} + R_{t-2} + R_{t-3}\dots R_{t-n})/n$

In [72]:

```python
def MA_R_Predictions(ratios,month):
    predicted_ratio=(ratios['Ratios'].values)[0]
    error=[]
    predicted_values=[]
    window_size=3
    predicted_ratio_values=[]
    for i in range(0,4464*40):
        if i%4464==0:
            predicted_ratio_values.append(0)
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_ratio_values.append(predicted_ratio)
        predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
        error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(ratios['Prediction']
        if i+1>=window_size:
            predicted_ratio=sum((ratios['Ratios'].values)[(i+1)-window_size:(i+1)])/window_size
        else:
            predicted_ratio=sum((ratios['Ratios'].values)[0:(i+1)])/(i+1)


    ratios['MA_R_Predicted'] = predicted_values
    ratios['MA_R_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
    mse_err = sum([e**2 for e in error])/len(error)
    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 $P_t = (P_{t-1} + P_{t-2} + P_{t-3} \ldots P_{t-n})/n$

```
In [73]:
```

```python
def MA_P_Predictions(ratios,month):
    predicted_value=(ratios['Prediction'].values)[0]
    error=[]
    predicted_values=[]
    window_size=1
    predicted_ratio_values=[]
    for i in range(0,4464*40):
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
        if i+1>=window_size:
            predicted_value=int(sum((ratios['Prediction'].values)[(i+1)-window_size:(i+1)])/window_size)
        else:
            predicted_value=int(sum((ratios['Prediction'].values)[0:(i+1)])/(i+1))

    ratios['MA_P_Predicted'] = predicted_values
    ratios['MA_P_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
    mse_err = sum([e**2 for e in error])/len(error)
    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 - $R_t = (N * R_{t-1} + (N-1) * R_{t-2} + (N-2) * R_{t-3}\dots. 1 * R_{t-n})/(N * (N+1)/2)$

```python
In [60]: 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)-(ratios['Prediction'].v
        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['Prediction'].values))
        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 - $P_t = (N * P_{t-1} + (N-1) * P_{t-2} + (N-2) * P_{t-3} \ldots 1 * P_{t-n})/(N * (N+1)/2)$

```
In [61]:    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['Prediction'].values))
           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 (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 weight of 0.18 which keeps exponentially decreasing for the subsequent values.

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

```
In [62]:   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)-(ratios['Prediction']
          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['Prediction'].values))
          21      mse_err = sum([e**2 for e in error])/len(error)
          22      return ratios,mape_err,mse_err
```

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

```
In [63]:   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'].values)[i]))
          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['Prediction'].values))
          18      mse_err = sum([e**2 for e in error])/len(error)
          19      return ratios,mape_err,mse_err
```

```
In [74]:   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

```
In [75]:  print ("Error Metric Matrix (Forecasting Methods) - MAPE & MSE")
          print ("----------------------------------------------------------
          print ("Moving Averages (Ratios) -                  MAPE: ",mean_err[0]," MSE: ",median_err[0])
          print ("Moving Averages (2016 Values) -             MAPE: ",mean_err[1]," MSE: ",median_err[1]
          print ("----------------------------------------------------------
          print ("Weighted Moving Averages (Ratios) -         MAPE: ",mean_err[2]," MSE: ",median_err[2])
          print ("Weighted Moving Averages (2016 Values) -    MAPE: ",mean_err[3]," MSE: ",median_err[3])
          print ("----------------------------------------------------------
          print ("Exponential Moving Averages (Ratios) -      MAPE: ",mean_err[4]," MSE: ",median_err[4])
          print ("Exponential Moving Averages (2016 Values) - MAPE: ",mean_err[5]," MSE: ",median_err[5])
```

```
Error Metric Matrix (Forecasting Methods) - MAPE & MSE
-----------------------------------------------------------------------------
Moving Averages (Ratios) -                  MAPE:  0.1821155173392136      MSE:  400.062550403225
8
Moving Averages (2016 Values) -             MAPE:  0.14292849686975506      MSE:  174.8490199372
7598
-----------------------------------------------------------------------------
Weighted Moving Averages (Ratios) -         MAPE:  0.1784869254376018      MSE:  384.015787410394
24
Weighted Moving Averages (2016 Values) -    MAPE:  0.13551088436182082      MSE:  162.46707549283
155
-----------------------------------------------------------------------------
Exponential Moving Averages (Ratios) -      MAPE:  0.17783550194861494     MSE:  378.34610215053766
Exponential Moving Averages (2016 Values) - MAPE:  0.1350915263669572      MSE:  159.73614471326164
```

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

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

# Fourier Features :

```
In [76]:   1  # For each cluster from 0 to 39 i.e total clusters
           2  # Fourier features dataframe - Stores fourier features for all clusters.
           3  fourier_features = pd.DataFrame(['A1', 'A2', 'A3', 'A4', 'A5', 'F1', 'F2', 'F3', 'F4', 'F5'])
           4  ans = []
           5  for i in range(0,40):
           6
           7      # for each month calculate fft and get frequency
           8      # regions cum hold data for each cluster in format jan,feb,mar. first 4464 values are for jan, next 4176
           9      janfft_data = regions_cum[i][0:4464]
          10      febfft_data = regions_cum[i][4464:4464+4176]
          11      marfft_data = regions_cum[i][4464+4176: 4464+4176+4464]
          12
          13      # calculate fft i.e Amplitude ......
          14      janfft_amp = np.fft.fft(janfft_data)
          15      janfft_freq = np.fft.fftfreq(4464, 1)
          16
          17      febfft_amp = np.fft.fft(febfft_data)
          18      febfft_freq = np.fft.fftfreq(4176, 1)
          19
          20      marfft_amp = np.fft.fft(marfft_data)
          21      marfft_freq = np.fft.fftfreq(4464, 1)
          22
          23      # Sort the amps and frequency and take only top 5 values..
          24      janfft_amp = sorted(janfft_amp, reverse = True)[:5]
          25      janfft_freq = sorted(janfft_freq, reverse = True)[:5]
          26
          27      febfft_amp = sorted(febfft_amp, reverse = True)[:5]
          28      febfft_freq = sorted(febfft_freq, reverse = True)[:5]
          29
          30      marfft_amp = sorted(marfft_amp, reverse = True)[:5]
          31      marfft_freq = sorted(marfft_freq, reverse = True)[:5]
          32
          33      # Each Cluster contains 4464 values of jan , 4176 values of feb, 4464 values of march.
          34      # For eahc value of a month F1, A1 do not change sowe replicate these f1, a1 values as follows;
          35      x = janfft_amp
          36      y = febfft_amp
          37      z = marfft_amp
          38      u = janfft_freq
          39      v = febfft_freq
          40      w = marfft_freq
          41      for f in range(5):
```

```python
42          janfft_amp[f] = [x[f]] * 4464
43          febfft_amp[f] = [y[f]] * 4176
44          marfft_amp[f] = [z[f]] * 4464
45
46          janfft_freq[f] = [u[f]] * 4464
47          febfft_freq[f] = [v[f]] * 4176
48          marfft_freq[f] = [w[f]] * 4464
49
50      # Converting to numpy array and Transpose to get right dimension.
51      janfft_amp = np.array(janfft_amp).T
52      febfft_amp = np.array(febfft_amp).T
53      marfft_amp = np.array(marfft_amp).T
54
55      janfft_freq = np.array(janfft_freq).T
56      febfft_freq = np.array(febfft_freq).T
57      marfft_freq = np.array(marfft_freq).T
58
59
60      # Joining amplitude and frequency of same month and combining different months together.
61      jan_clus = np.hstack((janfft_amp, janfft_freq))
62      feb_clus = np.hstack((febfft_amp, febfft_freq))
63      mar_clus = np.hstack((marfft_amp, marfft_freq))
64
65      clus = np.vstack((jan_clus, feb_clus))
66      clus = np.vstack((clus, mar_clus))
67
68      #Cluster Frame stores the features for a single cluster
69      cluster_features = pd.DataFrame(clus, columns=['A1', 'A2', 'A3', 'A4', 'A5', 'F1', 'F2', 'F3', 'F4', 'F5
70      cluster_features = cluster_features.astype(np.float)
71      ans.append(cluster_features)
72
73
74  # Combining 40 dataframes of fourier features belonging to each cluster into one dataframe
75  print(len(ans))
76  print(type(ans[0]))
77  fourier_features = ans[0]
78  for i in range(1, len(ans)):
79      fourier_features = pd.concat([fourier_features, ans[i]], ignore_index=True)
80  fourier_features = fourier_features.fillna(0)
81  print("Shape of fourier transformed features for all points - ", fourier_features.shape)
82  fourier_features = fourier_features.astype(np.float)
83  fourier_features.tail(3)
```

```
40
<class 'pandas.core.frame.DataFrame'>
Shape of fourier transformed features for all points -  (524160, 10)
```

Out[76]:

| | A1 | A2 | A3 | A4 | A5 | F1 | F2 | F3 | F4 | F5 |
|---|---|---|---|---|---|---|---|---|---|---|
| **524157** | 315146.0 | 11112.786226 | 11112.786226 | 6932.193758 | 6932.193758 | 0.499776 | 0.499552 | 0.499328 | 0.499104 | 0.49888 |

| | A1 | A2 | A3 | A4 | A5 | F1 | F2 | F3 | F4 | F5 |
|---|---|---|---|---|---|---|---|---|---|---|
| **524158** | 315146.0 | 11112.786226 | 11112.786226 | 6932.193758 | 6932.193758 | 0.499776 | 0.499552 | 0.499328 | 0.499104 | 0.49888 |
| **524159** | 315146.0 | 11112.786226 | 11112.786226 | 6932.193758 | 6932.193758 | 0.499776 | 0.499552 | 0.499328 | 0.499104 | 0.49888 |

```
In [77]:   1  # Preparing data to be split into train and test, The below prepares data in cumulative form which will be
           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 which represents the
           7  # that are happened for three months in 2016 data
           8
           9  # print(len(regions_cum))
          10  # 40
          11  # print(len(regions_cum[0]))
          12  # 12960
          13
          14  # we take number of pickups that are happened in last 5 10min intravels
          15  number_of_time_stamps = 5
          16
          17  # output varaible
          18  # it is list of lists
          19  # it will contain number of pickups 13099 for each cluster
          20  output = []
          21
          22
          23  # tsne_lat will contain 13104-5=13099 times lattitude of cluster center for every cluster
          24  # Ex: [[cent_lat 13099times],[cent_lat 13099times], [cent_lat 13099times].... 40 lists]
          25  # it is list of lists
          26  tsne_lat = []
          27
          28
          29  # tsne_lon will contain 13104-5=13099 times logitude of cluster center for every cluster
          30  # Ex: [[cent_long 13099times],[cent_long 13099times], [cent_long 13099times].... 40 lists]
          31  # it is list of lists
          32  tsne_lon = []
          33
          34  # we will code each day
          35  # sunday = 0, monday=1, tue = 2, wed=3, thur=4, fri=5,sat=6
          36  # for every cluster we will be adding 13099 values, each value represent to which day of the week that pick
          37  # it is list of lists
          38  tsne_weekday = []
          39
          40  # its an numbpy array, of shape (523960, 5)
          41  # each row corresponds to an entry in out data
```

```
42   # for the first row we will have [f0,f1,f2,f3,f4] fi=number of pickups happened in i+1th 10min intravel(bin
43   # the second row will have [f1,f2,f3,f4,f5]
44   # the third row will have [f2,f3,f4,f5,f6]
45   # and so on...
46   tsne_feature = []
47
48
49   tsne_feature = [0]*number_of_time_stamps
50   for i in range(0,40):
51       tsne_lat.append([kmeans.cluster_centers_[i][0]]*13099)
52       tsne_lon.append([kmeans.cluster_centers_[i][1]]*13099)
53       # jan 1st 2016 is thursday, so we start our day from 4: "(int(k/144))%7+4"
54       # our prediction start from 5th 10min intravel since we need to have number of pickups that are happened
55       tsne_weekday.append([int(((int(k/144))%7+4)%7) for k in range(5,4464+4176+4464)])
56       # regions_cum is a list of lists [[x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3
57       tsne_feature = np.vstack((tsne_feature, [regions_cum[i][r:r+number_of_time_stamps] for r in range(0,len
58       output.append(regions_cum[i][5:])
59   tsne_feature = tsne_feature[1:]
```
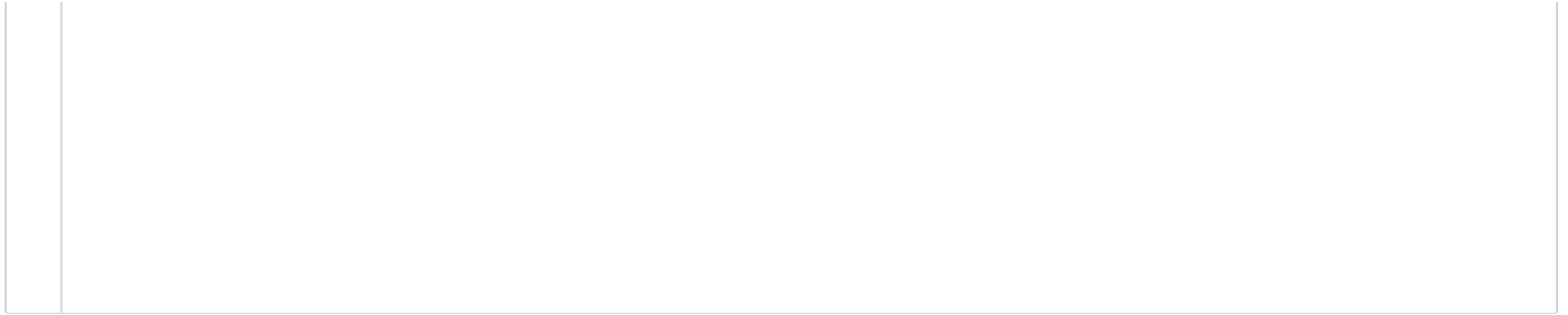
In [78]: `at) == tsne_feature.shape[0] == len(tsne_weekday)*len(tsne_weekday[0]) == 40*13099 == len(output)*len(output[0])`

Out[78]: True

```
In [79]:    1  # Getting the predictions of exponential moving averages to be used as a feature in cumulative form
            2
            3  # upto now we computed 8 features for every data point that starts from 50th min of the day
            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 best error
           14  # we will try to add the same exponential weighted moving avarage at t as a feature to our data
           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 intravel,
           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..x13104], [x5,x6,x7..x
           26  predict_list = []
           27  tsne_flat_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=[]
```

**Holts Winter Triple exponential smoothing :**

References - https://grisha.org/blog/2016/02/17/triple-exponential-smoothing-forecasting-part-iii/ (https://grisha.org/blog/2016/02/17/triple-exponential-smoothing-forecasting-part-iii/)

```python
def initial_trend(series, slen):
    sum = 0.0
    for i in range(slen):
        sum += float(series[i+slen] - series[i]) / slen
    return sum / slen

def initial_seasonal_components(series, slen):
    seasonals = {}
    season_averages = []
    n_seasons = int(len(series)/slen)
    # compute season averages
    for j in range(n_seasons):
        season_averages.append(sum(series[slen*j:slen*j+slen])/float(slen))
    # compute initial values
    for i in range(slen):
        sum_of_vals_over_avg = 0.0
        for j in range(n_seasons):
            sum_of_vals_over_avg += series[slen*j+i]-season_averages[j]
        seasonals[i] = sum_of_vals_over_avg/n_seasons
    return seasonals
```

```python
def triple_exponential_smoothing(series, slen, alpha, beta, gamma, n_preds):
    result = []
    seasonals = initial_seasonal_components(series, slen)
    for i in range(len(series)+n_preds):
        if i == 0: # initial values
            smooth = series[0]
            trend = initial_trend(series, slen)
            result.append(series[0])
            continue
        if i >= len(series): # we are forecasting
            m = i - len(series) + 1
            result.append((smooth + m*trend) + seasonals[i%slen])
        else:
            val = series[i]
            last_smooth, smooth = smooth, alpha*(val-seasonals[i%slen]) + (1-alpha)*(smooth+trend)
            trend = beta * (smooth-last_smooth) + (1-beta)*trend
            seasonals[i%slen] = gamma*(val-smooth) + (1-gamma)*seasonals[i%slen]
            result.append(smooth+trend+seasonals[i%slen])
    return result
```

```
In [82]:    1  alpha = 0.2
            2  beta = 0.15
            3  gamma = 0.2
            4  season_len = 24
            5
            6  predict_values_2 =[]
            7  predict_list_2 = []
            8  tsne_flat_exp_avg_2 = []
            9  for r in range(0,40):
           10      predict_values_2 = triple_exponential_smoothing(regions_cum[r][0:13104], season_len, alpha, beta, gamma
           11      predict_list_2.append(predict_values_2[5:])
```

```
In [ ]:     1
```

```
In [ ]:     1
```

```
In [ ]:     1
```

```
In [ ]:     1  # train, test split : 70% 30% split
            2  # Before we start predictions using the tree based regression models we take 3 months of 2016 pickup data
            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))
```

```python
In [83]:    1  # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our training data
            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)]
            5
            6  # Extracting the same for fourier features -->
            7
            8  fourier_features_train = pd.DataFrame(columns=['A1', 'A2', 'A3', 'A4', 'A5', 'F1', 'F2', 'F3', 'F4', 'F5'])
            9  fourier_features_test = pd.DataFrame(columns=['A1', 'A2', 'A3', 'A4', 'A5', 'F1', 'F2', 'F3', 'F4', 'F5'])
           10
           11  for i in range(40):
           12      fourier_features_train = fourier_features_train.append(fourier_features[i*13099 : 13099*i + 9169])
           13
           14  fourier_features_train.reset_index(inplace = True)
           15
           16
           17  for i in range(40):
           18      fourier_features_test = fourier_features_test.append(fourier_features[i*13099 + 9169 : 13099*(i+1)])
           19
           20  fourier_features_test.reset_index(inplace = True)
```

```python
In [84]:    1  print("Number of data clusters",len(train_features), "Number of data points in trian data", len(train_featu
            2  print("Number of data clusters",len(train_features), "Number of data points in test data", len(test_feature
```

```
Number of data clusters 40 Number of data points in trian data 9169 Each data point contains 5 features
Number of data clusters 40 Number of data points in test data 3930 Each data point contains 5 features
```

```
In [92]:  1  # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our training data
          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
          8  tsne_train_flat_triple_avg = [i[:9169] for i in predict_list_2]
```

```
In [93]:  1  # extracting the rest of the timestamp values i.e 30% of 12956 (total timestamps) for our test data
          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
          8  tsne_test_flat_triple_avg = [i[9169:] for i in predict_list_2]
```

```
In [94]:  1  # the above contains values in the form of list of lists (i.e. list of values of each region), here we make
          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])
```

```
In [95]:    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
           12  tsne_train_triple_avg = sum(tsne_train_flat_triple_avg,[])
```

```
In [96]:    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
           12  tsne_test_triple_avg = sum(tsne_test_flat_triple_avg,[])
```

In [97]:
```python
 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
 9  df_train['3EXP'] = tsne_train_triple_avg
10
11  print(df_train.shape)
```

(366760, 10)

In [98]:
```python
 1  # Preparing the data frame for our train data
 2  df_test = pd.DataFrame(data=test_new_features, columns=columns)
 3  df_test['lat'] = tsne_test_lat
 4  df_test['lon'] = tsne_test_lon
 5  df_test['weekday'] = tsne_test_weekday
 6  df_test['exp_avg'] = tsne_test_exp_avg
 7
 8  df_test['3EXP'] = tsne_test_triple_avg
 9
10  print(df_test.shape)
11  print(df_test.shape)
```

(157200, 10)
(157200, 10)

In [99]:
```
1 df_test.head()
```

Out[99]:

| | ft_5 | ft_4 | ft_3 | ft_2 | ft_1 | lat | lon | weekday | exp_avg | 3EXP |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 118 | 106 | 104 | 93 | 102 | 40.776228 | -73.982119 | 4 | 100 | 97.296682 |
| 1 | 106 | 104 | 93 | 102 | 101 | 40.776228 | -73.982119 | 4 | 100 | 105.445923 |
| 2 | 104 | 93 | 102 | 101 | 120 | 40.776228 | -73.982119 | 4 | 114 | 115.044145 |
| 3 | 93 | 102 | 101 | 120 | 131 | 40.776228 | -73.982119 | 4 | 125 | 132.975561 |
| 4 | 102 | 101 | 120 | 131 | 164 | 40.776228 | -73.982119 | 4 | 152 | 142.108910 |

## Merging the fourier features :

In [100]:
```
1 df_train_2 = df_train
2 df_test_2 = df_test
3 df_train = pd.concat([df_train, fourier_features_train], axis = 1)
4 df_test = pd.concat([df_test, fourier_features_test], axis = 1)
```

In [101]:
```
1 print("Shape of Train Data Now - ", df_train.shape)
2 df_train.drop(['index'], axis = 1, inplace=True)
3 df_train.head()
```

Shape of Train Data Now -  (366760, 21)

Out[101]:

| | ft_5 | ft_4 | ft_3 | ft_2 | ft_1 | lat | lon | weekday | exp_avg | 3EXP | A1 | A2 | A3 | A4 | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 63 | 217 | 189 | 137 | 40.776228 | -73.982119 | 4 | 150 | 126.474978 | 369774.0 | 24998.122651 | 24998.122651 | 15434.851794 | 15434.85179 |
| 1 | 63 | 217 | 189 | 137 | 135 | 40.776228 | -73.982119 | 4 | 139 | 136.988688 | 369774.0 | 24998.122651 | 24998.122651 | 15434.851794 | 15434.85179 |
| 2 | 217 | 189 | 137 | 135 | 129 | 40.776228 | -73.982119 | 4 | 132 | 153.426260 | 369774.0 | 24998.122651 | 24998.122651 | 15434.851794 | 15434.85179 |
| 3 | 189 | 137 | 135 | 129 | 150 | 40.776228 | -73.982119 | 4 | 144 | 168.323089 | 369774.0 | 24998.122651 | 24998.122651 | 15434.851794 | 15434.85179 |
| 4 | 137 | 135 | 129 | 150 | 164 | 40.776228 | -73.982119 | 4 | 158 | 175.333204 | 369774.0 | 24998.122651 | 24998.122651 | 15434.851794 | 15434.85179 |

In [102]:
```python
1  print("Shape of Test Data Now - ", df_test.shape)
2  df_test.drop(['index'], axis = 1, inplace=True)
3  df_test.head()
```

Shape of Test Data Now -  (157200, 21)

Out[102]:

|   | ft_5 | ft_4 | ft_3 | ft_2 | ft_1 | lat | lon | weekday | exp_avg | 3EXP | A1 | A2 | A3 | A4 | A |
|---|------|------|------|------|------|-----|-----|---------|---------|------|----|----|----|----|---|
| 0 | 118 | 106 | 104 | 93 | 102 | 40.776228 | -73.982119 | 4 | 100 | 97.296682 | 391598.0 | 10930.478599 | 10930.478599 | 10662.395979 | 10662.39597 |
| 1 | 106 | 104 | 93 | 102 | 101 | 40.776228 | -73.982119 | 4 | 100 | 105.445923 | 391598.0 | 10930.478599 | 10930.478599 | 10662.395979 | 10662.39597 |
| 2 | 104 | 93 | 102 | 101 | 120 | 40.776228 | -73.982119 | 4 | 114 | 115.044145 | 391598.0 | 10930.478599 | 10930.478599 | 10662.395979 | 10662.39597 |
| 3 | 93 | 102 | 101 | 120 | 131 | 40.776228 | -73.982119 | 4 | 125 | 132.975561 | 391598.0 | 10930.478599 | 10930.478599 | 10662.395979 | 10662.39597 |
| 4 | 102 | 101 | 120 | 131 | 164 | 40.776228 | -73.982119 | 4 | 152 | 142.108910 | 391598.0 | 10930.478599 | 10930.478599 | 10662.395979 | 10662.39597 |

In [ ]:
```python
1
```

In [ ]:
```python
1
```

## Using Linear Regression

```
In [90]:   1  # find more about LinearRegression function here http://scikit-learn.org/stable/modules/generated/sklearn.l
           2  # ------------------------
           3  # default paramters
           4  # sklearn.linear_model.LinearRegression(fit_intercept=True, normalize=False, copy_X=True, n_jobs=1)
           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 prediction.
          11  # set_params(**params)  Set the parameters of this estimator.
          12  # ----------------------
          13  # video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-intuition-1
          14  # ----------------------
          15
          16  from sklearn.linear_model import LinearRegression
          17  from sklearn.model_selection import GridSearchCV
          18
          19  intercept=[True,False]
          20  normalize=[True, False]
          21  copyX = [True, False]
          22  param_grid = dict(fit_intercept=intercept,normalize=normalize, copy_X = copyX)
          23
          24  lr_reg=LinearRegression()
          25  grid = GridSearchCV(estimator=lr_reg, param_grid=param_grid, cv = 2, n_jobs=-1)
          26  grid_result = grid.fit(df_train, tsne_train_output)
          27
          28  print("Best: %f using %s" % (grid_result.best_score_, grid_result.best_params_))
          29
          30  #
          31
          32
```

Best: 0.957718 using {'copy_X': True, 'fit_intercept': True, 'normalize': True}

In [92]:
```python
lr_reg = LinearRegression(normalize=True, n_jobs=-1).fit(df_train, tsne_train_output)
y_pred = lr_reg.predict(df_test)
lr_test_predictions = [round(value) for value in y_pred]
y_pred = lr_reg.predict(df_train)
lr_train_predictions = [round(value) for value in y_pred]

```

## Using Random Forest Regressor

```
In [94]:   1  # Training a hyper-parameter tuned random forest regressor on our train data
           2  # find more about LinearRegression function here http://scikit-learn.org/stable/modules/generated/sklearn.e
           3  # -------------------------
           4  # default paramters
           5  # sklearn.ensemble.RandomForestRegressor(n_estimators=10, criterion='mse', max_depth=None, min_samples_spli
           6  # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_
           7  # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_st
           8
           9  # some of methods of RandomForestRegressor()
          10  # apply(X)  Apply trees in the forest to X, return leaf indices.
          11  # decision_path(X)  Return the decision path in the forest
          12  # fit(X, y[, sample_weight])    Build a forest of trees from the training set (X, y).
          13  # get_params([deep])    Get parameters for this estimator.
          14  # predict(X)    Predict regression target for X.
          15  # score(X, y[, sample_weight])  Returns the coefficient of determination R^2 of the prediction.
          16  # -----------------------
          17  # video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-using-dec
          18  # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
          19  # -----------------------
          20  from sklearn.model_selection import RandomizedSearchCV
          21  estimator=[10,15,25,50,100]
          22  max_depth = [5,10,15,20,25]
          23  min_split = [2,5,10]
          24
          25  param_grid = dict(n_estimators=estimator,max_depth=max_depth, min_samples_split = min_split)
          26  rf = RandomForestRegressor()
          27  start_time = time.time()
          28  random_result = RandomizedSearchCV(estimator=rf, param_distributions=param_grid, cv = 2, n_jobs=-1)
          29  random_result.fit(df_train, tsne_train_output)
          30
          31  print("Best: %f using %s" % (random_result.best_score_, random_result.best_params_))
          32  print("Execution time: " + str((time.time() - start_time)) + ' ms')
          33  # regr1 = RandomForestRegressor(max_features='sqrt',min_samples_leaf=4,min_samples_split=3,n_estimators=40,
          34
          35  # regr1.fit(df_train, tsne_train_output)
```

```
Best: 0.957092 using {'n_estimators': 25, 'min_samples_split': 5, 'max_depth': 10}
Execution time: 399.7216980457306 ms
```

In [95]:
```
 1  ng on test data using our trained random forest model
 2
 3  ls regr1 is already hyper parameter tuned
 4  neters that we got above are found using grid search
    ndomForestRegressor(n_estimators = 25, max_depth = 10, min_samples_split= 5, n_jobs= -1).fit(df_train, tsne_trair
    egr1.predict(df_test)
    predictions = [round(value) for value in y_pred]
    egr1.predict(df_train)
    _predictions = [round(value) for value in y_pred]
```

In [96]:
```
1  #feature importances based on analysis using random forest
2  print (df_train.columns)
3  print (regr1.feature_importances_)
```

```
Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'lat', 'lon', 'weekday',
       'exp_avg'],
      dtype='object')
[8.28441103e-04 4.50865451e-04 5.86220341e-04 6.40066992e-04
 4.37792340e-04 2.10477769e-04 2.70639367e-04 8.07960836e-05
 9.96494701e-01]
```

## Using XgBoost Regressor

In [105]:

```python
# Training a hyper-parameter tuned Xg-Boost regressor on our train data

# find more about XGBRegressor function here http://xgboost.readthedocs.io/en/latest/python/python_api.html
# -------------------------
# default paramters
# xgboost.XGBRegressor(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True, objective='reg:linear
# booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1, max_delta_step=0, subsample=1, col
# colsample_bylevel=1, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, base_score=0.5, random_state=0, seed=
# missing=None, **kwargs)

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True,
# get_params([deep])    Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread
# get_score(importance_type='weight') -> get the feature importance
# ----------------------
# video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-using-dec
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
# ----------------------

params = {"learning_rate"    : [0.05, 0.10, 0.15, 0.20, 0.25, 0.30 ] ,
 "max_depth"        : [ 3, 4, 5, 6, 8, 10, 12, 15],
 "min_child_weight" : [ 1, 3, 5, 7 ],
 "gamma"            : [ 0.0, 0.1, 0.2 , 0.3, 0.4 ],
 "colsample_bytree" : [ 0.3, 0.4, 0.5 , 0.7, 0.8 ] }

xgb_r = xgb.XGBRegressor()
start_time = time.time()
random_result = RandomizedSearchCV(estimator=xgb_r, param_distributions=params, cv = 2, n_jobs=-1)
random_result.fit(df_train, tsne_train_output)

print("Best: %f using %s" % (random_result.best_score_, random_result.best_params_))
print("Execution time: " + str((time.time() - start_time)) + ' ms')
```

```
Best: 0.957433 using {'min_child_weight': 7, 'max_depth': 6, 'learning_rate': 0.1, 'gamma': 0.2, 'colsample_b
ytree': 0.8}
Execution time: 459.62532210350037 ms
```

In [107]:
```python
#predicting with our trained Xg-Boost regressor
# the models x_model is already hyper parameter tuned
# the parameters that we got above are found using grid search
x_model = xgb.XGBRegressor(
 learning_rate =0.1,
 n_estimators=1000,
 max_depth=6,
 min_child_weight=7,
 gamma=0.2,
 subsample=0.8,
 reg_alpha=200, reg_lambda=200,
 colsample_bytree=0.8,nthread=4)
x_model.fit(df_train, tsne_train_output)

y_pred = x_model.predict(df_test)
xgb_test_predictions = [round(value) for value in y_pred]
y_pred = x_model.predict(df_train)
xgb_train_predictions = [round(value) for value in y_pred]
```

```
In [108]:    1  #feature importances
             2  x_model.get_booster().get_score(importance_type='weight')
```

```
Out[108]:  {'exp_avg': 4849,
            'ft_3': 6169,
            'ft_4': 6389,
            'ft_2': 6542,
            'ft_5': 6891,
            'ft_1': 6018,
            'lon': 4382,
            'weekday': 2845,
            'lat': 4712}
```

## Calculating the error metric values for various models

```
In [109]:    1
             2
             3
         end((mean_absolute_error(tsne_train_output,df_train['ft_1'].values))/(sum(tsne_train_output)/len(tsne_train_outp
         end((mean_absolute_error(tsne_train_output,df_train['exp_avg'].values))/(sum(tsne_train_output)/len(tsne_train_d
         end((mean_absolute_error(tsne_train_output,rndf_train_predictions))/(sum(tsne_train_output)/len(tsne_train_outp
         end((mean_absolute_error(tsne_train_output, xgb_train_predictions))/(sum(tsne_train_output)/len(tsne_train_outp
         end((mean_absolute_error(tsne_train_output, lr_train_predictions))/(sum(tsne_train_output)/len(tsne_train_outp
             9
         end((mean_absolute_error(tsne_test_output, df_test['ft_1'].values))/(sum(tsne_test_output)/len(tsne_test_output)
         end((mean_absolute_error(tsne_test_output, df_test['exp_avg'].values))/(sum(tsne_test_output)/len(tsne_test_outp
         end((mean_absolute_error(tsne_test_output, rndf_test_predictions))/(sum(tsne_test_output)/len(tsne_test_output))
         end((mean_absolute_error(tsne_test_output, xgb_test_predictions))/(sum(tsne_test_output)/len(tsne_test_output)))
         end((mean_absolute_error(tsne_test_output, lr_test_predictions))/(sum(tsne_test_output)/len(tsne_test_output)))
```

```
In [115]:   1  print ("Error Metric Matrix (Tree Based Regression Methods) -  MAPE")
            2  print ("-----------------------------------------------------------------------------
            3  print ("Baseline Model -                    Train: ",train_mape[0],"      Test: ",test_mape[0])
            4  print ("Exponential Averages Forecasting -   Train: ",train_mape[1],"      Test: ",test_mape[1])
            5  print ("Linear Regression -                 Train: ",train_mape[4],"     Test: ",test_mape[4])
            6  print ("Random Forest Regression -           Train: ",train_mape[2],"      Test: ",test_mape[2])
            7  print ("XGB Regression -                     Train: ",train_mape[3],"      Test: ",test_mape[3])
```

```
Error Metric Matrix (Tree Based Regression Methods) -  MAPE

-----------------------------------------------------------------------------------------
Baseline Model -                    Train:  0.14005275878666593    Test:  0.13653125704827038
Exponential Averages Forecasting -  Train:  0.13289968436017227    Test:  0.12936180420430524
Linear Regression -                Train:  0.13331572016045437    Test:  0.1291202994009687
Random Forest Regression -          Train:  0.12876988917496496    Test:  0.12769984594153422
XGB Regression -                    Train:   0.12455820870483163    Test:  0.12569947513083501
```

## Error Metric Matrix

```
In [111]:   1  print ("Error Metric Matrix (Tree Based Regression Methods) -  MAPE")
            2  print ("-----------------------------------------------------------------------------
            3  print ("Baseline Model -                       Train: ",train_mape[0]," Test: ",test_mape[0])
            4  print ("Exponential Averages Forecasting -      Train: ",train_mape[1]," Test: ",test_mape[1])
            5  print ("Linear Regression -                     Train: ",train_mape[4]," Test: ",test_mape[4])
            6  print ("Random Forest Regression -              Train: ",train_mape[2]," Test: ",test_mape[2])
            7  print ("XgBoost Regression -                    Train: ",train_mape[3]," Test: ",test_mape[3])
            8  print ("-----------------------------------------------------------------------------
```

```
Error Metric Matrix (Tree Based Regression Methods) -  MAPE
-------------------------------------------------------------------------------
Baseline Model -                      Train:  0.14005275878666593      Test:  0.13653125704827038
Exponential Averages Forecasting -    Train:  0.13289968436017227      Test:  0.12936180420430524
Linear Regression -                   Train:  0.13331572016045437      Test:  0.1291202994009687
Random Forest Regression -            Train:  0.12876988917496496      Test:  0.12769984594153422
XgBoost Regression -                  Train:  0.12455820870483163      Test:  0.12569947513083501
-------------------------------------------------------------------------------
```

# Assignments

```
In [112]:   1  '''
            2  Task 1: Incorporate Fourier features as features into Regression models and measure MAPE. <br>
            3
            4  Task 2: Perform hyper-parameter tuning for Regression models.
            5          2a. Linear Regression: Grid Search
            6          2b. Random Forest: Random Search
            7          2c. Xgboost: Random Search
            8  Task 3: Explore more time-series features using Google search/Quora/Stackoverflow
            9  to reduce the MAPE to < 12%
           10  '''
```

Out[112]: '\nTask 1: Incorporate Fourier features as features into Regression models and measure MAPE. <br>\n\nTask 2: Perform hyper-parameter tuning for Regression models.\n         2a. Linear Regression: Grid Search\n          2b. Random Forest: Random Search \n        2c. Xgboost: Random Search\nTask 3: Explore more time-series features using Google search/Quora/Stackoverflow\nto reduce the MAPE to < 12%\n'

```
In [113]:   1  # TAsk 2 Is done Above.
```

```
In [114]:   1  # Task 1 Incorporating Fourier Features into regression model
```

```
In [ ]:     1
```

```
In [ ]:     1  # USing RF Regressor To see the MAPE
```

```python
In [103]:    1  # Reference : All the grid/random search help is taken from Datacamp
             2  from sklearn.model_selection import RandomizedSearchCV
             3  estimator=[10,15,25,50,100]
             4  max_depth = [5,10,15,20,25]
             5  min_split = [2,5,10]
             6
             7  param_grid = dict(n_estimators=estimator,max_depth=max_depth, min_samples_split = min_split)
             8  rf = RandomForestRegressor()
             9  start_time = time.time()
            10  random_result = RandomizedSearchCV(estimator=rf, param_distributions=param_grid, cv = 2, n_jobs=-1)
            11  random_result.fit(df_train, tsne_train_output)
            12
            13  print("Best: %f using %s" % (random_result.best_score_, random_result.best_params_))
            14  print("Execution time: " + str((time.time() - start_time)) + ' ms')
```

```
Best: 0.974961 using {'n_estimators': 50, 'min_samples_split': 2, 'max_depth': 15}
Execution time: 826.3263411521912 ms
```

```python
In [104]:    1  regr1 = RandomForestRegressor(n_estimators = 25, max_depth = 10, min_samples_split= 5, n_jobs= -1).fit(df_t:
             2  y_pred = regr1.predict(df_test)
             3  rndf_test_predictions = [round(value) for value in y_pred]
             4  y_pred = regr1.predict(df_train)
             5  rndf_train_predictions = [round(value) for value in y_pred]
```

```
In [106]:   1  train_fourier_rf = (mean_absolute_error(tsne_train_output,rndf_train_predictions))/(sum(tsne_train_output)/
            2  test_fourier_rf = (mean_absolute_error(tsne_test_output, rndf_test_predictions))/(sum(tsne_test_output)/len
            3
            4  print ("Random Forest Regression -        Train: ",train_fourier_rf,"     Test: ",test_fourier_rf)
            5
```

```
Random Forest Regression -        Train:  0.09704870085180659    Test:  0.09680185328893799
```

In [ ]:    1

In [ ]:    1

In [ ]:    1

In [ ]:    1

In [ ]:    1