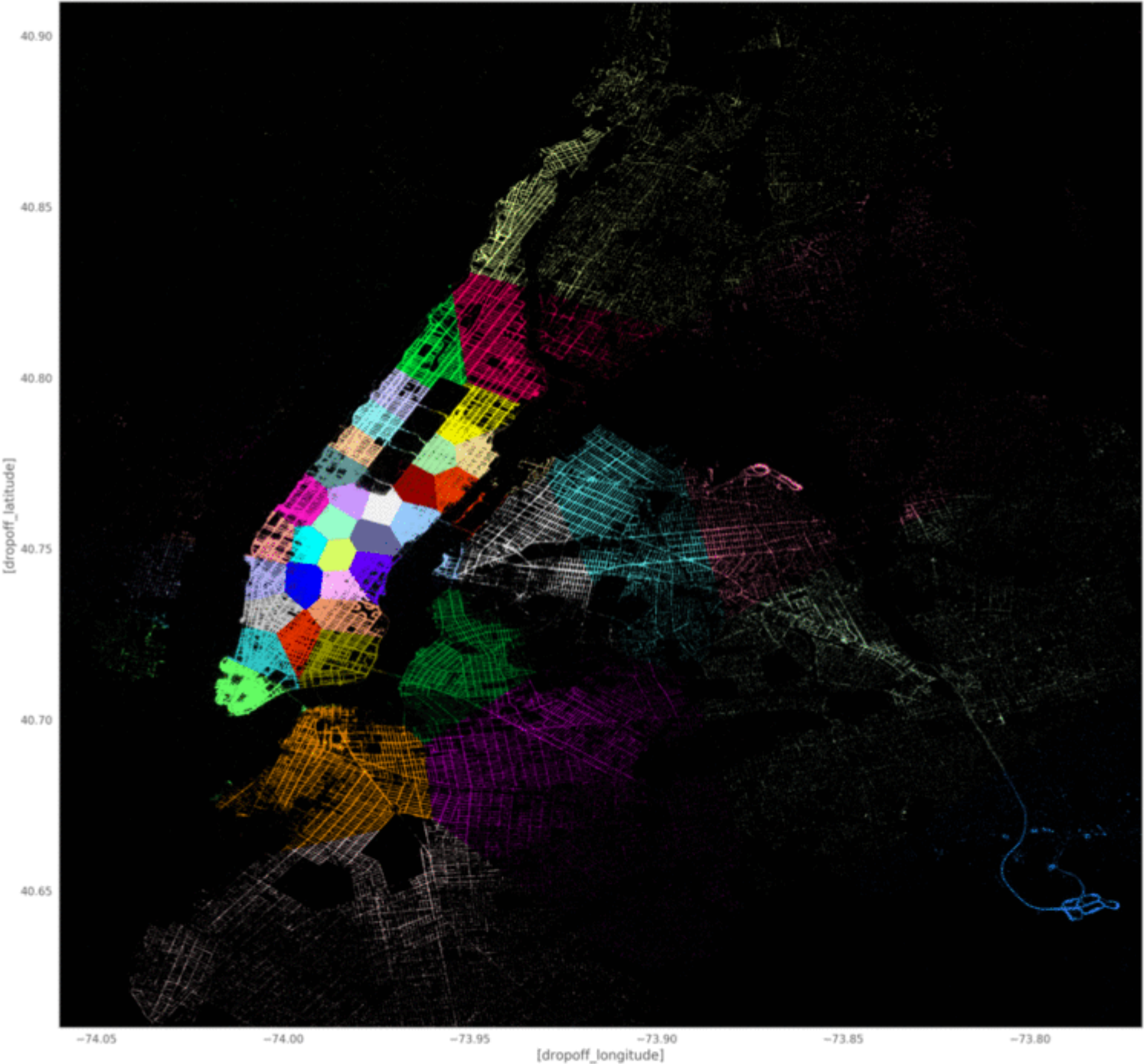


# Taxi demand prediction in New York City



In [1]:

```
1  #Importing Libraries
2  # pip3 install graphviz
3  #pip3 install dask
4  #pip3 install toolz
5  #pip3 install cloudpickle
6  # https://www.youtube.com/watch?v=iEW3G7ZzRZ0
7  # https://github.com/dask/dask-tutorial
8  # please do go through this python notebook: https://github.com/dask/dask-tutorial/blob/master/07_dataframe.ipynb
9  import dask.dataframe as dd#similar to pandas
10
11 import pandas as pd#pandas to create small dataframes
12
13 # pip3 install folium
14 # if this doesnt work refere install_folium.JPG in drive
15 import folium #open street map
16
17 # unix time: https://www.unixtimestamp.com/
18 import datetime #Convert to unix time
19
20 import time #Convert to unix time
21
22 # if numpy is not installed already : pip3 install numpy
23 import numpy as np#Do aritmetic operations on arrays
24
25 # matplotlib: used to plot graphs
26 import matplotlib
27 # matplotlib.use('nbagg') : matplotlib uses this protocall which makes plots more user intractive like zoom in and zoom out
28 matplotlib.use('nbagg')
29 import matplotlib.pyplot as plt
30 import seaborn as sns#Plots
31 from matplotlib import rcParams#Size of plots
32
33 # this lib is used while we calculate the stight line distance between two (Lat,Lon) pairs in miles
34 import gpxpy.geo #Get the haversine distance
35
36 from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
37 import math
38 import pickle
39 import os
40
41 # download mingwin: 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 import graphviz
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](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)

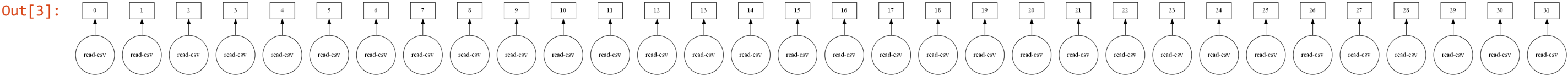
file name	file name size	number of records	number of features
yellow_tripdata_2016-01	1. 59G	10906858	19
yellow_tripdata_2016-02	1. 66G	11382049	19
yellow_tripdata_2016-03	1. 78G	12210952	19
yellow_tripdata_2016-04	1. 74G	11934338	19
yellow_tripdata_2016-05	1. 73G	11836853	19
yellow_tripdata_2016-06	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

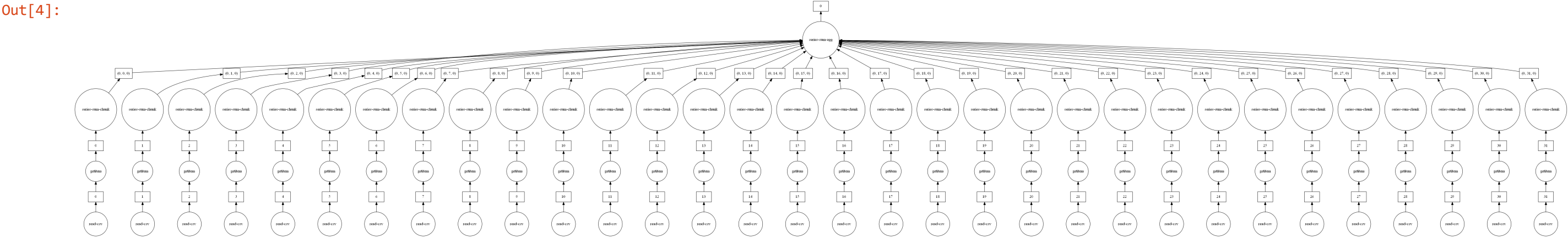
```
In [2]: 1 #Looking at the features
2 # dask dataframe : # https://github.com/dask/dask-tutorial/blob/master/07_dataframe.ipynb
3 month = dd.read_csv('yellow_tripdata_2015-01.csv')
4 print(month.columns)
```

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()
```



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



Features in the dataset:

```

<tr>
  <td>Dropoff_longitude</td>
  <td>Longitude where the meter was disengaged.</td>
</tr>
<tr>
  <td>Dropoff_ latitude</td>
  <td>Latitude where the meter was disengaged.</td>
</tr>
<tr>
  <td>Payment_type</td>
  <td>A numeric code signifying how the passenger paid for the trip.
  <ol>
    <li> Credit card </li>
    <li> Cash </li>
    <li> No charge </li>
    <li> Dispute</li>
    <li> Unknown </li>
    <li> Voided trip</li>
  </ol>
</td>
</tr>
<tr>
  <td>Fare_amount</td>
  <td>The time-and-distance fare calculated by the meter.</td>
</tr>
<tr>
  <td>Extra</td>
  <td>Miscellaneous extras and surcharges. Currently, this only includes. the $0.50 and $1 rush 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 being 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 included.</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	1. 2.	A code indicating the TPEP provider that provided the record. Creative Mobile Technologies 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	1. 2. 3. 4. 5. 6.	The final rate code in effect at the end of the trip. Standard rate JFK Newark Nassau or Westchester Negotiated fare Group ride
Store_and_fwd_flag		This flag indicates whether the trip record was held in vehicle memory before sending to the vendor, aka "store and forward," because the vehicle did not have a connection to the server. Y= store and forward trip N= not a store and forward trip

ML Problem Formulation

Time-series forecasting and Regression

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

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

Performance metrics

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

Data Cleaning

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

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_fwd_flag	dropoff_longitude	dropoff_latitude	payment_type	fare_amou
0	2	2015-01-15 19:05:39	2015-01-15 19:23:42	1	1.59	-73.993896	40.750111	1	N	-73.974785	40.750618	1	12
1	1	2015-01-10 20:33:38	2015-01-10 20:53:28	1	3.30	-74.001648	40.724243	1	N	-73.994415	40.759109	1	14
2	1	2015-01-10 20:33:38	2015-01-10 20:43:41	1	1.80	-73.963341	40.802788	1	N	-73.951820	40.824413	2	9
3	1	2015-01-10 20:33:39	2015-01-10 20:35:31	1	0.50	-74.009087	40.713818	1	N	-74.004326	40.719986	2	3
4	1	2015-01-10 20:33:39	2015-01-10 20:52:58	1	3.00	-73.971176	40.762428	1	N	-74.004181	40.742653	2	15

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.

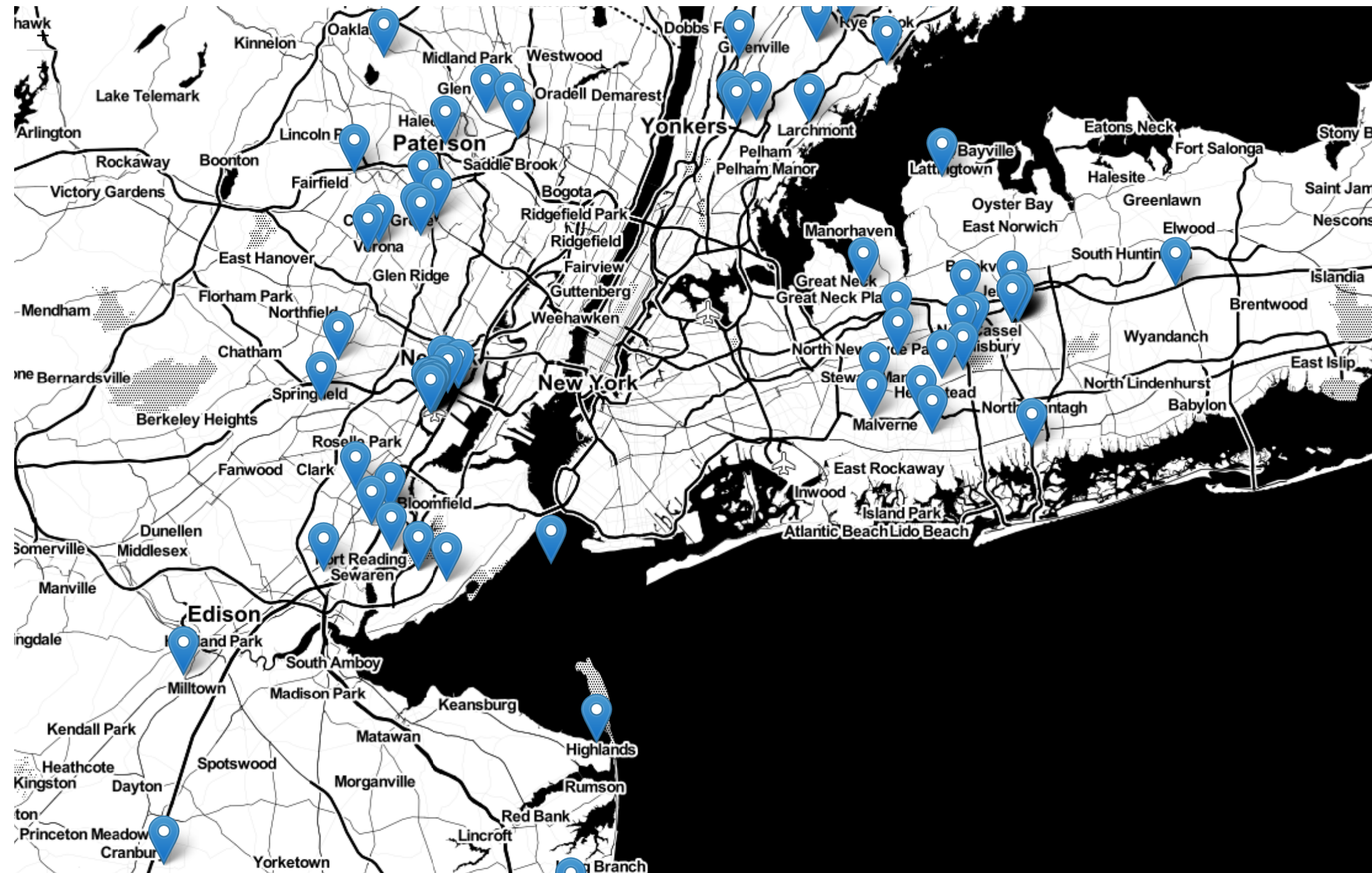


```

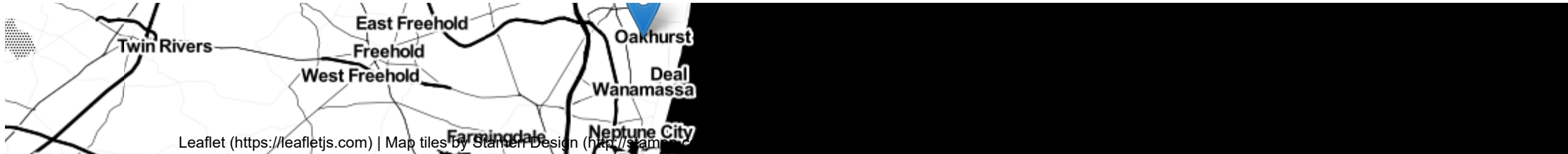
In [292]: 1 # Plotting pickup coordinates which are outside the bounding box of New-York
2 # we will collect all the points outside the bounding box of newyork city to outlier_locations
3 outlier_locations = month[((month.pickup_longitude <= -74.15) | (month.pickup_latitude <= 40.5774) | \
4 (month.pickup_longitude >= -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 indepth 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[292]:







**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

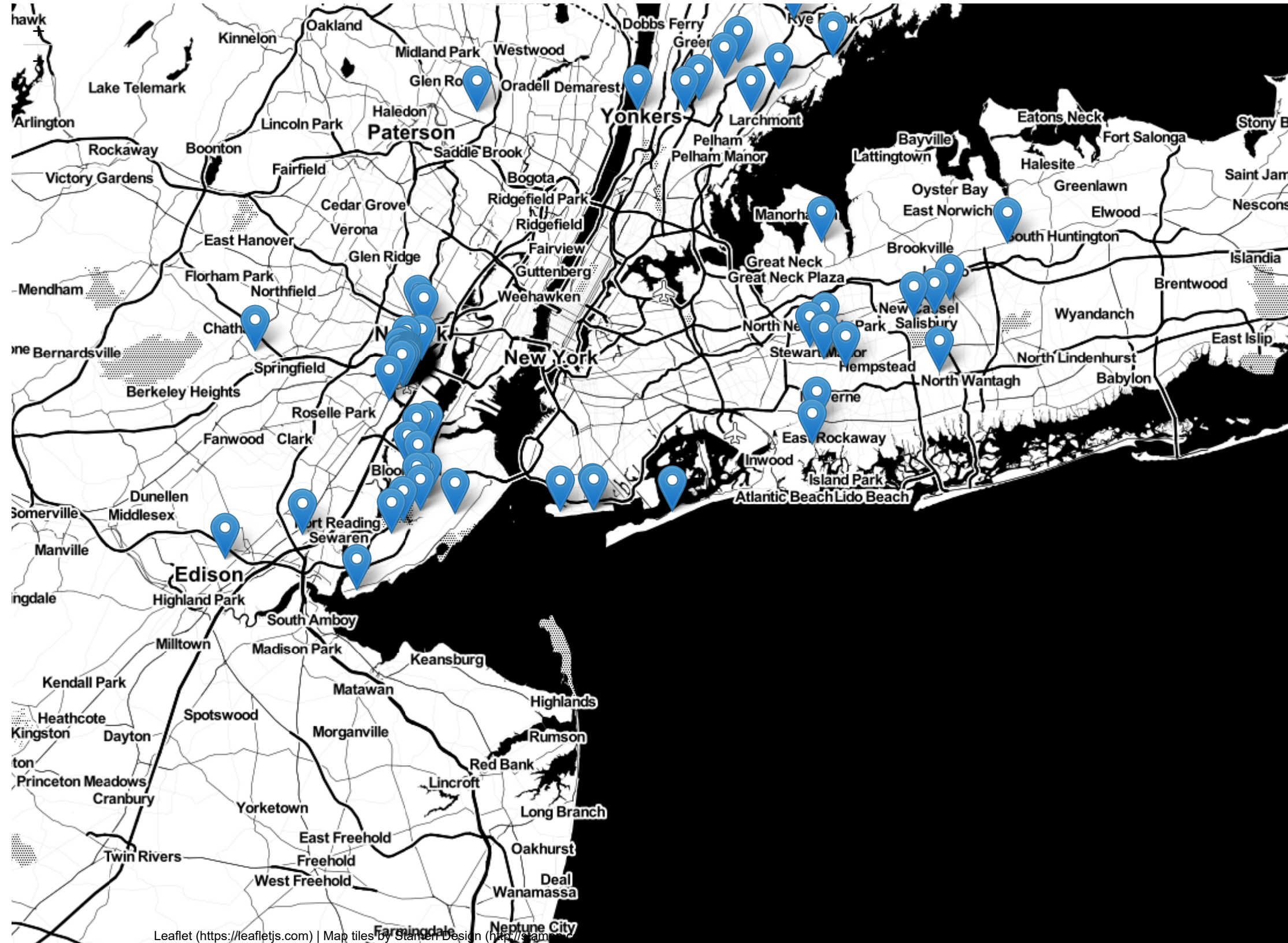
It is inferred from the source <https://www.flickr.com/places/info/2459115> (<https://www.flickr.com/places/info/2459115>) that New York is bounded by the location cordinates(lat,long) - (40.5774, -74.15) & (40.9176,-73.7004) so hence any cordinates not within these cordinates are not considered by us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

```

In [293]: 1 outlier_locations = month[((month.dropoff_longitude <= -74.15) | (month.dropoff_latitude <= 40.5774) | \
2         (month.dropoff_longitude >= -73.7004) | (month.dropoff_latitude >= 40.9176))]
3  ## consider sample outlier points
4  sample_outleir = outlier_locations.head(1000)
5  map_ = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')
6
7  for i,j in zip(sample_outleir['dropoff_latitude'],sample_outleir['dropoff_longitude']):
8      folium.Marker([i,j]).add_to(map_)
9  map_

```

Out[293]:



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

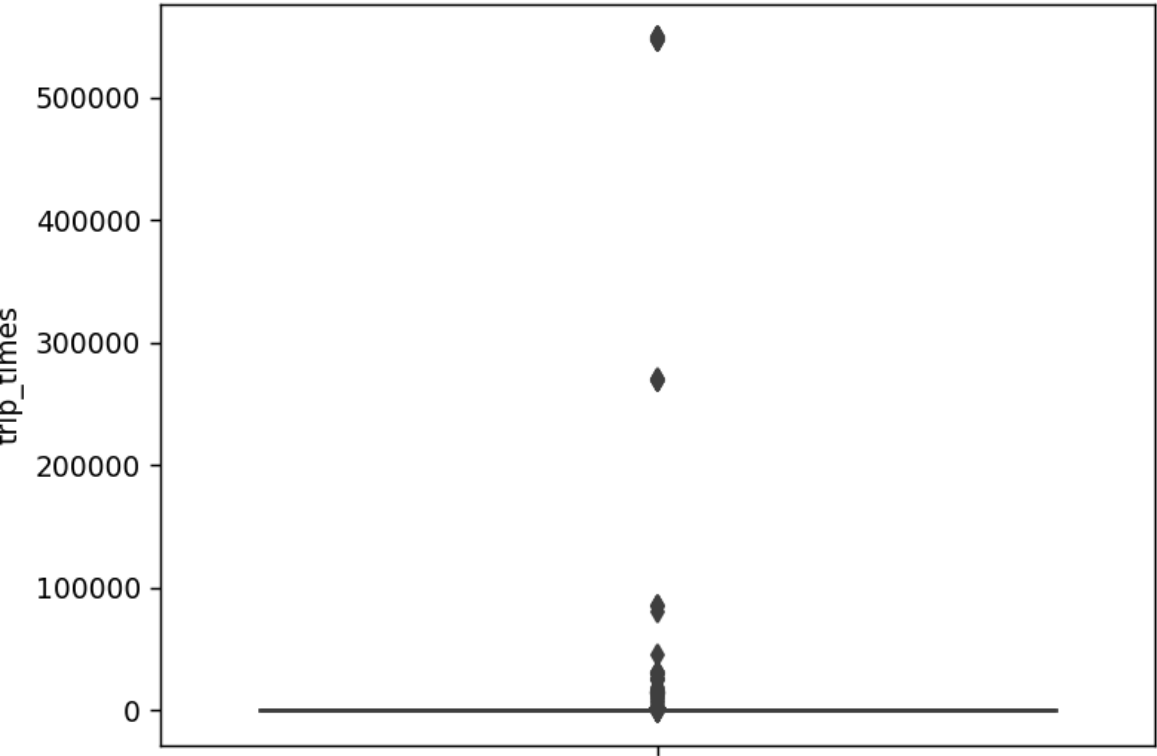
### 3. Trip Durations:

According to NYC Taxi & Limousine Commission 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 are used while binning
2
3 # in out data we have time in the formate "YYYY-MM-DD HH:MM:SS" we convert this sting to python time formate and then into unix time stamp
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 # we return a data frame which contains the columns
9 # 1.'passenger_count' : self explanatory
10 # 2.'trip_distance' : self explanatory
11 # 3.'pickup_longitude' : self explanatory
12 # 4.'pickup_latitude' : self explanatory
13 # 5.'dropoff_longitude' : self explanatory
14 # 6.'dropoff_latitude' : self explanatory
15 # 7.'total_amount' : total fair that was paid
16 # 8.'trip_times' : duration of each trip
17 # 9.'pickup_times' : pickup time converted into unix time
18 # 10.'Speed' : velocity of each trip
19
20 def return_with_trip_times(month):
21     duration = month[['tpep_pickup_datetime', 'tpep_dropoff_datetime']].compute()
22     #pickups and dropoffs to unix time
23     duration_pickup = [convert_to_unix(x) for x in duration['tpep_pickup_datetime'].values]
24     duration_drop = [convert_to_unix(x) for x in duration['tpep_dropoff_datetime'].values]
25     #calculate duration of tripsduration_drop
26     durations = (np.array(duration_drop) - np.array(duration_pickup))/float(60)
27
28     #append durations of trips and speed in miles/hr to a new dataframe
29     new_frame = month[['passenger_count', 'trip_distance', 'pickup_longitude', 'pickup_latitude', 'dropoff_longitude', 'dropoff_latitude', 'total_amount']].compute()
30
31     new_frame['trip_times'] = durations
32     new_frame['pickup_times'] = duration_pickup
33     new_frame['Speed'] = 60*(new_frame['trip_distance']/new_frame['trip_times'])
34
35     return new_frame
36
37 frame_with_durations = return_with_trip_times(month)
```

```
In [9]: 1 # the skewed box plot shows us the presence of outliers
        2 sns.boxplot(y="trip_times", data =frame_with_durations)
        3 plt.show()
```

Figure 1 



x= y=429256

```
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 5.3833333333333334
30 percentile value is 6.8166666666666666
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
```

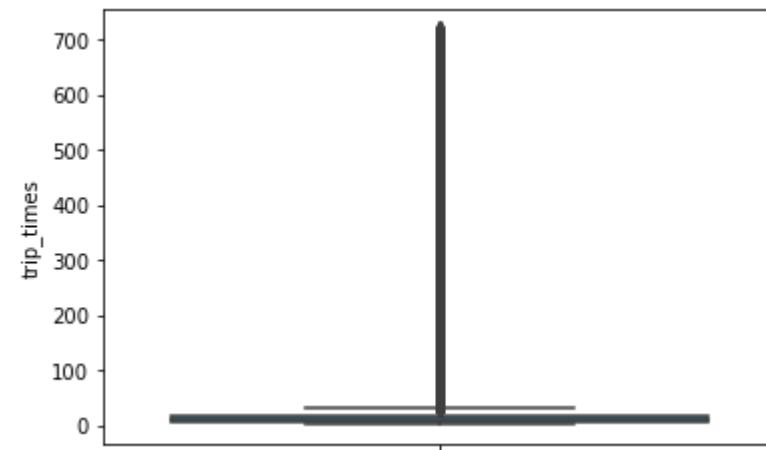
```
In [11]: 1 #Looking further from the 99th percetnile
2 for i in range(90,100):
3     var =frame_with_durations["trip_times"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6     print("100 percentile value is ",var[-1])
```

```
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 34.466666666666667
98 percentile value is 38.716666666666667
99 percentile value is 46.75
100 percentile value is 548555.6333333333
```

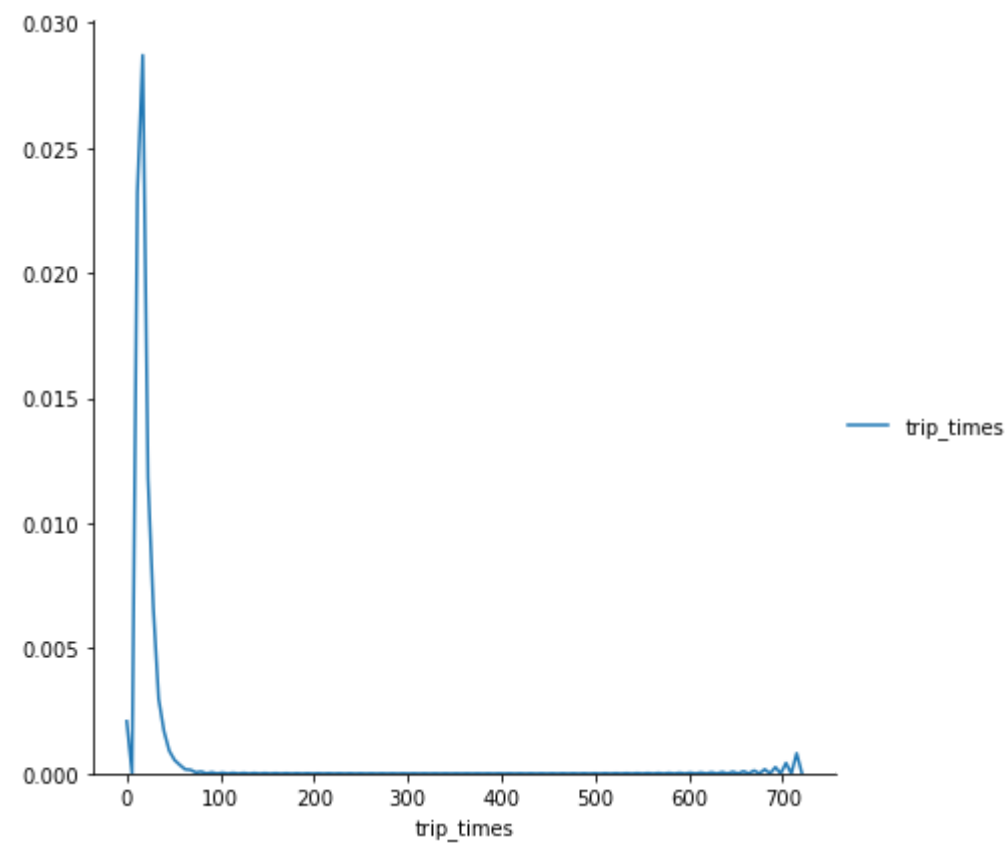
```
In [12]: 1 #removing data based on our analysis and TLC regulations
2 frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_times>1) & (frame_with_durations.trip_times< 720)]
```



```
In [13]: 1 #box-plot after removal of outliers
2 %matplotlib inline
3 sns.boxplot(y="trip_times", data =frame_with_durations_modified)
4 plt.show()
```

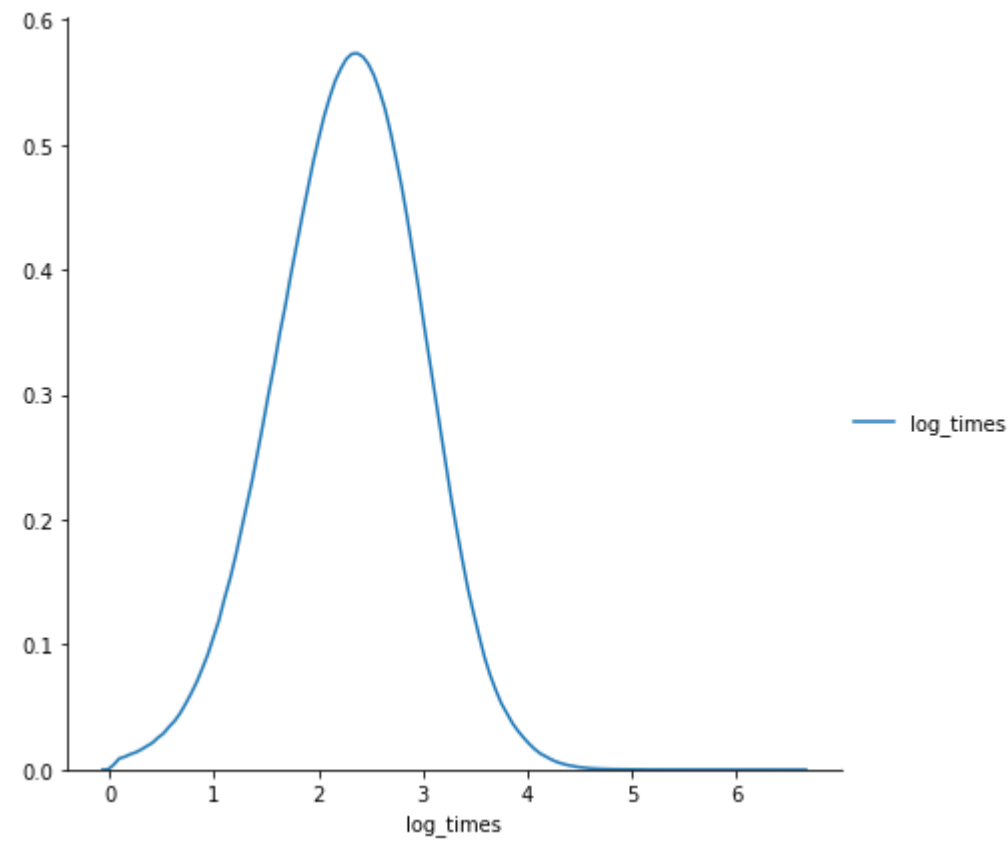


```
In [14]: 1 #pdf of trip-times after removing the outliers
2 sns.FacetGrid(frame_with_durations_modified,size=6) \
3     .map(sns.kdeplot,"trip_times") \
4     .add_legend();
5 plt.show();
```

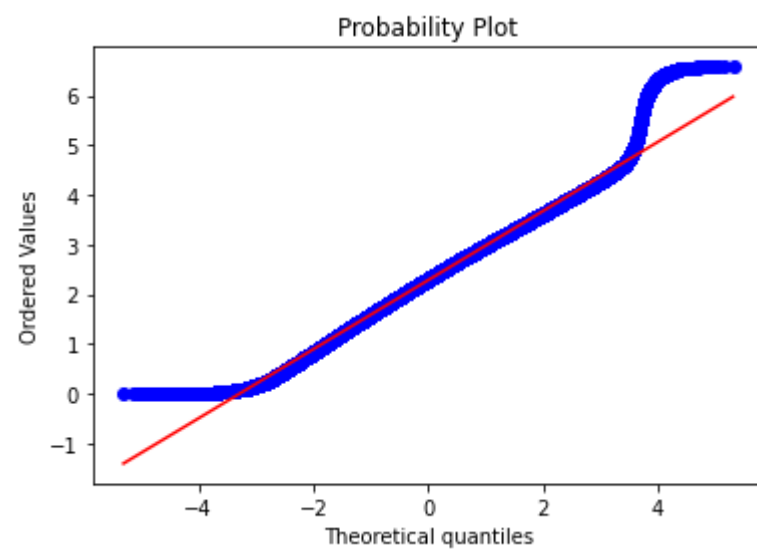


```
In [15]: 1 #converting the values to log-values to chec for log-normal
2 import math
3 frame_with_durations_modified['log_times']=[math.log(i) for i in frame_with_durations_modified['trip_times'].values]
```

```
In [16]: 1 #pdf of Log-values
2 sns.FacetGrid(frame_with_durations_modified,size=6) \
3     .map(sns.kdeplot,"log_times") \
4     .add_legend();
5 plt.show();
```

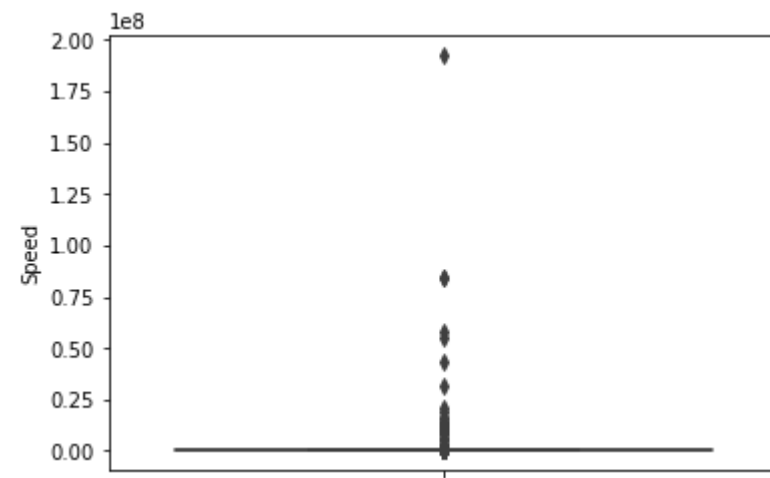


```
In [17]: 1 #Q-Q plot for checking if trip-times is log-normal
2 import scipy
3 scipy.stats.probplot(frame_with_durations_modified['log_times'].values, plot=plt)
4 plt.show()
```



#### 4. Speed

```
In [18]: 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_durations_modified['trip_times'])
4 sns.boxplot(y="Speed", data =frame_with_durations_modified)
5 plt.show()
```



```
In [19]: 1 #calculating speed values at each percentile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var =frame_with_durations_modified["Speed"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6     print("100 percentile value is ",var[-1])
```

```
0 percentile value is 0.0
10 percentile value is 6.409495548961425
20 percentile value is 7.80952380952381
30 percentile value is 8.929133858267717
40 percentile value is 9.98019801980198
50 percentile value is 11.06865671641791
60 percentile value is 12.286689419795222
70 percentile value is 13.796407185628745
80 percentile value is 15.963224893917962
90 percentile value is 20.186915887850468
100 percentile value is 192857142.85714284
```

```
In [20]: 1 #calculating speed values at each percentile 90,91,92,93,94,95,96,97,98,99,100
2 for i in range(90,100):
3     var =frame_with_durations_modified["Speed"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6     print("100 percentile value is ",var[-1])
```

```
90 percentile value is 20.186915887850468
91 percentile value is 20.91645569620253
92 percentile value is 21.752988047808763
93 percentile value is 22.721893491124263
94 percentile value is 23.844155844155843
95 percentile value is 25.182552504038775
96 percentile value is 26.80851063829787
97 percentile value is 28.84304932735426
98 percentile value is 31.591128254580514
99 percentile value is 35.7513566847558
100 percentile value is 192857142.85714284
```

```
In [21]: 1 #calculating speed values at each percentile 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])
```

```
99.0 percentile value is 35.7513566847558
99.1 percentile value is 36.31084727468969
99.2 percentile value is 36.91470054446461
99.3 percentile value is 37.588235294117645
99.4 percentile value is 38.33035714285714
99.5 percentile value is 39.17580340264651
99.6 percentile value is 40.15384615384615
99.7 percentile value is 41.338301043219076
99.8 percentile value is 42.86631016042781
99.9 percentile value is 45.3107822410148
100 percentile value is 192857142.85714284
```

```
In [22]: 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.Speed<45.31)]
```

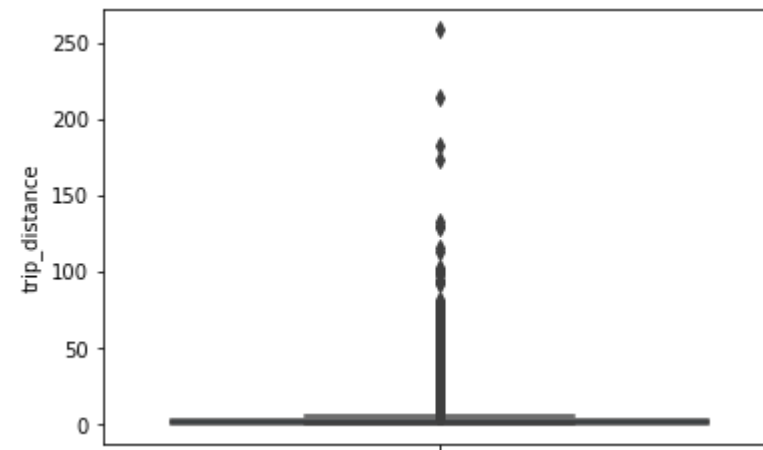
```
In [23]: 1 #avg.speed of cabs in New-York
2 sum(frame_with_durations_modified["Speed"]) / float(len(frame_with_durations_modified["Speed"]))
```

Out[23]: 12.450173996027528

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

#### 4. Trip Distance

```
In [24]: 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 [25]: 1 #calculating trip distance values at each percntile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var =frame_with_durations_modified["trip_distance"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6 print("100 percentile value is ",var[-1])
```

```
0 percentile value is 0.01
10 percentile value is 0.66
20 percentile value is 0.9
30 percentile value is 1.1
40 percentile value is 1.39
50 percentile value is 1.69
60 percentile value is 2.07
70 percentile value is 2.6
80 percentile value is 3.6
90 percentile value is 5.97
100 percentile value is 258.9
```



```
In [26]: 1 #calculating trip distance values at each percntile 90,91,92,93,94,95,96,97,98,99,100
2 for i in range(90,100):
3     var =frame_with_durations_modified["trip_distance"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6     print("100 percentile value is ",var[-1])
```

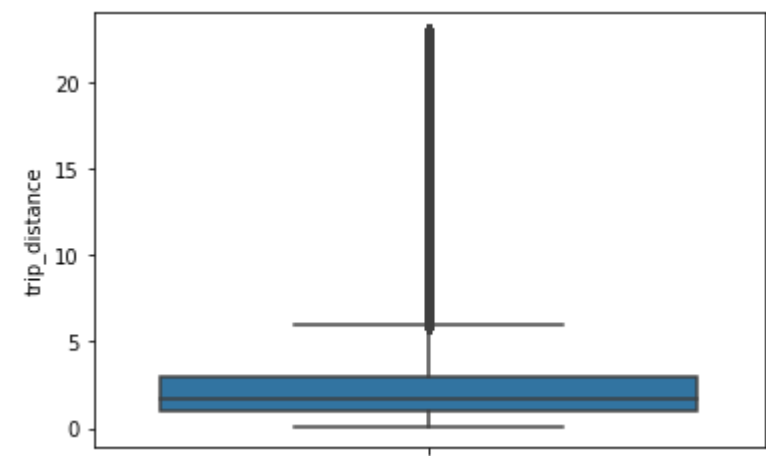
```
90 percentile value is 5.97
91 percentile value is 6.45
92 percentile value is 7.07
93 percentile value is 7.85
94 percentile value is 8.72
95 percentile value is 9.6
96 percentile value is 10.6
97 percentile value is 12.1
98 percentile value is 16.03
99 percentile value is 18.17
100 percentile value is 258.9
```

```
In [27]: 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])
```

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

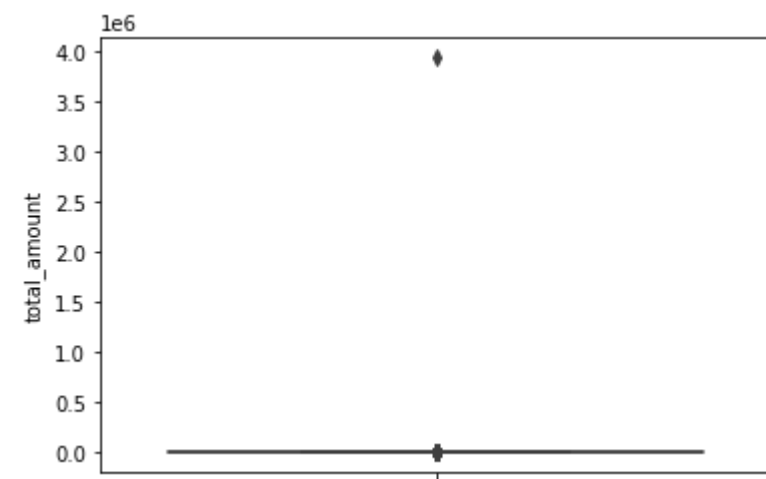
```
In [28]: 1 #removing further outliers based on the 99.9th percentile value
2 frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_distance>0) & (frame_with_durations.trip_distance<23)]
```

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



5. Total Fare

```
In [30]: 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()
```



```
In [31]: 1 #calculating total fare amount values at each percentile 0,10,20,30,40,50,60,70,80,90,100
2 for i in range(0,100,10):
3     var = frame_with_durations_modified["total_amount"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6     print("100 percentile value is ",var[-1])
```

```
0 percentile value is -242.55
10 percentile value is 6.3
20 percentile value is 7.8
30 percentile value is 8.8
40 percentile value is 9.8
50 percentile value is 11.16
60 percentile value is 12.8
70 percentile value is 14.8
80 percentile value is 18.3
90 percentile value is 25.8
100 percentile value is 3950611.6
```

```
In [32]: 1 #calculating total fare amount values at each percentile 90,91,92,93,94,95,96,97,98,99,100
2 for i in range(90,100):
3     var = frame_with_durations_modified["total_amount"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
6     print("100 percentile value is ",var[-1])
```

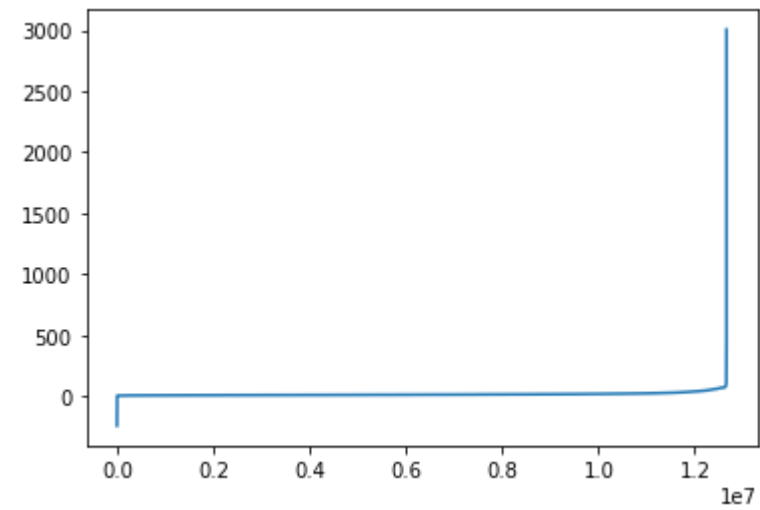
```
90 percentile value is 25.8
91 percentile value is 27.3
92 percentile value is 29.3
93 percentile value is 31.8
94 percentile value is 34.8
95 percentile value is 38.53
96 percentile value is 42.6
97 percentile value is 48.13
98 percentile value is 58.13
99 percentile value is 66.13
100 percentile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

```
In [33]: 1 #calculating total fare amount values at each percentile 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["total_amount"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
6     print("100 percentile value is ",var[-1])
```

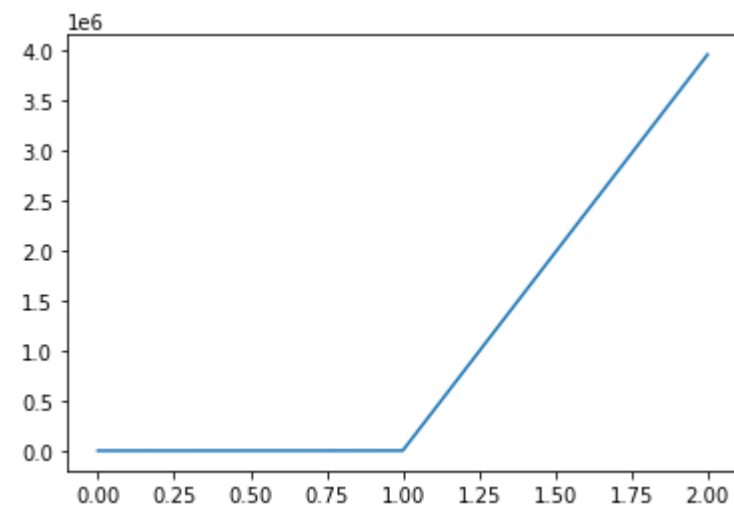
```
99.0 percentile value is 66.13
99.1 percentile value is 68.13
99.2 percentile value is 69.6
99.3 percentile value is 69.6
99.4 percentile value is 69.73
99.5 percentile value is 69.75
99.6 percentile value is 69.76
99.7 percentile value is 72.58
99.8 percentile value is 75.35
99.9 percentile value is 88.28
100 percentile value is 3950611.6
```

**Observation:-** As even the 99.9th percentile value doesn't 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 analysis

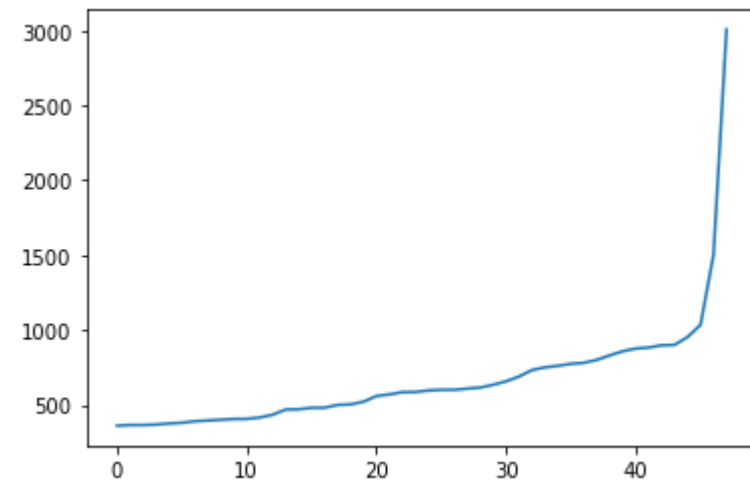
```
In [34]: 1 #below plot shows us the fare values(sorted) to find a sharp increase to remove those values as outliers
2 # plot the fare amount excluding last two values in sorted data
3 plt.plot(var[:-2])
4 plt.show()
```



```
In [35]: 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 [36]: 1 # now looking at values not including the last two points we again find a drastic increase at around 1000 fare value  
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 [37]: 1 #removing all outliers based on our univariate analysis above
2 def remove_outliers(new_frame):
3     a = new_frame.shape[0]
4     print ("Number of pickup records = ",a)
5     temp_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_frame.dropoff_longitude <= -73.7004) & \
6                             (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitude <= 40.9176)) & \
7                             ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_latitude >= 40.5774) & \
8                             (new_frame.pickup_longitude <= -73.7004) & (new_frame.pickup_latitude <= 40.9176))]
9     b = temp_frame.shape[0]
10    print ("Number of outlier coordinates lying outside NY boundaries:",(a-b))
11
12
13    temp_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times < 720)]
14    c = temp_frame.shape[0]
15    print ("Number of outliers from trip times analysis:",(a-c))
16
17
18    temp_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_distance < 23)]
19    d = temp_frame.shape[0]
20    print ("Number of outliers from trip distance analysis:",(a-d))
21
22    temp_frame = new_frame[(new_frame.Speed <= 65) & (new_frame.Speed >= 0)]
23    e = temp_frame.shape[0]
24    print ("Number of outliers from speed analysis:",(a-e))
25
26    temp_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_amount >0)]
27    f = temp_frame.shape[0]
28    print ("Number of outliers from fare analysis:",(a-f))
29
30
31    new_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_frame.dropoff_longitude <= -73.7004) & \
32                            (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitude <= 40.9176)) & \
33                            ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_latitude >= 40.5774) & \
34                            (new_frame.pickup_longitude <= -73.7004) & (new_frame.pickup_latitude <= 40.9176))]
35
36    new_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times < 720)]
37    new_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_distance < 23)]
38    new_frame = new_frame[(new_frame.Speed < 45.31) & (new_frame.Speed > 0)]
39    new_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_amount >0)]
40
41    print ("Total outliers removed",a - new_frame.shape[0])
42    print ("---")
43    return new_frame
```

```
In [38]: ► 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_outliers_removed))/len(frame_with_durations))
```

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 [39]: 1  ### we are gonna create different clusters first from 10 to 100
2  ## then we will send the cluster len and cluster_centers to our functions to find points within distance of 2 miles
3  coords = frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']].values
4  neighbors = []
5  def find_min_distance(cluster_centers,cluster_len):
6      ## function that finds min distance between two coordinates
7      ## And average good neighbors between two points and average bad neighbors between two points
8      less_,more_ = [],[]
9      min_dist = 1000
10     for i in range(0,cluster_len):
11         nice_points,wrong_points = 0,0
12         for j in range(0,cluster_len):
13             if j!=i:
14                 distance = gpxpy.geo.haversine_distance(cluster_centers[i][0],cluster_centers[i][1],
15                                                         cluster_centers[j][0],cluster_centers[j][1])
16                 distance = distance/(1.60934*1000)
17                 min_dist = min(distance,min_dist)
18                 if distance <= 2:
19                     nice_points+=1
20                 else:
21                     wrong_points+=1
22             less_.append(nice_points)
23             more_.append(wrong_points)
24
25     neighbors.append(less_)
26     print('After choosing the clusters with length ',cluster_len)
27     print('***10)
28     print('Average number of points with intercluster distance <= 2 ',np.ceil(sum(less_)/len(less_)))
29     print('Average number of points with intercluster distance > 2 ',np.ceil(sum(more_)/len(more_)))
30     print('Minimum distance calculated to be :',min_dist)
31     print('***10)
32
33     def find_centroids(increment):
34         kmeans = MiniBatchKMeans(n_clusters=increment,batch_size=1000,random_state=42).fit(coords)
35         frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(coords)
36         cluster_centers = kmeans.cluster_centers_
37         cluster_len = len(cluster_centers)
38         return cluster_centers,cluster_len
39
40     for increment in range(10,100,10):
41         cluster_centers,cluster_len = find_centroids(increment)
42         find_min_distance(cluster_centers,cluster_len)

```

After choosing the clusters with length 10

\*\*\*\*\*

Average number of points with intercluster distance <= 2 2.0

Average number of points with intercluster distance > 2 8.0

Minimum distance calculated to be : 1.1192643504338826

\*\*\*\*\*

After choosing the clusters with length 20

\*\*\*\*\*

Average number of points with intercluster distance <= 2 5.0

Average number of points with intercluster distance > 2 15.0

Minimum distance calculated to be : 0.5068296118318109

\*\*\*\*\*

After choosing the clusters with length 30

\*\*\*\*\*

Average number of points with intercluster distance <= 2 8.0

Average number of points with intercluster distance > 2 22.0

```
Minimum distance calculated to be : 0.2977082261311233
*****
After choosing the clusters with length 40
*****
Average number of points with intercluster distance <= 2 10.0
Average number of points with intercluster distance > 2 30.0
Minimum distance calculated to be : 0.3486854241772553
*****
After choosing the clusters with length 50
*****
Average number of points with intercluster distance <= 2 12.0
Average number of points with intercluster distance > 2 38.0
Minimum distance calculated to be : 0.2645808955207251
*****
After choosing the clusters with length 60
*****
Average number of points with intercluster distance <= 2 16.0
Average number of points with intercluster distance > 2 43.0
Minimum distance calculated to be : 0.23245490125677482
*****
After choosing the clusters with length 70
*****
Average number of points with intercluster distance <= 2 19.0
Average number of points with intercluster distance > 2 51.0
Minimum distance calculated to be : 0.05834656875677777
*****
After choosing the clusters with length 80
*****
Average number of points with intercluster distance <= 2 22.0
Average number of points with intercluster distance > 2 58.0
Minimum distance calculated to be : 0.17930841488997812
*****
After choosing the clusters with length 90
*****
Average number of points with intercluster distance <= 2 28.0
Average number of points with intercluster distance > 2 62.0
Minimum distance calculated to be : 0.11384781001034917
*****
```

## 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 20.
- We need minimum distance to be atleast 0.5 miles.

In [40]:

▶

1 frame\_with\_durations\_outliers\_removed

Out[40]:

	passenger_count	trip_distance	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	total_amount	trip_times	pickup_times	Speed	pickup_cluster
0	1	1.59	-73.993896	40.750111	-73.974785	40.750618	17.05	18.050000	1.421329e+09	5.285319	0
1	1	3.30	-74.001648	40.724243	-73.994415	40.759109	17.80	19.833333	1.420902e+09	9.983193	34
2	1	1.80	-73.963341	40.802788	-73.951820	40.824413	10.80	10.050000	1.420902e+09	10.746269	8
3	1	0.50	-74.009087	40.713818	-74.004326	40.719986	4.80	1.866667	1.420902e+09	16.071429	54
4	1	3.00	-73.971176	40.762428	-74.004181	40.742653	16.30	19.316667	1.420902e+09	9.318378	44
...	...	...	...	...	...	...	...	...	...	...	...
12615	2	1.00	-73.951988	40.786217	-73.953735	40.775162	7.55	3.933333	1.420897e+09	15.254237	33
12616	2	0.80	-73.982742	40.728184	-73.974976	40.720013	8.80	5.700000	1.420897e+09	8.421053	23
12617	1	3.40	-73.979324	40.749550	-73.969101	40.787800	14.30	13.283333	1.420897e+09	15.357591	45
12618	1	1.30	-73.999565	40.738483	-73.981819	40.737652	13.55	15.316667	1.420897e+09	5.092492	11
12619	1	0.70	-73.960350	40.766399	-73.968643	40.760777	6.30	5.800000	1.420897e+09	7.241379	24

12371076 rows × 11 columns

In [41]:

▶

1 kmeans = MiniBatchKMeans(n\_clusters=20, batch\_size=1000, random\_state=42).fit(coords)  
2 frame\_with\_durations\_outliers\_removed['pickup\_cluster'] = kmeans.predict(frame\_with\_durations\_outliers\_removed[['pickup\_latitude', 'pickup\_longitude']])

Plotting the cluster centers:



```

In [294]: 1 # Plotting the cluster centers on OSM
2 cluster_centers = kmeans.cluster_centers_
3 cluster_len = len(cluster_centers)
4 map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')
5 for i in range(cluster_len):
6     folium.Marker(list((cluster_centers[i][0], cluster_centers[i][1])),
7                   popup=(str(cluster_centers[i][0]) + str(cluster_centers[i][1]))).add_to(map_osm)
8 map_osm

```

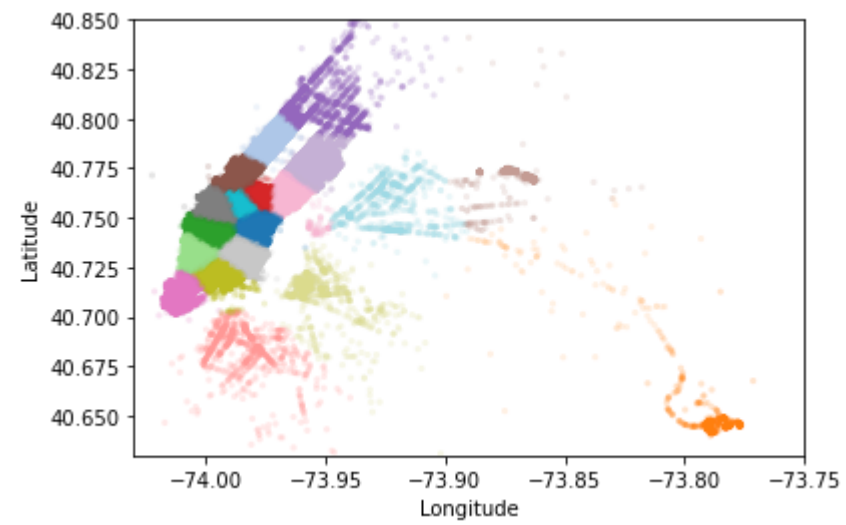
Out[294]:



Leaflet (<https://leafletjs.com>) | Map tiles by Stamen Design (<http://stamen.com>), under CC BY 3.0 (<http://creativecommons.org/licenses/by/3.0>). Data by © OpenStreetMap (<http://openstreetmap.org>), under ODbL (<http://www.openstreetmap.org/copyright>).

## Plotting the cluster

```
In [43]: ▶ 1 def plot_cluster(frame):  
2     city_long_border = (-74.03, -73.75)  
3     city_lat_border = (40.63, 40.85)  
4     fig, ax = plt.subplots(ncols=1, rows=1)  
5     plt.scatter(frame.pickup_longitude.values[:100000], frame.pickup_latitude.values[:100000], s=10, lw=0  
6                 , c=frame.pickup_cluster.values[:100000], cmap='tab20', alpha=0.2)  
7     ax.set_xlim(city_long_border)  
8     ax.set_ylim(city_lat_border)  
9     ax.set_xlabel('Longitude')  
10    ax.set_ylabel('Latitude')  
11    plt.show()  
12  
13 plot_cluster(frame_with_durations_outliers_removed)
```



## Time-binning

In [44]:

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

In [45]:

```
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_removed[['pickup_latitude', 'pickup_longitude']])
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_cluster','pickup_bins']).count()
```

In [46]:

```
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[46]:

		trip_distance
pickup_cluster	pickup_bins	
	0	1
0	1	167
	2	340
	3	432
	4	514
	5	520

```

In [47]: 1 # upto now we cleaned data and prepared data for the month 2015,
2
3 # now do the same operations for months Jan, Feb, March of 2016
4 # 1. get the dataframe which includes only required columns
5 # 2. adding trip times, speed, unix time stamp of pickup_time
6 # 4. remove the outliers based on trip_times, speed, trip_duration, total_amount
7 # 5. add pickup_cluster to each data point
8 # 6. add pickup_bin (index of 10min intravel to which that trip belongs to)
9 # 7. group by data, based on 'pickup_cluster' and 'pickup_bin'
10
11 # Data Preparation for the months of Jan, Feb and March 2016
12 def datapreparation(month, kmeans, month_no, year_no):
13
14     print ("Return with trip times..")
15
16     frame_with_durations = return_with_trip_times(month)
17
18     print ("Remove outliers..")
19     frame_with_durations_outliers_removed = remove_outliers(frame_with_durations)
20
21     print ("Estimating clusters..")
22     frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']])
23     #frame_with_durations_outliers_removed_2016['pickup_cluster'] = kmeans.predict(frame_with_durations_outliers_removed_2016[['pickup_latitude', 'pickup_longitude']])
24
25     print ("Final groupbying..")
26     final_updated_frame = add_pickup_bins(frame_with_durations_outliers_removed, month_no, year_no)
27     final_groupby_frame = final_updated_frame[['pickup_cluster', 'pickup_bins', 'trip_distance']].groupby(['pickup_cluster', 'pickup_bins']).count()
28
29     return final_updated_frame, final_groupby_frame
30
31 month_jan_2016 = dd.read_csv('yellow_tripdata_2016-01.csv')
32 month_feb_2016 = dd.read_csv('yellow_tripdata_2016-02.csv')
33 month_mar_2016 = dd.read_csv('yellow_tripdata_2016-03.csv')
34
35 jan_2016_frame, jan_2016_groupby = datapreparation(month_jan_2016, kmeans, 1, 2016)
36 feb_2016_frame, feb_2016_groupby = datapreparation(month_feb_2016, kmeans, 2, 2016)
37 mar_2016_frame, mar_2016_groupby = datapreparation(month_mar_2016, kmeans, 3, 2016)

```

```

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

```

```

---
Estimating clusters..
Final groupbying..
Return with trip times..
Remove outliers..
Number of pickup records = 12210952
Number of outlier coordinates lying outside NY boundaries: 232444
Number of outliers from trip times analysis: 30868
Number of outliers from trip distance analysis: 87318
Number of outliers from speed analysis: 23889
Number of outliers from fare analysis: 5859
Total outliers removed 324635
---
Estimating clusters..
Final groupbying..

```

## Smoothing

```

In [48]: ► 1 # Gets the unique bins where pickup values are present for each each reigion
2 # for each cluster region we will collect all the indices of 10min intravels in which the pickups are happened
3 # we got an observation that there are some pickpbins that doesnt have any pickups
4 def return_unq_pickupbins(frame):
5     values=[]
6     for i in range(0,20):
7         new = frame[frame['pickup_cluster'] == i ]
8         list_unq = list(set(new['pickup_bins']))
9         list_unq.sort()
10        values.append(list_unq)
11    return values

```

```

In [49]: ► 1 ## unique_pickup bins for jan feb march
2 jan_2015_unique = return_unq_pickupbins(jan_2015_frame)
3 jan_2016_unique = return_unq_pickupbins(jan_2016_frame)
4
5 #feb
6 feb_2016_unique = return_unq_pickupbins(feb_2016_frame)
7
8 #march
9 march_2016_unique = return_unq_pickupbins(mar_2016_frame)

```



```
In [50]: 1 for i in range(20):  
2         print('Number of bins with zero pickups of 10mins interval for the cluster ',i,'is',4464 - len(set(jan_2015_unique[i])))  
3         print(' '*100)
```

```
Number of bins with zero pickups of 10mins interval for the cluster  0 is 29  
*****  
Number of bins with zero pickups of 10mins interval for the cluster  1 is 36  
*****  
Number of bins with zero pickups of 10mins interval for the cluster  2 is 143  
*****  
Number of bins with zero pickups of 10mins interval for the cluster  3 is 461  
*****  
Number of bins with zero pickups of 10mins interval for the cluster  4 is 32  
*****  
Number of bins with zero pickups of 10mins interval for the cluster  5 is 27  
*****  
Number of bins with zero pickups of 10mins interval for the cluster  6 is 30  
*****  
Number of bins with zero pickups of 10mins interval for the cluster  7 is 41  
*****  
Number of bins with zero pickups of 10mins interval for the cluster  8 is 18  
*****  
Number of bins with zero pickups of 10mins interval for the cluster  9 is 15  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 10 is 29  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 11 is 58  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 12 is 34  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 13 is 18  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 14 is 36  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 15 is 31  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 16 is 29  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 17 is 82  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 18 is 31  
*****  
Number of bins with zero pickups of 10mins interval for the cluster 19 is 27  
*****
```

```
In [51]: ► 1 ### now Lets define functions for filling the missing values with zeros
2 def fill_missing(count_values,values):
3     smoothed_regions=[]
4     ind=0
5     for r in range(0,20):
6         smoothed_bins=[]
7         for i in range(4464):
8             if i in values[r]:
9                 smoothed_bins.append(count_values[ind])
10                ind+=1
11            else:
12                smoothed_bins.append(0)
13        smoothed_regions.extend(smoothed_bins)
14    return smoothed_regions
15
```



```

In [52]: 1 def smoothening_using_avg(pickupvalues,uniquetimebins):
2         idx = 0
3         smoothed_region = []
4         repeat = 0
5         for clstr in range(0,20): ## using optimal number of clusters
6             smoothed_bins = []
7             for i in range(4464): ### looping over total number of bins
8                 if repeat != 0: ## to skip if the pickup bin is already resolved
9                     repeat-=1
10                    continue
11                if i in uniquetimebins[clstr]: ## check if a pick value exists for the bin
12                    smoothed_bins.append(pickupvalues[idx])
13                else:
14                    ## checking the condition if zero pickups are happening at the beginning
15                    right_hand_limit = 0
16                    for j in range(i,4464): ## set right hand limit to the index where a pickupbin with value exists after 0's
17                        if j not in uniquetimebins[clstr]:
18                            continue
19                        else:
20                            right_hand_limit=j
21                            break
22                    if i==0 :
23                        #-----case 1 :pickups missing in the beginning -----
24                        smoothed_value = pickupvalues[idx]*1.0/((right_hand_limit-i)+1)*1.0
25                        for j in range(i,right_hand_limit+1):
26                            smoothed_bins.append(math.ceil(smoothed_value))
27                        repeat = right_hand_limit-i
28                    else :
29                        # #-----case2: pickups missing at the end -----
30                        if right_hand_limit == 0 :
31                            smoothed_value = pickupvalues[idx-1]*1.0/((4464-i)+1)
32                            for j in range(i,4464):
33                                smoothed_bins.append(math.ceil(smoothed_value))
34                                smoothed_bins[i-1] = math.ceil(smoothed_value)
35                                repeat = right_hand_limit-i
36                        ##-----case3:pickups missing in the middle -----
37                        else:
38                            smoothed_value = (pickupvalues[idx]+pickupvalues[idx-1])*1.0/((right_hand_limit - i)+2)*1.0
39                            for j in range(i,right_hand_limit+1):
40                                smoothed_bins.append(math.ceil(smoothed_value))
41                                smoothed_bins[i-1] = math.ceil(smoothed_value)
42                                repeat = right_hand_limit-i
43                            idx+=1
44                            smoothed_region.extend(smoothed_bins)
45                        return smoothed_region

```

```

In [53]: 1 #Filling Missing values of Jan-2015 with 0
2 # here in jan_2015_groupby dataframe the trip_distance represents the number of pickups that are happened
3 jan_2015_fill = fill_missing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)
4 ## let's fill the missing values with average values
5 jan_2015_smoothed_values = smoothening_using_avg(jan_2015_groupby['trip_distance'].values,jan_2015_unique)

```

```
In [54]: 1 ### Lets do a sanity check for missing values
2 ## jan2015 data split to 10 mins interval will have 4464 bins totally , (i.e 6*24*31 = 4464)
3 ## there are 20 clusters totally hence number of values should be 20 * 4464 = 89280
4 print('total number of zero filled values :', len(jan_2015_fill))
5 print('total number of zero filled values :', len(jan_2015_smoothed_values))
```

total number of zero filled values : 89280  
total number of zero filled values : 89280

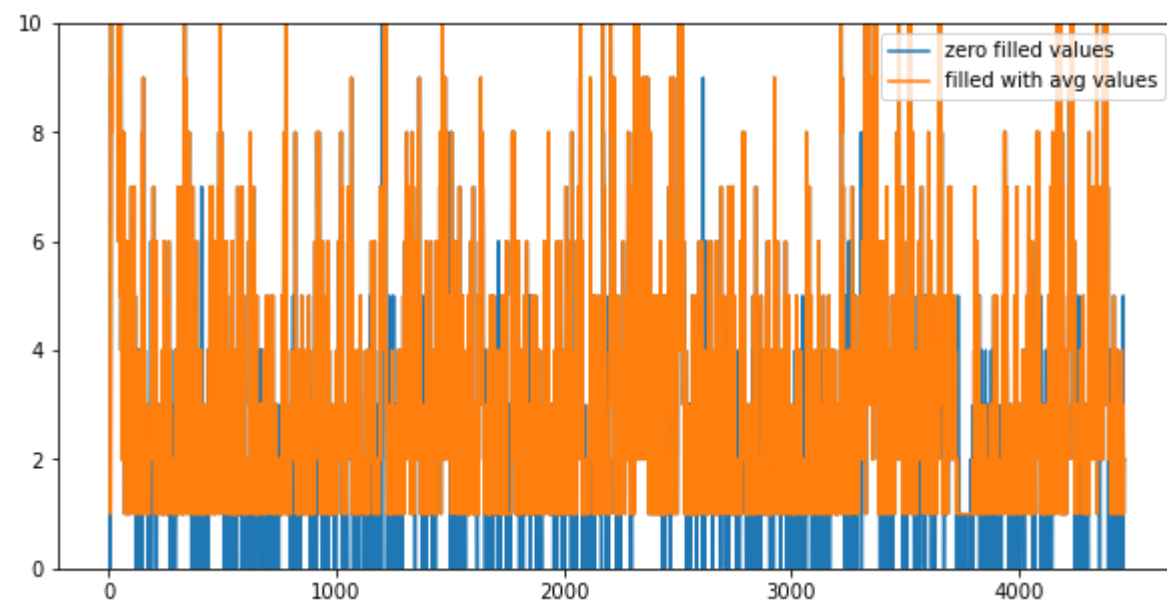
```
In [55]: 1 ## Lets check if there are any zeros after smoothening
2 count=0
3 val = [count for i in jan_2015_fill if i==0 ]
4 print('Number of pickups having zero values after smoothening using zero fill is: ', len(val))
```

Number of pickups having zero values after smoothening using zero fill is: 1207

```
In [56]: 1 ## Lets check if there are any zeros after smoothening using average
2 count=0
3 val = [count for i in jan_2015_smoothed_values if i==0 ]
4 print('Number of pickups having zero values after smoothening is: ', len(val))
```

Number of pickups having zero values after smoothening is: 0

```
In [57]: 1 ### plot smoothing vs filling
2 plt.figure(figsize=(10,5))
3 plt.plot(jan_2015_fill[13384:17848], label="zero filled values")
4 plt.plot(jan_2015_smoothed_values[13384:17848], label="filled with avg values")
5 plt.ylim(0,10)
6 plt.legend()
7 plt.show()
```



```
In [58]: 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 happened in 1st
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 the future values
9 # when you are using smoothing we are looking at the future number of pickups which might cause a data leakage.
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 [59]: 1 # Jan-2015 data is smoothed, Jan, Feb & March 2016 data missing values are filled with zero
2 jan_2015_smooth = smoothening_using_avg(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, march_2016_unique)
```

```
In [60]: 1 # Making list of all the values of pickup data in every bin for a period of 3 months and storing them region-wise
2 regions_cum = []
3
4 ## example :
5 ## a = [1,2,3]
6 ## b = [4,5,6]
7 ## a+b = [1,2,3,4,5,6]
8
9 ## number of 10 min interval in the month of jan 2016 is 24*31*60/10 = 4464
10 ## number of 10 min interval in the month of feb 2016 is 24*29*60/10 = 4176
11 ## number of 10 min interval in the month of march 2016 is 24*31*60/10 = 4464
12 for i in range(20):
13     regions_cum.append(jan_2016_smooth[4464*i:4464*(i+1)]+feb_2016_smooth[4176*i:4176*(i+1)]+mar_2016_smooth[4464*i:4464*(i+1)])
```

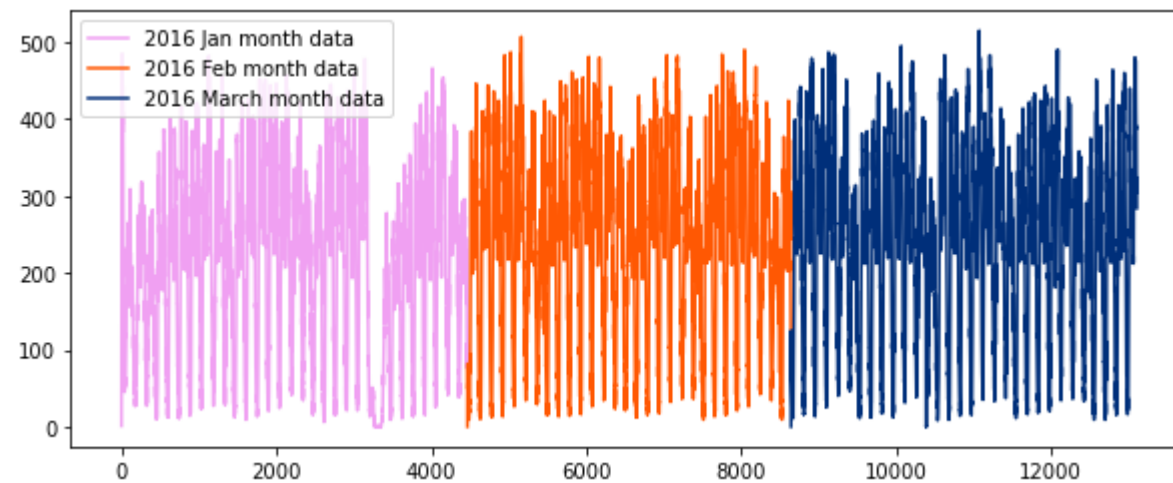
## Time series and Fourier Transforms

```

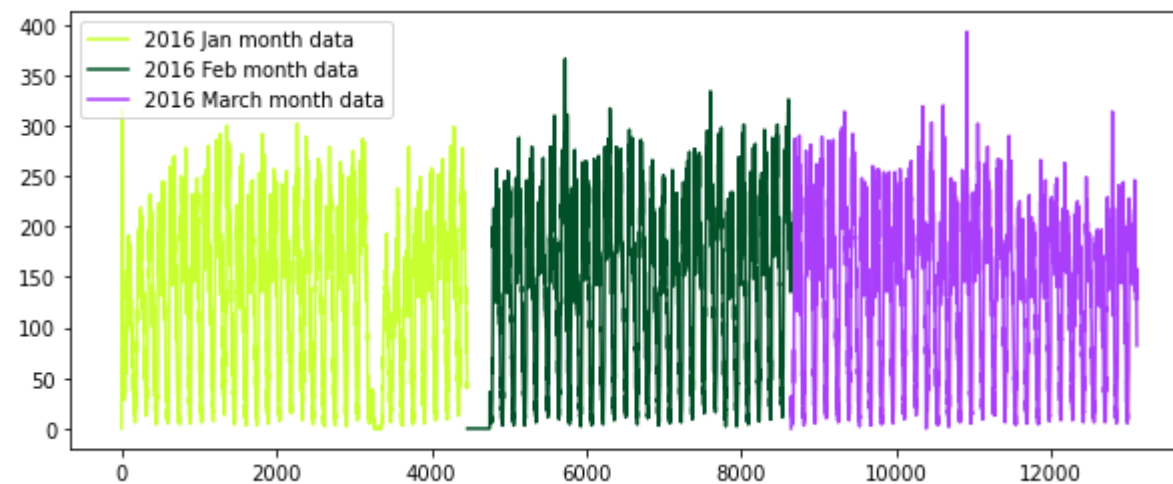
In [61]: 1 def unique_colors():
2         """ There're better ways to generate unique colors"""
3         return plt.cm.gist_ncar(np.random.random())
4
5     ### Let's generate unique colors
6     first_x = list(range(0,4464))
7     second_x = list(range(4464,8640))
8     third_x = list(range(8640,13104))
9
10    ## Lets plot the monthly data for each cluster
11    for i in range(20):
12        print('----- Month data for cluster ',i,'-----')
13        plt.figure(figsize=(10,4))
14        plt.plot(first_x,regions_cum[i][:4464],color=unique_colors(),label='2016 Jan month data')
15        plt.plot(second_x,regions_cum[i][4464:8640],color=unique_colors(),label='2016 Feb month data')
16        plt.plot(third_x,regions_cum[i][8640:13104],color=unique_colors(),label='2016 March month data')
17        plt.legend()
18        plt.show()

```

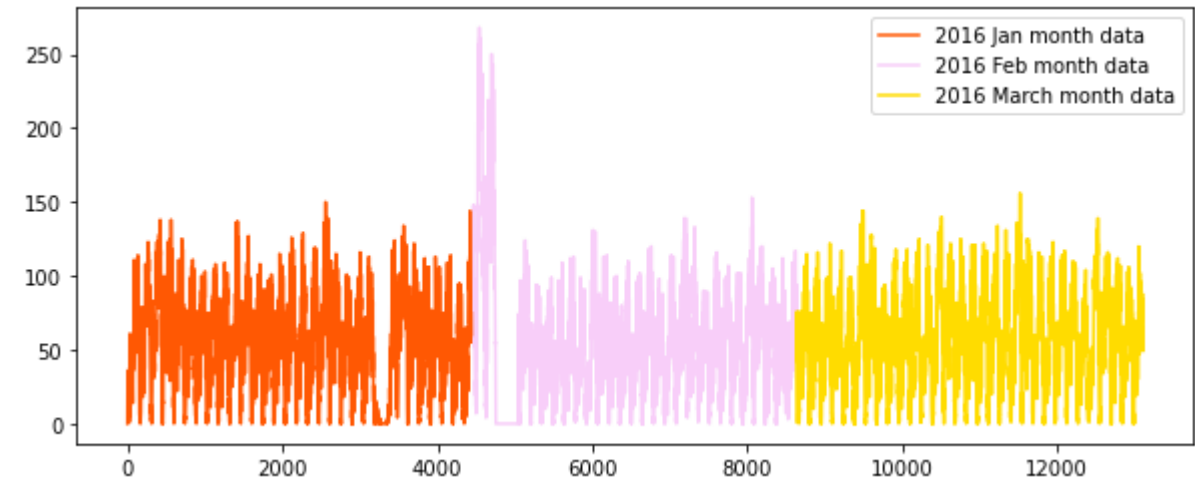
----- Month data for cluster 0 -----



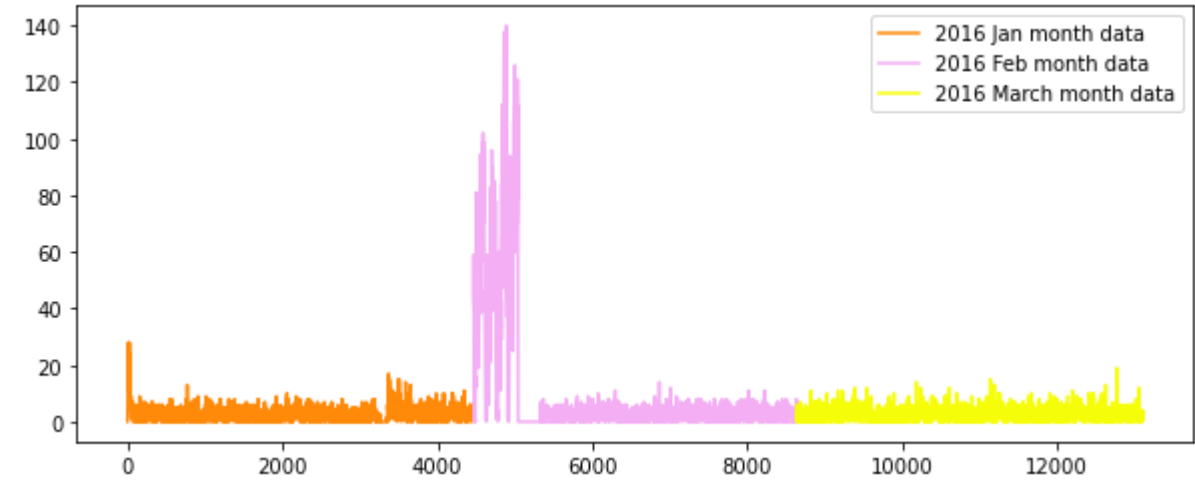
----- Month data for cluster 1 -----



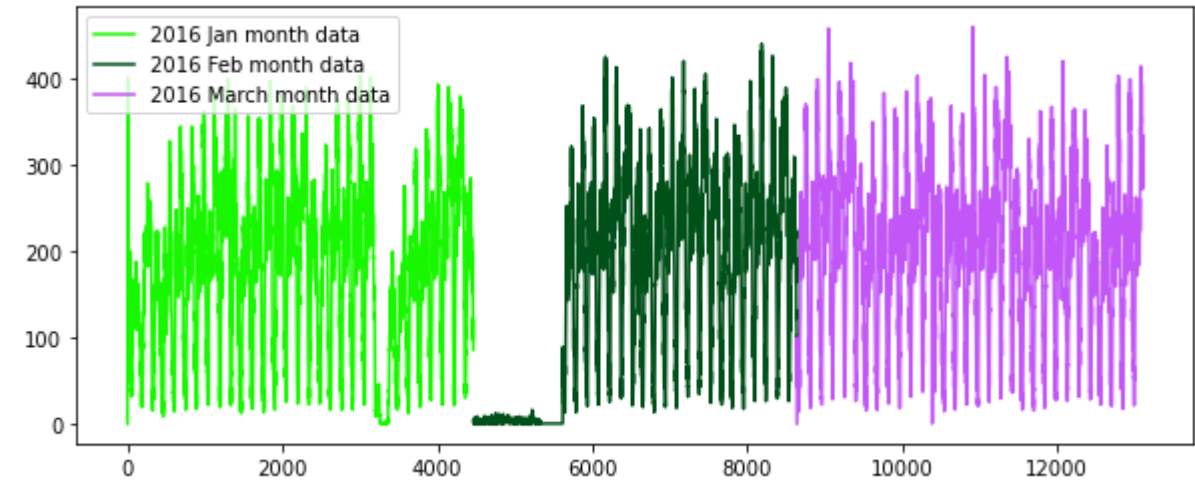
----- Month data for cluster 2 -----



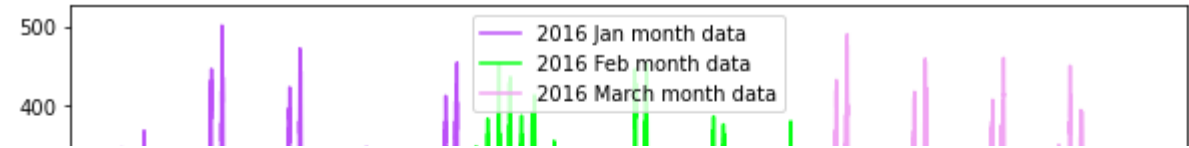
----- Month data for cluster 3 -----



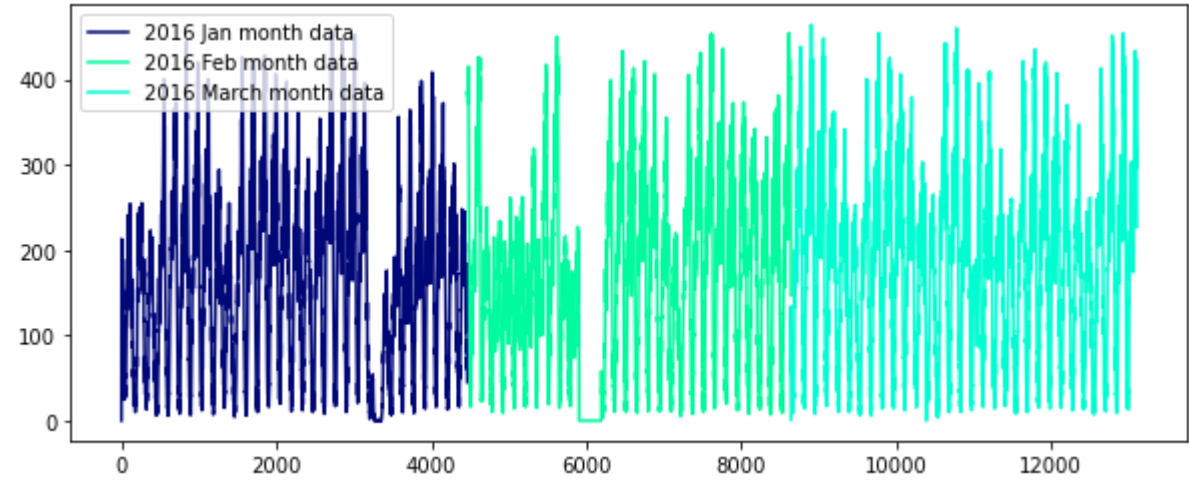
----- Month data for cluster 4 -----



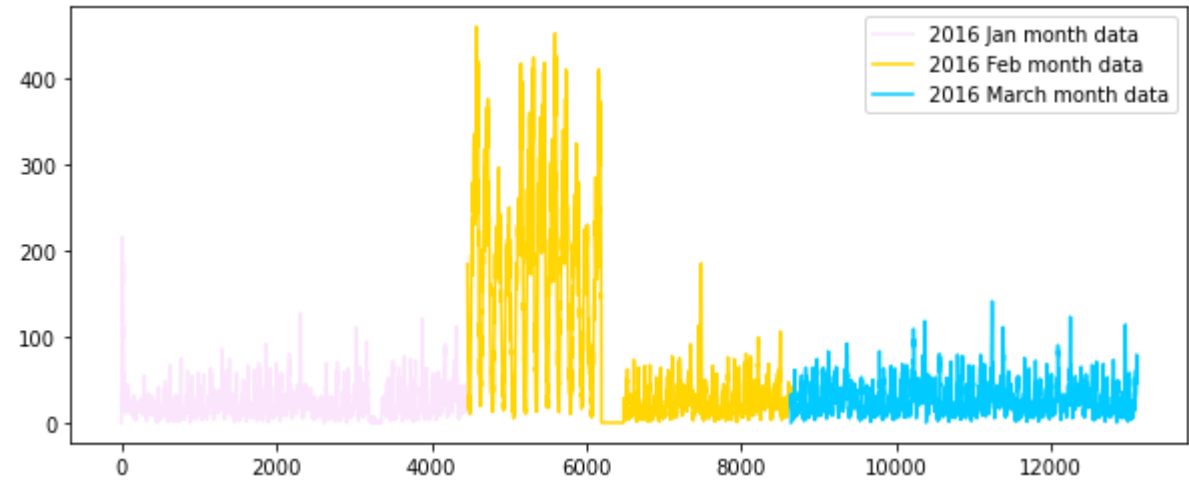
----- Month data for cluster 5 -----



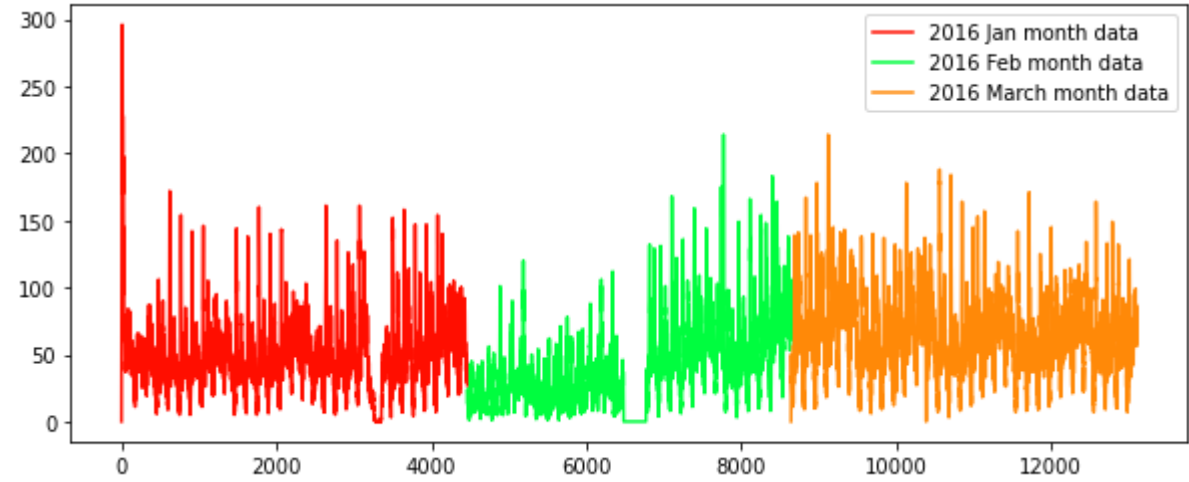
----- Month data for cluster 6 -----



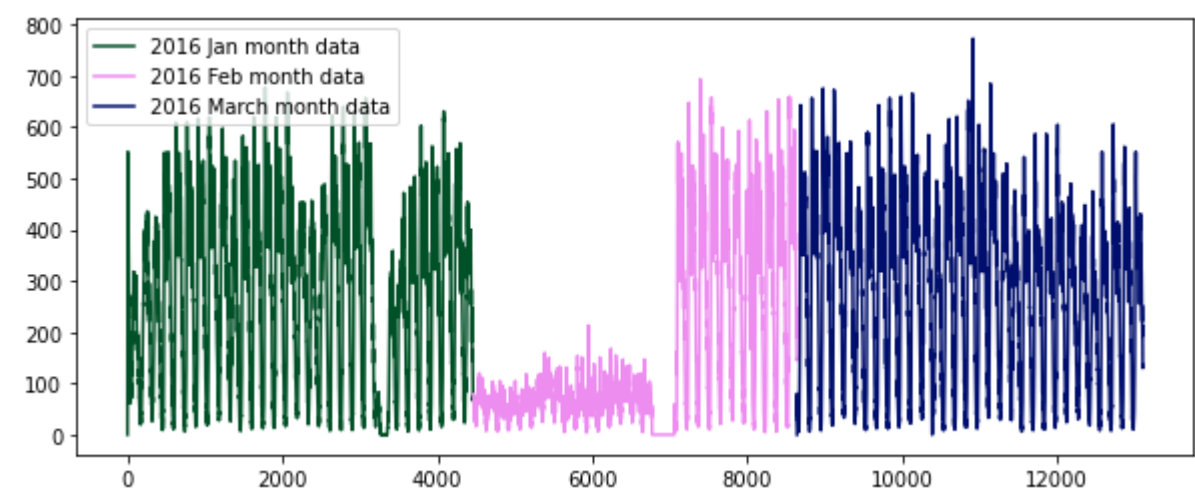
----- Month data for cluster 7 -----



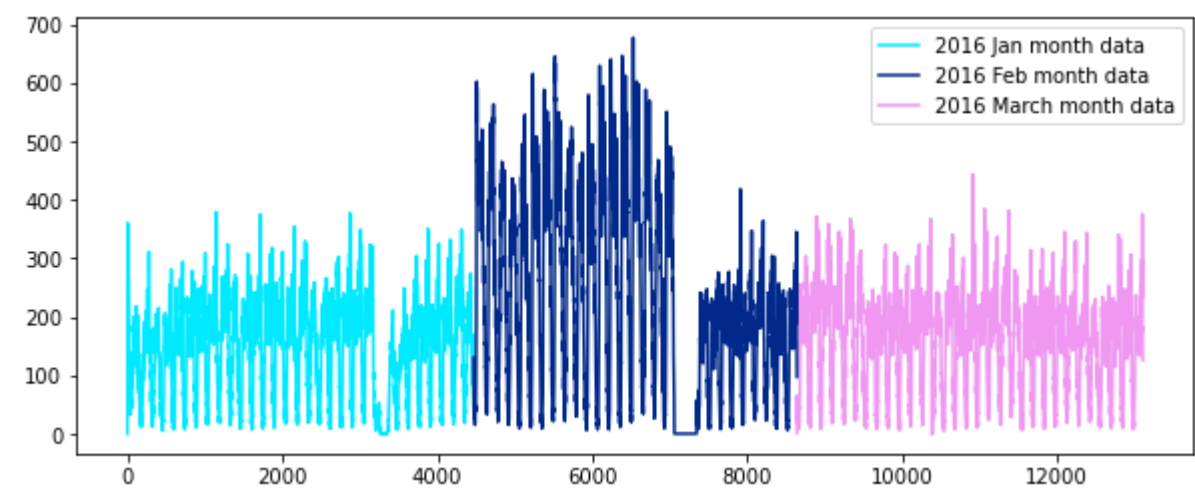
----- Month data for cluster 8 -----



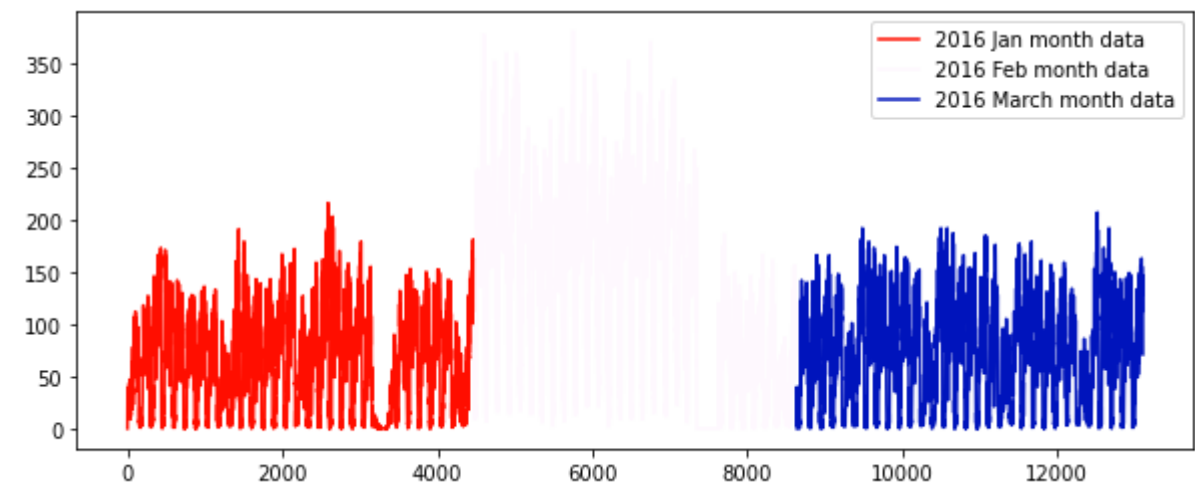
----- Month data for cluster 9 -----



----- Month data for cluster 10 -----



----- Month data for cluster 11 -----

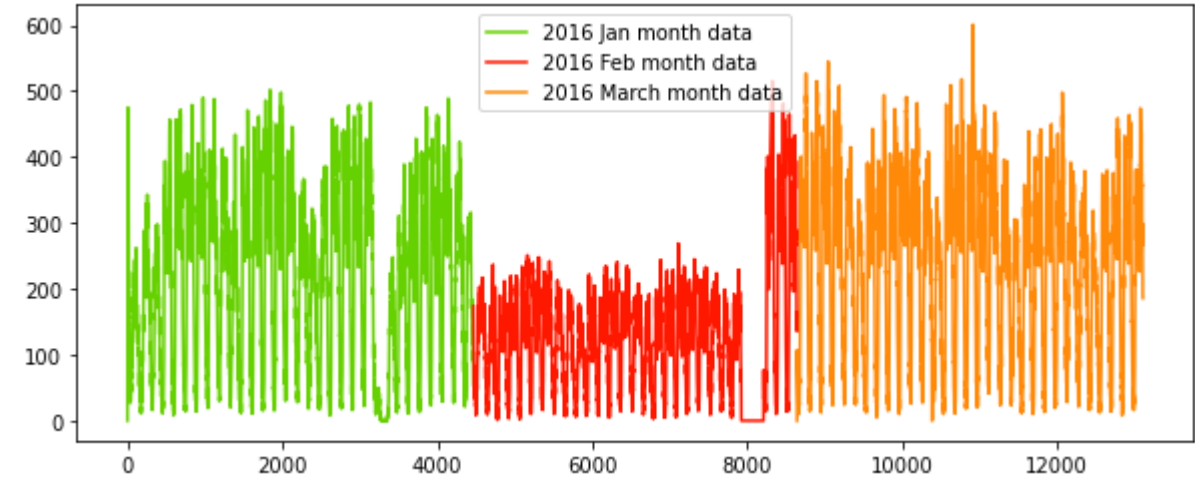


----- Month data for cluster 12 -----

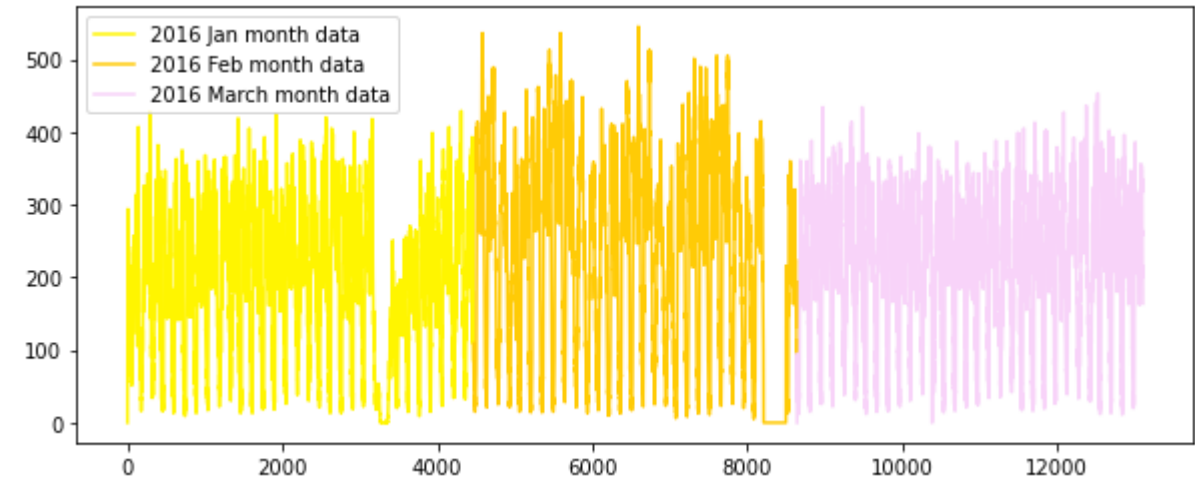




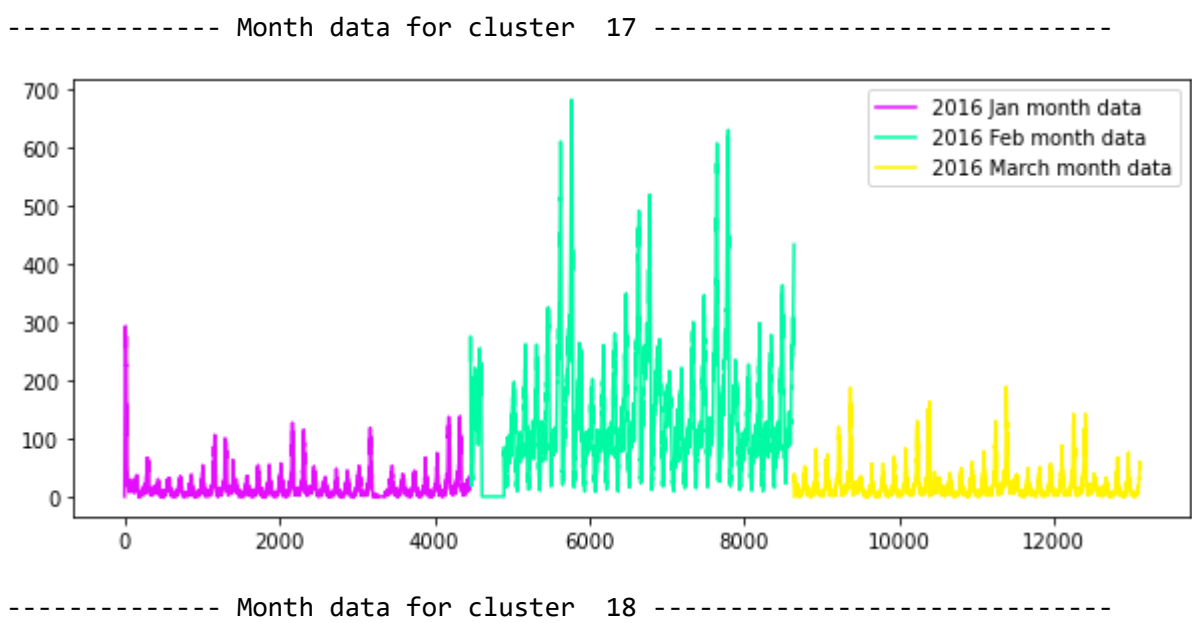
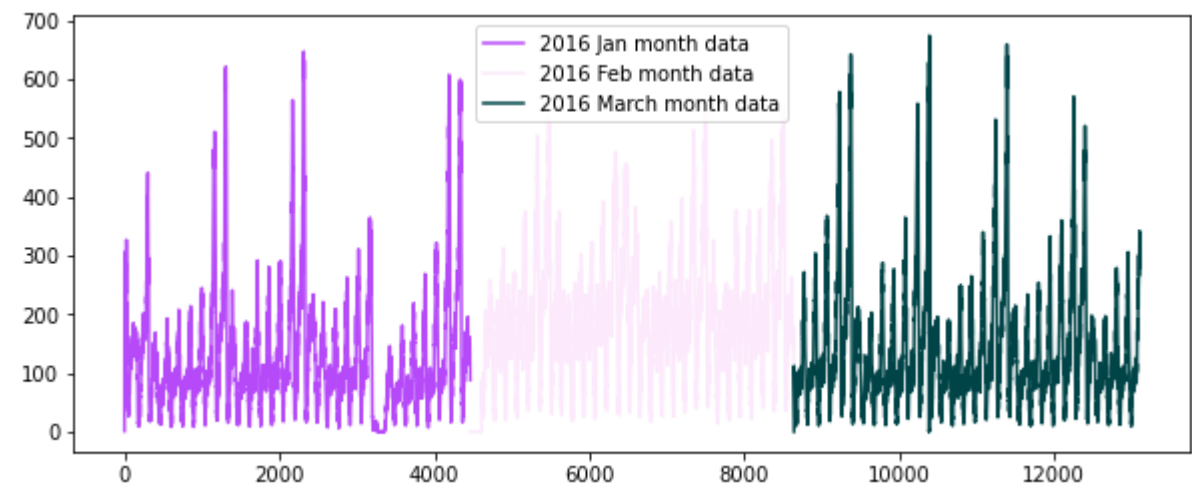
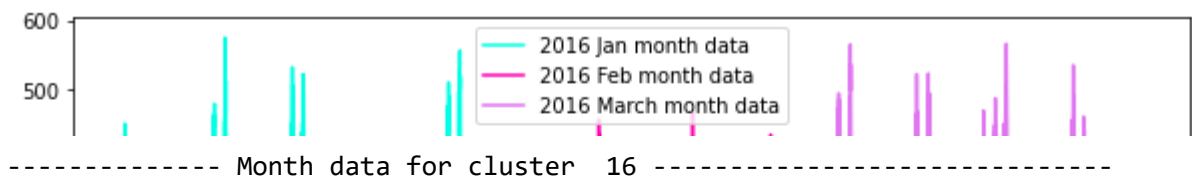
----- Month data for cluster 13 -----

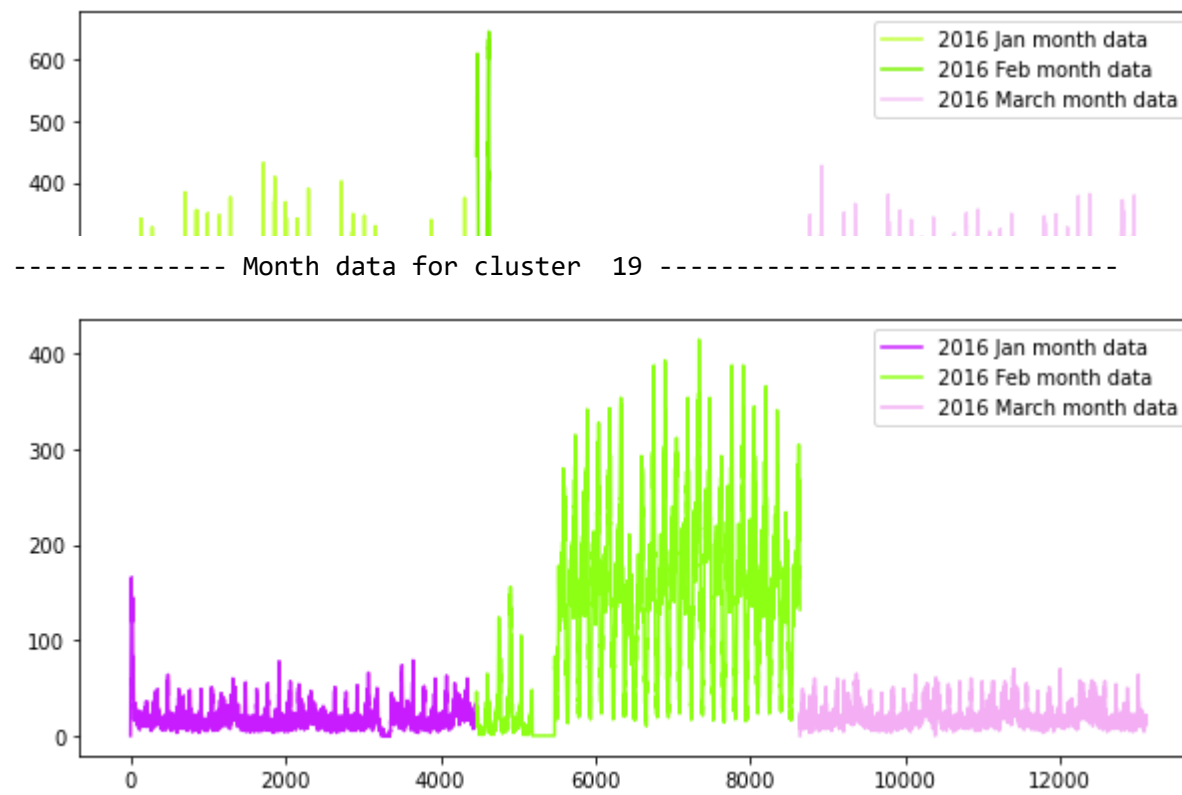


----- Month data for cluster 14 -----

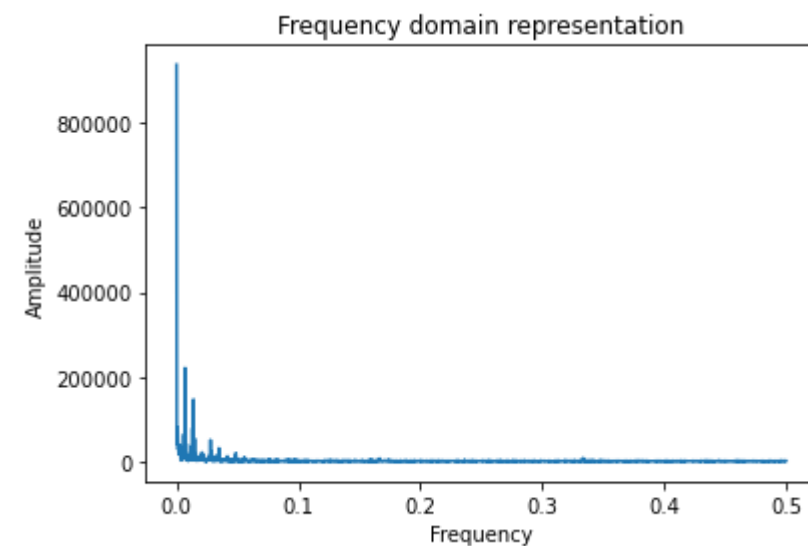


----- Month data for cluster 15 -----





```
In [62]: ▶ 1 ## Plot the amplitude and frequency of n/2 sample frequencies
2 amplitude = np.fft.fft(np.array(jan_2016_smooth)[0:4460])
3 frequency = np.fft.fftfreq(4460,1)
4 n=len(frequency)
5 plt.figure()
6 plt.plot(frequency[:int(n/2)],np.abs(amplitude)[:int(n/2)])
7 plt.title('Frequency domain representation')
8 plt.xlabel('Frequency')
9 plt.ylabel('Amplitude')
10 plt.show()
```



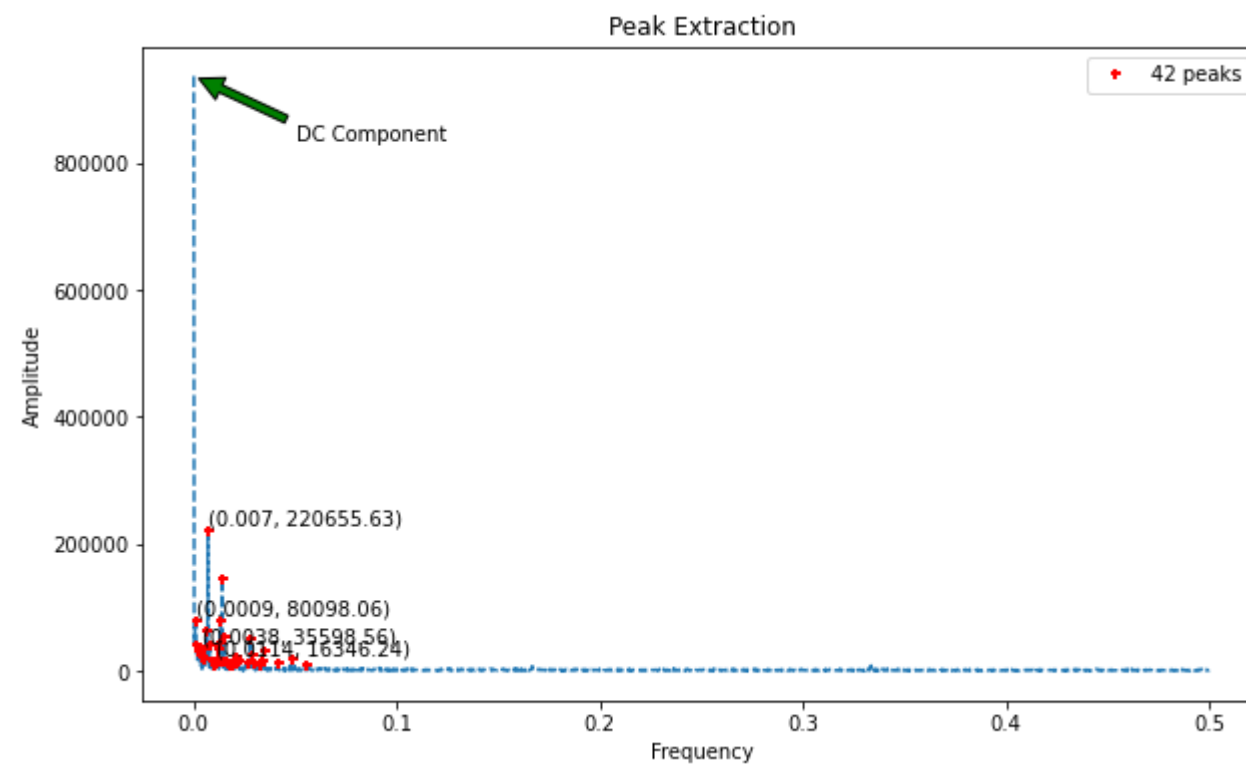
- DC component : A regular wave without a DC component's mean equals zero. With a DC component the mean of the sine wave is not equal to zero. That is a bias is added to the signal.
- $x(t) = D + B \cdot \sin(2 \cdot \pi \cdot f \cdot t)$
- In the above equation D is DC component added and we will not consider its amplitude and frequency. Frequency and amplitude we will consider from the second peak

```

In [63]: 1 import peakutils
2 from peakutils.plot import plot as pp
3 def detect_peaks(y,thres):
4     index = peakutils.indexes(np.abs(y),thres,min_dist=1,thres_abs=True)
5     return index
6
7 thres = 10000
8 index = detect_peaks(amplitude[:int(n/2)],thres)
9 plt.figure(figsize=(10,6))
10 pp(frequency[:int(n/2)],np.abs(amplitude[:int(n/2)]),index)
11
12 plt.annotate('DC Component', xy = (frequency[:int(n/2)][0],np.abs(amplitude[:int(n/2)])[0]),
13             xytext = (frequency[:int(n/2)][0]+0.05,np.abs(amplitude[:int(n/2)])[0]-100000),
14             arrowprops = dict(facecolor = 'green',
15                               shrink = 0.05),)
16 cnt=0
17 for i,j in zip(np.round(frequency[index][:20],4),np.round(np.abs(amplitude)[index][:20],2)):
18     if cnt%5==0:
19         plt.annotate((i,j),xy=(i,j),xytext=(i+0.000029,j+10000))
20     cnt+=1
21 plt.title('Peak Extraction')
22 plt.xlabel('Frequency')
23 plt.ylabel('Amplitude')

```

Out[63]: Text(0, 0.5, 'Amplitude')



- There is a total of 42 peaks above the threshold of 10000
- Annotated and checked some peak's frequency and amplitude

## 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

```
In [64]: 1  ## ratio feature tells us how much is one compared to the other!!!
2  ## It will tell us how many pickups happen in 2016 compqred to 2015
3  ## Ratio feature: rt = pt(2016) / pt(2015)
4  ## ratio between yearly patterns of two connsequive years(2015,2016)
5  ##
6
7  ratio_jan = pd.DataFrame()
8  ratio_jan['prediction_year_pickups'] = jan_2016_smooth
9  ratio_jan['given_year_pickups'] = jan_2015_smooth
10 ratio_jan['Ratios'] = ratio_jan['prediction_year_pickups']*1.0 / ratio_jan['given_year_pickups']*1.0
```

```
In [65]: 1  ## look into the ratio feature
2  ratio_jan.head(5)
```

Out[65]:

	prediction_year_pickups	given_year_pickups	Ratios
0	1	84	0.011905
1	168	84	2.000000
2	363	340	1.067647
3	370	432	0.856481
4	361	514	0.702335

```
In [66]: 1  print('Total time bins in the month of jan : ',len(ratio_jan))
```

Total time bins in the month of jan : 89280

```
In [67]: 1  print('Number of times more pickups happened in year 2015 than in year 2016: ',len(ratio_jan[ratio_jan['Ratios']<1]))
```

Number of times more pickups happened in year 2015 than in year 2016: 55595

```
In [68]: 1  print('Number of times more pickups happened in year 2016 than in year 2015: ',len(ratio_jan[ratio_jan['Ratios']>1]))
```

Number of times more pickups happened in year 2016 than in year 2015: 31148

```
In [69]: 1  print('Number of times same number of pickups happened in year 2016 and 2015: ',len(ratio_jan[ratio_jan['Ratios']==1]))
```

Number of times same number of pickups happened in year 2016 and 2015: 2537

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

- MA values calculation:  $P_{t(2016)} = (P_{t2015} * R_t)$

```

In [70]: 1 def simple_moving_average_ratios(ratios):
2         # here predicted ratio values for i is the average of the ratios of previous three values
3         #  $rt = (rt-1 + rt-2 + rt-3)/3$ 
4         # here predicted ma values is the product of rt and pickups at the given year
5         #  $ma = rt * pt(\text{given})$ 
6         #
7         predicted_ratio = (ratios['Ratios'].values)[0]
8         predicted_ma = []
9         predicted_ratio_values=[]
10        absolute_error = []
11        squared_error = []
12        window_size=3
13        for i in range(0,4464*20):
14            #-----first values for sma and ratio is going to be 0 -----
15            if i%4464==0:
16                predicted_ma.append(0)
17                predicted_ratio_values.append(0)
18                absolute_error.append(0)
19                squared_error.append(0)
20                continue
21
22            predicted_ratio_values.append(predicted_ratio)
23            sma_value = int((ratios['given_year_pickups'].values)[i])*predicted_ratio
24            predicted_ma.append(sma_value)
25            err = abs(sma_value - ratios['prediction_year_pickups'][i])
26            absolute_error.append(err)
27            squared_error.append(math.pow(err,2))
28            if i+1 >= window_size:
29                predicted_ratio = sum((ratios['Ratios'].values)[(i+1)-window_size:i+1])/window_size
30            else:
31                predicted_ratio = sum((ratios['Ratios'].values)[0:i+1])/(i+1)
32
33            ratios['SMA_ratios_predictions'] = predicted_ma
34            ratios['SMA_ratios_absolute_error'] = absolute_error
35            mape_error = (sum(absolute_error)/len(absolute_error)) / (sum(ratios['prediction_year_pickups'].values)/len(ratios['prediction_year_pickups'].values))
36            mse_error = sum(squared_error)/len(squared_error)
37            return ratios,mape_error,mse_error
38

```

```

In [71]: 1 #here, if we calculate absolute percentage error by this formulae:
2         # "error = (abs(int(predicted_ratio_values[i] * ratios["Given"].values[i]) - ratios["Prediction"].values[i])) / ratios["Prediction"].values[i]"
3         #then it will lead to divide by zero problem because many of the values in " ratios["Prediction"].values[i]" are zeros.
4         # so we used this method to calculate mean absolute percentage error: "mean of error/mean of real values"

```

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

```

In [72]: 1 def simple_moving_average_pickups(ratios):
2         predicted_value = (ratios['prediction_year_pickups'].values)[0]
3         predicted_ma = []
4         absolute_error = []
5         squared_error = []
6         window_size=1
7         for i in range(0,4464*20):
8             predicted_ma.append(predicted_value)
9             err = abs(predicted_value - ratios['prediction_year_pickups'][i])
10            absolute_error.append(err)
11            squared_error.append(math.pow(err,2))
12
13            if i+1 >= window_size:
14                predicted_value = int(sum(ratios['prediction_year_pickups'][(i+1)-window_size:i+1])/window_size)
15            else:
16                predicted_value = int(sum(ratios['prediction_year_pickups'][0:i+1])/(i+1))
17
18            ratios['SMA_pickups_predictions'] = predicted_ma
19            ratios['SMA_pickups_absolute_error'] = absolute_error
20            mape_error = (sum(absolute_error)/len(absolute_error)) / (sum(ratios['prediction_year_pickups'].values)/len(ratios['prediction_year_pickups'].values))
21            mse_error = sum(squared_error)/len(squared_error)
22
23            return ratios,mape_error,mse_error

```

```

In [73]: 1 x,y,c = simple_moving_average_pickups(ratio_jan)

```

```

In [74]: 1 y*100

```

Out[74]: 10.870779108525399

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



```

In [75]: 1 def weighted_moving_average_ratios(ratios):
2     predicted_ratio = (ratios['Ratios'].values)[0]
3     absolute_error = []
4     squared_error = []
5     predicted_ratio_values = []
6     predicted_values = []
7     window_size, sum_coeff = 5, 0
8     for i in range(window_size, 0, -1):
9         sum_coeff += i
10    for i in range(0, 4464*20) :
11        if i%4464 == 0:
12            predicted_ratio_values.append(0)
13            absolute_error.append(0)
14            squared_error.append(0)
15            predicted_values.append(0)
16            continue
17
18        value = ratios['given_year_pickups'][i]*predicted_ratio
19        predicted_values.append(value)
20        predicted_ratio_values.append(predicted_ratio)
21        err = abs(value - ratios['prediction_year_pickups'][i])
22        absolute_error.append(err)
23        squared_error.append(math.pow(err,2))
24        if (i+1) >= window_size:
25            sum_ = 0
26            for j in range(window_size, 0, -1):
27                sum_ += j*(ratios['Ratios'][i-window_size+j])
28            predicted_ratio = sum_/sum_coeff
29        else:
30            sum_ = 0
31            for j in range(i, 0, -1):
32                sum_ += j*(ratios['Ratios'][j])
33            predicted_ratio = sum_/sum_coeff
34
35        ratios['WMA_ratios_predictions'] = predicted_values
36        ratios['WMA_ratios_absolute_error'] = absolute_error
37
38        mape_error = (sum(absolute_error)/len(absolute_error)) / (sum(ratios['prediction_year_pickups'].values)/len(ratios['prediction_year_pickups'].values))
39        mse_error = sum(squared_error)/len(squared_error)
40
41    return ratios, mape_error, mse_error

```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 5 is optimal for getting the best results using Weighted Moving Averages using previous Ratio values therefore we get  $R_t = (5 * R_{t-1} + 4 * R_{t-2} + 3 * R_{t-3} + 2 * R_{t-4} + R_{t-5})/15$

Weighted Moving Averages using Previous 2016 Values -  $P_t = (N * P_{t-1} + (N - 1) * P_{t-2} + (N - 2) * P_{t-3} \dots 1 * P_{t-n})/(N * (N + 1)/2)$

```

In [76]: 1 def weighted_moving_average_pickups(ratios):
2     predicted_value = (ratios['prediction_year_pickups'].values)[0]
3     absolute_error = []
4     squared_error = []
5     predicted_values = []
6     window_size, sum_coeff = 2, 0
7     for i in range(window_size, 0, -1):
8         sum_coeff += i
9     for i in range(0, 4464*20) :
10         if i%4464 == 0:
11             absolute_error.append(0)
12             squared_error.append(0)
13             predicted_values.append(0)
14             continue
15
16         predicted_values.append(predicted_value)
17         err = abs(predicted_value - ratios['prediction_year_pickups'][i])
18         absolute_error.append(err)
19         squared_error.append(math.pow(err, 2))
20         if (i+1) >= window_size:
21             sum_ = 0
22             for j in range(window_size, 0, -1):
23                 sum_ += j*(ratios['prediction_year_pickups'][i-window_size+j])
24             predicted_value = sum_/sum_coeff
25         else:
26             sum_ = 0
27             for j in range(i, 0, -1):
28                 sum_ += j*(ratios['prediction_year_pickups'][j])
29             predicted_value = sum_/sum_coeff
30
31         ratios['WMA_ratios_predictions'] = predicted_values
32         ratios['WMA_ratios_absolute_error'] = absolute_error
33
34         mape_error = (sum(absolute_error)/len(absolute_error)) / (sum(ratios['prediction_year_pickups'].values)/len(ratios['prediction_year_pickups'].values))
35         mse_error = sum(squared_error)/len(squared_error)
36
37     return ratios, mape_error, mse_error

```

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) ([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 infinitely 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  $\sim 1/(1 - \alpha) = 10$  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 [77]: 1 def exp_moving_average_ratios(ratios):
2         predicted_ratio=(ratios['Ratios'].values)[0]
3         alpha=0.6
4         absolute_error,squared_error=[],[]
5         predicted_values=[]
6         predicted_ratio_values=[]
7         for i in range(0,4464*20):
8             if i%4464==0:
9                 predicted_ratio_values.append(0)
10                predicted_values.append(0)
11                absolute_error.append(0)
12                squared_error.append(0)
13                continue
14
15                predicted_ratio_values.append(predicted_ratio)
16                predicted_values.append(int(((ratios['given_year_pickups'].values)[i])*predicted_ratio))
17                err = abs(((ratios['given_year_pickups'].values[i])*predicted_ratio)-(ratios['prediction_year_pickups'].values)[i]))
18                absolute_error.append(err)
19                squared_error.append(err**2)
20                predicted_ratio = (alpha*predicted_ratio) + (1-alpha)*((ratios['Ratios'].values)[i])
21
22            ratios['EMA_ratios_predictions'] = predicted_values
23            ratios['WMA_ratios_absolute_error'] = absolute_error
24            mape_err = (sum(absolute_error)/len(absolute_error))/(sum(ratios['prediction_year_pickups'].values)/len(ratios['prediction_year_pickups'].values))
25            mse_err = sum(squared_error)/len(squared_error)
26            return ratios,mape_err,mse_err

```

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

```

In [78]: 1 def exp_moving_average_pickups(ratios):
2         predicted_value=(ratios['prediction_year_pickups'].values)[0]
3         alpha=0.3
4         absolute_error,squared_error=[],[]
5         predicted_values=[]
6         for i in range(0,4464*20):
7             if i%4464==0:
8                 predicted_values.append(0)
9                 absolute_error.append(0)
10                squared_error.append(0)
11                continue
12                predicted_values.append(predicted_value)
13                err = abs((predicted_value)-(ratios['prediction_year_pickups'].values)[i]))
14                absolute_error.append(err)
15                squared_error.append(err**2)
16                predicted_value = int((alpha*predicted_value) + (1-alpha)*((ratios['prediction_year_pickups'].values)[i])))
17
18            ratios['EMA_ratios_predictions'] = predicted_values
19            ratios['WMA_ratios_absolute_error'] = absolute_error
20            mape_err = (sum(absolute_error)/len(absolute_error))/(sum(ratios['prediction_year_pickups'].values)/len(ratios['prediction_year_pickups'].values))
21            mse_err = sum(squared_error)/len(squared_error)
22            return ratios,mape_err,mse_err

```

In [79]: ▶

```
1 mean_err=[0]*10
2 median_err=[0]*10
3 r1,mean_err[0],median_err[0] = simple_moving_average_ratios(ratio_jan)
4 r2,mean_err[1],median_err[1] = simple_moving_average_pickups(ratio_jan)
5 r3,mean_err[2],median_err[2] = weighted_moving_average_ratios(ratio_jan)
6 r4,mean_err[3],median_err[3] = weighted_moving_average_pickups(ratio_jan)
7 r5,mean_err[4],median_err[4] = exp_moving_average_ratios(ratio_jan)
8 r6,mean_err[5],median_err[5] = exp_moving_average_pickups(ratio_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 [80]: ▶

```
1 err_table = pd.DataFrame(columns=['Model','MAPE%','MSE'])
2 a1 = ['Simple Moving Average using ratios',round(mean_err[0]*100,2),median_err[0]]
3 a2 = ['Simple Moving Average using 2016 predictions',round(mean_err[1]*100,2),median_err[1]]
4 a3 = ['Weighted Moving Average using ratios',round(mean_err[2]*100,2),median_err[2]]
5 a4 = ['Weighted Moving Average using 2016 predictions',round(mean_err[3]*100,2),median_err[3]]
6 a5 = ['Exponential Moving Average using ratios',round(mean_err[4]*100,2),median_err[4]]
7 a6 = ['Exponential Moving Average using 2016 predictions',round(mean_err[5]*100,2),median_err[5]]
8 err_table.append(pd.DataFrame([a1,a2,a3,a4,a5,a6],columns=['Model','MAPE%','MSE']))
```

Out[80]:

	Model	MAPE%	MSE
0	Simple Moving Average using ratios	13.97	834.442528
1	Simple Moving Average using 2016 predictions	10.87	381.956933
2	Weighted Moving Average using ratios	13.80	834.044410
3	Weighted Moving Average using 2016 predictions	10.60	372.088764
4	Exponential Moving Average using ratios	13.75	846.339660
5	Exponential Moving Average using 2016 predictions	10.60	370.280253

**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

```

In [81]: 1 # first we are creating some initial features
2 # feature 1 : latitude , feature 2 : Longitude, feature 3: day of the week
3 # longitude : it is the cluster longitude value
4 # We have jan_2016 data which has sublists of size 20 where each list corresponds to a clusternumber
5 # and element of each cluster is a pickup values at 10-minute intervals
6 # every cluster has 4464 values which refers to number of pickups at every 10 mins in the month of jan.
7 ## jan_2016 = [[81,23,45,33,56,.....4464 values],.....20 lists]
8 ## here 81,23,45 are pickup values at timestamp 1 ,2,3,4...
9 ## region_cum is a list of size 20 which contains sublists each having jan_2016-march2016 pickup values
10 ## Total number of values for region_cum would be [[4464+4176+4464],[4464+4176+4464],[4464+4176+4464]...20 lists]
11 ## where 4464 is the total number of 10 minute interval in the month of jan,march and 4176 is the total number of
12 ## 10 minute interval in the month of feb
13
14 # we are going to consider values from the 5th timestamp as we need first 5 timestamp to predict hence we will omit it
15 # in our features
16
17 ## building feature 1 : latitude
18 # latitude : it is the cluster latitude value , we will repeat this value till 13099 (4464+4176+4464 - 5) for each cluster
19 lat = []
20
21 ## building feature 2 : longitude
22 # longitude : it is the cluster longitude value , we will repeat this value till 13099 (4464+4176+4464 - 5) for each cluster
23 lon = []
24
25 ## building feature 3 : week_day
26 # weekday : it is the day of the week coded from 0-6(sun-sat) , we will repeat this value till 13099 (4464+4176+4464 - 5) for each cluster
27 week_day = []
28
29 ## output will contain the pickup values 13099 for each cluster
30
31 out=[]
32
33 # featue4 : tsne_feat
34 # its an numpy arrayhttp://localhost:8888/notebooks/Documents/appleidai/taxi_demand/test.ipynb#y, of shape (523960, 5)
35 # each row corresponds to an entry in out data
36 # for the first row we will have [f0,f1,f2,f3,f4] fi=number of pickups happened in i+1th 10min intravel(bin)
37 # the second row will have [f1,f2,f3,f4,f5]
38 # the third row will have [f2,f3,f4,f5,f6]
39 # and so on...
40
41 number_of_timestamps = 5
42 tsne_feat = [0]*5
43
44 for i in range(20):
45     lat.append([kmeans.cluster_centers_[i][0]]*13099)
46     lon.append([kmeans.cluster_centers_[i][1]]*13099)
47     ## jan 1 2016 is a friday hence we start with the code 5
48     week_day.append([int(((int(k/144))%7+5)%7) for k in range(5,4464+4176+4464)])
49     tsne_feat = np.vstack((tsne_feat,[regions_cum[i][r:r+number_of_timestamps] for r in range(0,len(regions_cum[i])- number_of_timestamps)]))
50     out.append(regions_cum[i][5:])
51
52 tsne_feat = tsne_feat[1:]

```

```

In [82]: 1 ## sanity check : if all the features is of size
2 len(lat[0])*len(lat) == tsne_feat.shape[0] == len(week_day)*len(week_day[0]) == 20*13099 == len(out)*len(out[0])

```

Out[82]: True

## Adding exponential moving average

- from the baseline models we said the exponential weighted moving average gives us the best error
- we will try to add the same exponential weighted moving average at  $t$  as a feature to our data
- exponential weighted moving average  $\Rightarrow p'(t) = \alpha p'(t-1) + (1-\alpha)P(t-1)$

**1. cluster center latitude**

**2. cluster center longitude**

**3. day of the week**

**4.  $f_{t-1}$ : number of pickups that are happened previous  $t-1$ th 10min intravel**

**5.  $f_{t-2}$ : number of pickups that are happened previous  $t-2$ th 10min intravel**

**6.  $f_{t-3}$ : number of pickups that are happened previous  $t-3$ th 10min intravel**

**7.  $f_{t-4}$ : number of pickups that are happened previous  $t-4$ th 10min intravel**

**8.  $f_{t-5}$ : number of pickups that are happened previous  $t-5$ th 10min intravel**

```
In [83]: ▶ 1 alpha = 0.3
2 predicted_values = []
3 predict_list = []
4 tsne_flat_exp_avg = []
5 for clstr in range(0,20):
6     for i in range(0,13104):
7         if i == 0:
8             predicted_value = regions_cum[clstr][0]
9             predicted_values.append(0)
10            continue
11            predicted_values.append(predicted_value)
12            predicted_value = int((alpha*predicted_value) + (1-alpha)*(regions_cum[clstr][i]))
13            predict_list.append(predicted_values[5:])
14            predicted_values=[]
```

```
In [84]: ▶ 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)))
```

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

```
In [85]: ▶ 1 # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our training data
2 train_features = [tsne_feat[i*13099:(13099*i+9169)] for i in range(0,20)]
3 # temp = [0]*(12955 - 9068)
4 test_features = [tsne_feat[(13099*(i))+9169:13099*(i+1)] for i in range(0,20)]
```

```
In [86]: 1 print("Number of data clusters",len(train_features), "Number of data points in trian data", len(train_features[0]), "Each data point contains", len(train_features[0][0]),"fe
2 print("Number of data clusters",len(train_features), "Number of data points in test data", len(test_features[0]), "Each data point contains", len(test_features[0][0]),"featu
```

Number of data clusters 20 Number of data points in trian data 9169 Each data point contains 5 features  
 Number of data clusters 20 Number of data points in test data 3930 Each data point contains 5 features

```
In [87]: 1 # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our training data
2 train_flat_lat = [i[:9169] for i in lat]
3 train_flat_lon = [i[:9169] for i in lon]
4 train_flat_weekday = [i[:9169] for i in week_day]
5 train_flat_output = [i[:9169] for i in out]
6 train_flat_exp_avg = [i[:9169] for i in predict_list]
```

```
In [88]: 1 # extracting the rest of the timestamp values i.e 30% of 12956 (total timestamps) for our test data
2 test_flat_lat = [i[9169:] for i in lat]
3 test_flat_lon = [i[9169:] for i in lon]
4 test_flat_weekday = [i[9169:] for i in week_day]
5 test_flat_output = [i[9169:] for i in out]
6 test_flat_exp_avg = [i[9169:] for i in predict_list]
```

```
In [89]: 1 # the above contains values in the form of list of lists (i.e. list of values of each region), here we make all of them in one list
2 train_new_features = []
3 for i in range(0,20):
4     train_new_features.extend(train_features[i])
5 test_new_features = []
6 for i in range(0,20):
7     test_new_features.extend(test_features[i])
```

```
In [90]: 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 train_lat = sum(train_flat_lat, [])
7 train_lon = sum(train_flat_lon, [])
8 train_weekday = sum(train_flat_weekday, [])
9 train_output = sum(train_flat_output, [])
10 train_exp_avg = sum(train_flat_exp_avg, [])
```

```
In [91]: 1 len(test_flat_exp_avg[0])
```

Out[91]: 3930



```
In [92]: 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 test_lat = sum(test_flat_lat, [])
7 test_lon = sum(test_flat_lon, [])
8 test_weekday = sum(test_flat_weekday, [])
9 test_output = sum(test_flat_output, [])
10 test_exp_avg = sum(test_flat_exp_avg,[])
```

```
In [93]: 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'] = train_lat
5 df_train['lon'] = train_lon
6 df_train['weekday'] = train_weekday
7 df_train['exp_avg'] = train_exp_avg
8
9 print(df_train.shape)
```

(183380, 9)

```
In [94]: 1 # Preparing the data frame for our train data
2 df_test = pd.DataFrame(data=test_new_features, columns=columns)
3 df_test['lat'] = test_lat
4 df_test['lon'] = test_lon
5 df_test['weekday'] = test_weekday
6 df_test['exp_avg'] = test_exp_avg
7 print(df_test.shape)
```

(78600, 9)

```
In [95]: 1 df_train.head(5)
```

Out[95]:

	ft_5	ft_4	ft_3	ft_2	ft_1	lat	lon	weekday	exp_avg
0	1	168	363	370	361	40.746118	-73.979062	5	356
1	168	363	370	361	413	40.746118	-73.979062	5	395
2	363	370	361	413	434	40.746118	-73.979062	5	422
3	370	361	413	434	459	40.746118	-73.979062	5	447
4	361	413	434	459	445	40.746118	-73.979062	5	445

## Adding fourier features

In [96]:

```

1  ## we will add top 5 frequency and amplitude values
2  ## for every region take top 5 frequency and amplitude
3  ## there are 20 clusters
4  final = []
5  for i in range(0,20):
6      jan_data = regions_cum[i][0:4464]
7      feb_data = regions_cum[i][4464:4464+4176]
8      mar_data = regions_cum[i][4464+4176:4464+4176+4464]
9
10     ## creating fourier transforms
11     jan_amp = np.fft.fft(jan_data,4464)
12     feb_amp = np.fft.fft(feb_data,4176)
13     mar_amp = np.fft.fft(mar_data,4464)
14
15     jan_freq = np.fft.fftfreq(4464,1)
16     feb_freq = np.fft.fftfreq(4176,1)
17     mar_freq = np.fft.fftfreq(4464,1)
18
19     jan_ffftamp = sorted(jan_amp,reverse=True)[:5]
20     feb_ffftamp = sorted(feb_amp,reverse=True)[:5]
21     mar_ffftamp = sorted(mar_amp,reverse=True)[:5]
22
23     jan_fftfreq = sorted(jan_freq,reverse=True)[:5]
24     feb_fftfreq = sorted(feb_freq,reverse=True)[:5]
25     mar_fftfreq = sorted(mar_freq,reverse=True)[:5]
26
27
28     m,n,o = [0]*5,[0]*5,[0]*5
29     a,b,c = [0]*5,[0]*5,[0]*5
30
31
32     for i in range(5):
33         m[i] = [jan_fftfreq[i]] * 4464
34         n[i] = [feb_fftfreq[i]] * 4176
35         o[i] = [mar_fftfreq[i]] * 4464
36
37         a[i] = [jan_ffftamp[i]] * 4464
38         b[i] = [feb_ffftamp[i]] * 4176
39         c[i] = [mar_ffftamp[i]] * 4464
40
41
42     jan_ffftamp = np.array(a).T
43     feb_ffftamp = np.array(b).T
44     mar_ffftamp = np.array(c).T
45
46     jan_fftfreq = np.array(m).T
47     feb_fftfreq = np.array(n).T
48     mar_fftfreq = np.array(o).T
49
50
51     jan = np.hstack((jan_ffftamp,jan_fftfreq))
52     feb = np.hstack((feb_ffftamp,feb_fftfreq))
53     mar = np.hstack((mar_ffftamp,mar_fftfreq))
54
55     all_ = np.vstack((jan , feb))
56     all_ = np.vstack((all_ , mar))
57
58     dt = pd.DataFrame(data = all_ ,columns=['A1' , 'A2' , 'A3' , 'A4' , 'A5' , 'F1' , 'F2' , 'F3' , 'F4' , 'F5'])
59     dt = dt.astype(np.float)
60     final.append(dt)

```

```
In [97]: 1 ## Lets concat the fourier features
2 fourier_feat = final[0]
3 for i in range(1,len(final)):
4     fourier_feat = pd.concat([fourier_feat,final[i]],ignore_index=True)
5 print("Shape of fourier transformed features for all points - ", fourier_feat.shape)
6 fourier_feat = fourier_feat.astype(np.float)
7 fourier_feat.tail(3)
```

Shape of fourier transformed features for all points - (262080, 10)

Out[97]:

	A1	A2	A3	A4	A5	F1	F2	F3	F4	F5
262077	91289.0	9746.506136	9746.506136	6436.649191	6436.649191	0.499776	0.499552	0.499328	0.499104	0.49888
262078	91289.0	9746.506136	9746.506136	6436.649191	6436.649191	0.499776	0.499552	0.499328	0.499104	0.49888
262079	91289.0	9746.506136	9746.506136	6436.649191	6436.649191	0.499776	0.499552	0.499328	0.499104	0.49888

Merging the fourier features

```
In [98]: 1 fourier_feat_train = pd.DataFrame(columns=['A1','A2','A3','A4','A5','F1','F2','F3','F4','F5'])
2 fourier_feat_test = pd.DataFrame(columns=['A1','A2','A3','A4','A5','F1','F2','F3','F4','F5'])
3
4 for i in range(20):
5     fourier_feat_train = fourier_feat_train.append(fourier_feat[i*13099 : 13099*i + 9169])
6
7 fourier_feat_train.reset_index(inplace=True)
8
9 for i in range(20):
10     fourier_feat_test = fourier_feat_test.append(fourier_feat[13099*i + 9169 : 13099*(i+1)])
11
12 fourier_feat_test.reset_index(inplace=True)
```

```
In [99]: 1 df_train = pd.concat([df_train,fourier_feat_train],axis=1)
```

```
In [100]: 1 df_train.drop(['index'],axis=1,inplace=True)
```

```
In [101]: 1 df_train.head(5)
```

Out[101]:

	ft_5	ft_4	ft_3	ft_2	ft_1	lat	lon	weekday	exp_avg	A1	A2	A3	A4	A5	F1	F2	F3	F4	F5
0	1	168	363	370	361	40.746118	-73.979062	5	356	936982.0	64630.623606	64630.623606	60044.740121	60044.740121	0.499776	0.499552	0.499328	0.499104	0.49888
1	168	363	370	361	413	40.746118	-73.979062	5	395	936982.0	64630.623606	64630.623606	60044.740121	60044.740121	0.499776	0.499552	0.499328	0.499104	0.49888
2	363	370	361	413	434	40.746118	-73.979062	5	422	936982.0	64630.623606	64630.623606	60044.740121	60044.740121	0.499776	0.499552	0.499328	0.499104	0.49888
3	370	361	413	434	459	40.746118	-73.979062	5	447	936982.0	64630.623606	64630.623606	60044.740121	60044.740121	0.499776	0.499552	0.499328	0.499104	0.49888
4	361	413	434	459	445	40.746118	-73.979062	5	445	936982.0	64630.623606	64630.623606	60044.740121	60044.740121	0.499776	0.499552	0.499328	0.499104	0.49888

```
In [102]: 1 df_test = pd.concat([df_test,fourier_feat_test],axis=1)
2 df_test.drop(['index'],axis=1,inplace=True)
```

In [103]:

1 df\_test.head(5)

Out[103]:

	ft_5	ft_4	ft_3	ft_2	ft_1	lat	lon	weekday	exp_avg	A1	A2	A3	A4	A5	F1	F2	F3	F4	F5
0	250	223	217	229	253	40.746118	-73.979062	5	244	1069666.0	32438.27824	32438.27824	26396.260518	26396.260518	0.499776	0.499552	0.499328	0.499104	0.49888
1	223	217	229	253	271	40.746118	-73.979062	5	262	1069666.0	32438.27824	32438.27824	26396.260518	26396.260518	0.499776	0.499552	0.499328	0.499104	0.49888
2	217	229	253	271	317	40.746118	-73.979062	5	300	1069666.0	32438.27824	32438.27824	26396.260518	26396.260518	0.499776	0.499552	0.499328	0.499104	0.49888
3	229	253	271	317	343	40.746118	-73.979062	5	330	1069666.0	32438.27824	32438.27824	26396.260518	26396.260518	0.499776	0.499552	0.499328	0.499104	0.49888
4	253	271	317	343	351	40.746118	-73.979062	5	344	1069666.0	32438.27824	32438.27824	26396.260518	26396.260518	0.499776	0.499552	0.499328	0.499104	0.49888

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/>)

In [104]:

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

In [105]:

1  
2 def initial\_seasonal\_components(series, slen):  
3 seasonals = {}  
4 season\_averages = []  
5 n\_seasons = int(len(series)/slen)  
6 # compute season averages  
7 for j in range(n\_seasons):  
8 season\_averages.append(sum(series[slen\*j:slen\*j+slen])/float(slen))  
9 # compute initial values  
10 for i in range(slen):  
11 sum\_of\_vals\_over\_avg = 0.0  
12 for j in range(n\_seasons):  
13 sum\_of\_vals\_over\_avg += series[slen\*j+i]-season\_averages[j]  
14 seasonals[i] = sum\_of\_vals\_over\_avg/n\_seasons  
15 return seasonals

```
In [141]: 1 def triple_exponential_smoothing(series,slen, alpha, beta, gamma, n_preds):
2     result = []
3     # n_preds is the number of predictions to be made
4     abs_err = []
5     seasonals = initial_seasonal_components(series, slen)
6     for i in range(len(series)+n_preds):
7         if i == 0: # initial values
8             smooth = series[0]
9             abs_err.append(0)
10            trend = initial_trend(series, slen)
11            result.append(series[0])
12            continue
13        if i >= len(series): # we are forecasting
14            m = i - len(series) + 1
15            result.append((smooth + m*trend) + seasonals[i%slen])
16        else:
17            val = series[i]
18            last_smooth, smooth = smooth, alpha*(val-seasonals[i%slen]) + (1-alpha)*(smooth+trend)
19            trend = beta * (smooth-last_smooth) + (1-beta)*trend
20            seasonals[i%slen] = gamma*(val-smooth) + (1-gamma)*seasonals[i%slen]
21            final_val = smooth+trend+seasonals[i%slen]
22            abs_err.append(abs(series[i] - final_val))
23            result.append(final_val)
24
25    return result,abs_err
```

## Hyperparameter Tuning Alpha,Beta and Gamma values

In [152]:

```

1 season_len = 24
2 predict_list_2 = []
3 mape = []
4 ## Lets manually fine tune alpha beta and gamma
5 alpha = [0.25,0.2,0.15,0.1]
6 beta = [0.1,0.15,0.20,0.25]
7 gamma = [0.1,0.15,0.20,0.08,0.05]
8 fine_tune = []
9 regions_sum = []
10 abs_sum = []
11 length = len(regions_cum[0]) * 20
12 for a in alpha :
13     for b in beta:
14         for g in gamma:
15             for r in range(0,20):
16                 predict_values_2,abs_val = triple_exponential_smoothing(regions_cum[r][0:13104], season_len, a, b, g, 0)
17                 predict_list_2.append(predict_values_2[5:])
18                 regions_sum.append(sum(regions_cum[r][0:13104]))
19                 abs_sum.append(sum(abs_val))
20             mape_err = (sum(abs_sum)/length)/(sum(regions_sum)/length)
21             print('Mape% with alpha : {} ,beta:{}, gamma:{} is '.format(a,b,g),mape_err*100)
22             fine_tune.append((a,b,g,mape_err*100))

```

```

Mape% with alpha : 0.25 ,beta:0.1, gamma:0.1 is 9.942204773276476
Mape% with alpha : 0.25 ,beta:0.1, gamma:0.15 is 9.735437880893153
Mape% with alpha : 0.25 ,beta:0.1, gamma:0.2 is 9.525150005361171
Mape% with alpha : 0.25 ,beta:0.1, gamma:0.08 is 9.671787529522046
Mape% with alpha : 0.25 ,beta:0.1, gamma:0.05 is 9.812423702725827
Mape% with alpha : 0.25 ,beta:0.15, gamma:0.1 is 9.780286682901682
Mape% with alpha : 0.25 ,beta:0.15, gamma:0.15 is 9.717639567671236
Mape% with alpha : 0.25 ,beta:0.15, gamma:0.2 is 9.640733292547917
Mape% with alpha : 0.25 ,beta:0.15, gamma:0.08 is 9.653398474042957
Mape% with alpha : 0.25 ,beta:0.15, gamma:0.05 is 9.684941926351097
Mape% with alpha : 0.25 ,beta:0.2, gamma:0.1 is 9.685905129856026
Mape% with alpha : 0.25 ,beta:0.2, gamma:0.15 is 9.711086043550349
Mape% with alpha : 0.25 ,beta:0.2, gamma:0.2 is 9.975406038670775
Mape% with alpha : 0.25 ,beta:0.2, gamma:0.08 is 9.962528393093713
Mape% with alpha : 0.25 ,beta:0.2, gamma:0.05 is 9.954722082587626
Mape% with alpha : 0.25 ,beta:0.25, gamma:0.1 is 13.51369634201155
Mape% with alpha : 0.25 ,beta:0.25, gamma:0.15 is 122.93096699272878
Mape% with alpha : 0.25 ,beta:0.25, gamma:0.2 is 8921.858556640067
Mape% with alpha : 0.25 ,beta:0.25, gamma:0.08 is 8453.858073832873
Mape% with alpha : 0.25 ,beta:0.25, gamma:0.05 is 8031.784623819195
Mape% with alpha : 0.2 ,beta:0.1, gamma:0.1 is 7649.887482144555
Mape% with alpha : 0.2 ,beta:0.1, gamma:0.15 is 7302.684154000607
Mape% with alpha : 0.2 ,beta:0.1, gamma:0.2 is 6985.6486019144895
Mape% with alpha : 0.2 ,beta:0.1, gamma:0.08 is 6695.086363134685
Mape% with alpha : 0.2 ,beta:0.1, gamma:0.05 is 6427.78198855943
Mape% with alpha : 0.2 ,beta:0.15, gamma:0.1 is 6180.999161684495
Mape% with alpha : 0.2 ,beta:0.15, gamma:0.15 is 5952.480771763048
Mape% with alpha : 0.2 ,beta:0.15, gamma:0.2 is 5740.268165478195
Mape% with alpha : 0.2 ,beta:0.15, gamma:0.08 is 5542.727767544533
Mape% with alpha : 0.2 ,beta:0.15, gamma:0.05 is 5358.365550025177
Mape% with alpha : 0.2 ,beta:0.2, gamma:0.1 is 5185.876248311459
Mape% with alpha : 0.2 ,beta:0.2, gamma:0.15 is 5024.161365814673
Mape% with alpha : 0.2 ,beta:0.2, gamma:0.2 is 4872.23312688974
Mape% with alpha : 0.2 ,beta:0.2, gamma:0.08 is 4729.264493707522
Mape% with alpha : 0.2 ,beta:0.2, gamma:0.05 is 4594.47107286817
Mape% with alpha : 0.2 ,beta:0.25, gamma:0.1 is 4467.163961223516
Mape% with alpha : 0.2 ,beta:0.25, gamma:0.15 is 4346.781140817232

```

```

Mape% with alpha : 0.2 ,beta:0.25, gamma:0.2 is 4232.71595606426
Mape% with alpha : 0.2 ,beta:0.25, gamma:0.08 is 4124.4756637853625
Mape% with alpha : 0.2 ,beta:0.25, gamma:0.05 is 4021.6489036947223
Mape% with alpha : 0.15 ,beta:0.1, gamma:0.1 is 3923.928624072198
Mape% with alpha : 0.15 ,beta:0.1, gamma:0.15 is 3830.844703144093
Mape% with alpha : 0.15 ,beta:0.1, gamma:0.2 is 3742.073661563652
Mape% with alpha : 0.15 ,beta:0.1, gamma:0.08 is 3657.376565397876
Mape% with alpha : 0.15 ,beta:0.1, gamma:0.05 is 3576.4533015981015
Mape% with alpha : 0.15 ,beta:0.15, gamma:0.1 is 3499.0154653648437
Mape% with alpha : 0.15 ,beta:0.15, gamma:0.15 is 3424.860208588546
Mape% with alpha : 0.15 ,beta:0.15, gamma:0.2 is 3353.782463148084
Mape% with alpha : 0.15 ,beta:0.15, gamma:0.08 is 3285.6348152085793
Mape% with alpha : 0.15 ,beta:0.15, gamma:0.05 is 3220.2201358229436
Mape% with alpha : 0.15 ,beta:0.2, gamma:0.1 is 3157.34984557462
Mape% with alpha : 0.15 ,beta:0.2, gamma:0.15 is 3096.8884594139654
Mape% with alpha : 0.15 ,beta:0.2, gamma:0.2 is 3038.6995944461814
Mape% with alpha : 0.15 ,beta:0.2, gamma:0.08 is 2982.6869587480587
Mape% with alpha : 0.15 ,beta:0.2, gamma:0.05 is 2928.7161433384736
Mape% with alpha : 0.15 ,beta:0.25, gamma:0.1 is 2876.6611893277855
Mape% with alpha : 0.15 ,beta:0.25, gamma:0.15 is 2826.4257978744226
Mape% with alpha : 0.15 ,beta:0.25, gamma:0.2 is 2777.9160622260642
Mape% with alpha : 0.15 ,beta:0.25, gamma:0.08 is 2731.0658781469833
Mape% with alpha : 0.15 ,beta:0.25, gamma:0.05 is 2685.780404441717
Mape% with alpha : 0.1 ,beta:0.1, gamma:0.1 is 2642.0989767266997
Mape% with alpha : 0.1 ,beta:0.1, gamma:0.15 is 2599.808289826351
Mape% with alpha : 0.1 ,beta:0.1, gamma:0.2 is 2558.8424394291933
Mape% with alpha : 0.1 ,beta:0.1, gamma:0.08 is 2519.199175698003
Mape% with alpha : 0.1 ,beta:0.1, gamma:0.05 is 2480.786440582889
Mape% with alpha : 0.1 ,beta:0.15, gamma:0.1 is 2443.4968859706023
Mape% with alpha : 0.1 ,beta:0.15, gamma:0.15 is 2407.3065364344084
Mape% with alpha : 0.1 ,beta:0.15, gamma:0.2 is 2372.1671652468226
Mape% with alpha : 0.1 ,beta:0.15, gamma:0.08 is 2338.0785474867216
Mape% with alpha : 0.1 ,beta:0.15, gamma:0.05 is 2304.9720235165987
Mape% with alpha : 0.1 ,beta:0.2, gamma:0.1 is 2272.770736553712
Mape% with alpha : 0.1 ,beta:0.2, gamma:0.15 is 2241.4537892011313
Mape% with alpha : 0.1 ,beta:0.2, gamma:0.2 is 2210.9851305748093
Mape% with alpha : 0.1 ,beta:0.2, gamma:0.08 is 2181.363407877965
Mape% with alpha : 0.1 ,beta:0.2, gamma:0.05 is 2152.5375943185245
Mape% with alpha : 0.1 ,beta:0.25, gamma:0.1 is 2124.4548756474983
Mape% with alpha : 0.1 ,beta:0.25, gamma:0.15 is 2097.094526110388
Mape% with alpha : 0.1 ,beta:0.25, gamma:0.2 is 2070.4291187319714
Mape% with alpha : 0.1 ,beta:0.25, gamma:0.08 is 2044.4550031142467
Mape% with alpha : 0.1 ,beta:0.25, gamma:0.05 is 2019.1344950566602

```

In [204]:

```

1  ## Lets get the parameters which gives minimum mape value
2  from operator import itemgetter
3  best_params = fine_tune[np.argmin(list(map(itemgetter(3),fine_tune)))]
4  print('For triple exponential smoothening the minimum mape value was found to be ',best_params[3])
5  print('The best alpha,beta and gamma parameters for triple exponential smoothening having minimum MAPE value are',best_params[:3])

```

For triple exponential smoothening the minimum mape value was found to be 9.525150005361171  
 The best alpha,beta and gamma parameters for triple exponential smoothening having minimum MAPE value are (0.25, 0.1, 0.2)



```
In [206]: 1 ## triple exponent smoothing with optimal parameters
2 alpha = 0.25
3 beta = 0.1
4 gamma = 0.2
5 season_len = 24
6 predict_list_2 = []
7 predict_values_2 = []
8 for r in range(20):
9     predict_values_2,abs_val = triple_exponential_smoothing(regions_cum[r][0:13104], season_len, 0.25, 0.1, 0.2, 0)
10     predict_list_2.append(predict_values_2[5:])
11
```

```
In [215]: 1 train_flat_exp_avg_2 = [i[:9169] for i in predict_list_2]
2 test_flat_exp_avg_2 = [i[9169:] for i in predict_list_2]
```

```
In [216]: 1 train_exp_avg_2 = sum(train_flat_exp_avg_2,[])
2 test_exp_avg_2 = sum(test_flat_exp_avg_2,[])
```

```
In [217]: 1 df_train['exp_avg_2'] = train_exp_avg_2
2 df_test['exp_avg_2'] = test_exp_avg_2
```

```
In [222]: 1 df_train.head(4)
```

Out[222]:

	ft_5	ft_4	ft_3	ft_2	ft_1	lat	lon	weekday	exp_avg	A1	A2	A3	A4	A5	F1	F2	F3	F4	F5	exp_avg_2
0	1	168	363	370	361	40.746118	-73.979062	5	356	936982.0	64630.623606	64630.623606	60044.740121	60044.740121	0.499776	0.499552	0.499328	0.499104	0.49888	325.786156
1	168	363	370	361	413	40.746118	-73.979062	5	395	936982.0	64630.623606	64630.623606	60044.740121	60044.740121	0.499776	0.499552	0.499328	0.499104	0.49888	368.221926
2	363	370	361	413	434	40.746118	-73.979062	5	422	936982.0	64630.623606	64630.623606	60044.740121	60044.740121	0.499776	0.499552	0.499328	0.499104	0.49888	417.987207
3	370	361	413	434	459	40.746118	-73.979062	5	447	936982.0	64630.623606	64630.623606	60044.740121	60044.740121	0.499776	0.499552	0.499328	0.499104	0.49888	441.461259

Using Linear Regression (hyperparameter tuning :GridSearchCV)

```

In [243]: ► 1 # find more about LinearRegression function here http://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.LinearRegression.html
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-1-2-copy-8/
14 # -----
15
16 from sklearn.linear_model import LinearRegression
17 from sklearn.model_selection import GridSearchCV
18 from sklearn.model_selection import RandomizedSearchCV
19 lr = LinearRegression()
20 params = {'fit_intercept':[True,False], 'normalize':[True,False], 'copy_X':[True,False]}
21 g_s = GridSearchCV(lr,params,cv=3,n_jobs=-1)
22 g_s.fit(df_train,train_output)
23 print(g_s.best_estimator_)

```

LinearRegression(fit\_intercept=False, normalize=True)

```

In [244]: ► 1 lr = LinearRegression(fit_intercept=False,normalize=True)
2 lr.fit(df_train,train_output)
3

```

Out[244]: LinearRegression(fit\_intercept=False, normalize=True)

```

In [245]: ► 1 train_pred = lr.predict(df_train)
2 lr_train_pred = [round(value) for value in train_pred]
3 test_pred = lr.predict(df_test)
4 lr_test_pred = [round(value) for value in test_pred]

```

## Using Random Forest (hyperparameter tuning :RandomSearchCV)

```

In [253]: 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.ensemble.RandomForestRegressor.html
3 # -----
4 # default paramters
5 # sklearn.ensemble.RandomForestRegressor(n_estimators=10, criterion='mse', max_depth=None, min_samples_split=2,
6 # min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
7 # min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False)
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-decision-trees-2/
18 # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
19 # -----
20
21 # import sklearn
22 # sorted(sklearn.metrics.SCORERS.keys())
23 rf = RandomForestRegressor(n_jobs=-1)
24 params = dict()
25 params['min_samples_leaf'] = [1, 2, 4]
26 params['min_samples_split'] = [2, 5, 10]
27 params['n_estimators'] = [int(x) for x in np.linspace(start = 40, stop = 1000, num = 10)]
28 params['max_features'] = ['auto', 'sqrt']
29 rs_rf = RandomizedSearchCV(rf, params, cv=3, scoring='neg_mean_absolute_error', verbose=3)
30 rs_rf.fit(df_train, train_output)
31 print(rs_rf.best_estimator_)

```

Fitting 3 folds for each of 10 candidates, totalling 30 fits

[CV] n\_estimators=360, min\_samples\_split=5, min\_samples\_leaf=4, max\_features=sqrt

[Parallel(n\_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[CV] n\_estimators=360, min\_samples\_split=5, min\_samples\_leaf=4, max\_features=sqrt, score=-8.693, total= 26.5s

[CV] n\_estimators=360, min\_samples\_split=5, min\_samples\_leaf=4, max\_features=sqrt

[Parallel(n\_jobs=1)]: Done 1 out of 1 | elapsed: 26.4s remaining: 0.0s

[CV] n\_estimators=360, min\_samples\_split=5, min\_samples\_leaf=4, max\_features=sqrt, score=-9.958, total= 24.8s

[CV] n\_estimators=360, min\_samples\_split=5, min\_samples\_leaf=4, max\_features=sqrt

[Parallel(n\_jobs=1)]: Done 2 out of 2 | elapsed: 51.2s remaining: 0.0s

[CV] n\_estimators=360, min\_samples\_split=5, min\_samples\_leaf=4, max\_features=sqrt, score=-9.783, total= 25.3s

[CV] n\_estimators=680, min\_samples\_split=10, min\_samples\_leaf=4, max\_features=auto

[CV] n\_estimators=680, min\_samples\_split=10, min\_samples\_leaf=4, max\_features=auto, score=-8.505, total= 2.9min

[CV] n\_estimators=680, min\_samples\_split=10, min\_samples\_leaf=4, max\_features=auto

[CV] n\_estimators=680, min\_samples\_split=10, min\_samples\_leaf=4, max\_features=auto, score=-9.783, total= 2.9min

[CV] n\_estimators=680, min\_samples\_split=10, min\_samples\_leaf=4, max\_features=auto

[CV] n\_estimators=680, min\_samples\_split=10, min\_samples\_leaf=4, max\_features=auto, score=-9.424, total= 2.8min

[CV] n\_estimators=360, min\_samples\_split=2, min\_samples\_leaf=2, max\_features=auto

[CV] n\_estimators=360, min\_samples\_split=2, min\_samples\_leaf=2, max\_features=auto, score=-8.553, total= 1.8min

[CV] n\_estimators=360, min\_samples\_split=2, min\_samples\_leaf=2, max\_features=auto

[CV] n\_estimators=360, min\_samples\_split=2, min\_samples\_leaf=2, max\_features=auto, score=-9.787, total= 1.8min

[CV] n\_estimators=360, min\_samples\_split=2, min\_samples\_leaf=2, max\_features=auto

```
[CV] n_estimators=360, min_samples_split=2, min_samples_leaf=2, max_features=auto, score=-9.453, total= 1.7min
[CV] n_estimators=573, min_samples_split=10, min_samples_leaf=1, max_features=sqrt
[CV] n_estimators=573, min_samples_split=10, min_samples_leaf=1, max_features=sqrt, score=-8.709, total= 42.8s
[CV] n_estimators=573, min_samples_split=10, min_samples_leaf=1, max_features=sqrt
[CV] n_estimators=573, min_samples_split=10, min_samples_leaf=1, max_features=sqrt, score=-9.918, total= 41.4s
[CV] n_estimators=573, min_samples_split=10, min_samples_leaf=1, max_features=sqrt
[CV] n_estimators=573, min_samples_split=10, min_samples_leaf=1, max_features=sqrt, score=-9.764, total= 41.1s
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=2, max_features=sqrt
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=2, max_features=sqrt, score=-8.703, total= 17.8s
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=2, max_features=sqrt
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=2, max_features=sqrt, score=-9.930, total= 17.8s
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=2, max_features=sqrt
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=2, max_features=sqrt, score=-9.768, total= 17.6s
[CV] n_estimators=40, min_samples_split=5, min_samples_leaf=2, max_features=sqrt
[CV] n_estimators=40, min_samples_split=5, min_samples_leaf=2, max_features=sqrt, score=-8.875, total= 3.8s
[CV] n_estimators=40, min_samples_split=5, min_samples_leaf=2, max_features=sqrt
[CV] n_estimators=40, min_samples_split=5, min_samples_leaf=2, max_features=sqrt, score=-9.987, total= 3.6s
[CV] n_estimators=40, min_samples_split=5, min_samples_leaf=2, max_features=sqrt
[CV] n_estimators=40, min_samples_split=5, min_samples_leaf=2, max_features=sqrt, score=-9.855, total= 3.4s
[CV] n_estimators=40, min_samples_split=10, min_samples_leaf=1, max_features=sqrt
[CV] n_estimators=40, min_samples_split=10, min_samples_leaf=1, max_features=sqrt, score=-8.793, total= 3.1s
[CV] n_estimators=40, min_samples_split=10, min_samples_leaf=1, max_features=sqrt
[CV] n_estimators=40, min_samples_split=10, min_samples_leaf=1, max_features=sqrt, score=-9.989, total= 3.6s
[CV] n_estimators=40, min_samples_split=10, min_samples_leaf=1, max_features=sqrt
[CV] n_estimators=40, min_samples_split=10, min_samples_leaf=1, max_features=sqrt, score=-9.939, total= 3.1s
[CV] n_estimators=253, min_samples_split=2, min_samples_leaf=1, max_features=auto
[CV] n_estimators=253, min_samples_split=2, min_samples_leaf=1, max_features=auto, score=-8.577, total= 1.5min
[CV] n_estimators=253, min_samples_split=2, min_samples_leaf=1, max_features=auto
[CV] n_estimators=253, min_samples_split=2, min_samples_leaf=1, max_features=auto, score=-9.813, total= 1.4min
[CV] n_estimators=253, min_samples_split=2, min_samples_leaf=1, max_features=auto
[CV] n_estimators=253, min_samples_split=2, min_samples_leaf=1, max_features=auto, score=-9.484, total= 1.4min
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=1, max_features=auto
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=1, max_features=auto, score=-8.556, total= 1.2min
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=1, max_features=auto
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=1, max_features=auto, score=-9.802, total= 1.2min
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=1, max_features=auto
[CV] n_estimators=253, min_samples_split=10, min_samples_leaf=1, max_features=auto, score=-9.456, total= 1.1min
[CV] n_estimators=786, min_samples_split=5, min_samples_leaf=1, max_features=sqrt
[CV] n_estimators=786, min_samples_split=5, min_samples_leaf=1, max_features=sqrt, score=-8.678, total= 1.1min
[CV] n_estimators=786, min_samples_split=5, min_samples_leaf=1, max_features=sqrt
[CV] n_estimators=786, min_samples_split=5, min_samples_leaf=1, max_features=sqrt, score=-9.893, total= 1.1min
[CV] n_estimators=786, min_samples_split=5, min_samples_leaf=1, max_features=sqrt
[CV] n_estimators=786, min_samples_split=5, min_samples_leaf=1, max_features=sqrt, score=-9.753, total= 1.1min
```

[Parallel(n\_jobs=1)]: Done 30 out of 30 | elapsed: 29.7min finished

```
RandomForestRegressor(min_samples_leaf=4, min_samples_split=10,
                      n_estimators=680, n_jobs=-1)
```

```
In [254]: ► 1 rf = RandomForestRegressor(min_samples_leaf=4, min_samples_split=10,
2                                     n_estimators=680, n_jobs=-1)
3 rf.fit(df_train, train_output)
4
```

```
Out[254]: RandomForestRegressor(min_samples_leaf=4, min_samples_split=10,
                               n_estimators=680, n_jobs=-1)
```

```
In [263]: ▶ 1 train_pred = rf.predict(df_train)
           2 rf_train_pred = [round(value) for value in train_pred]
           3 test_pred = rf.predict(df_test)
           4 rf_test_pred = [round(value) for value in test_pred]
```

## Using Xgboost (hyperparameter tuning :RandomSearchCV)

```

In [259]: 1 # Training a hyper-parameter tuned Xg-Boost regressor on our train data
2
3 # find more about XGBRegressor function here http://xgboost.readthedocs.io/en/latest/python/python_api.html?#module-xgboost.sklearn
4 # -----
5 # default paramters
6 # xgboost.XGBRegressor(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True, objective='reg:linear',
7 # booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1, max_delta_step=0, subsample=1, colsample_bytree=1,
8 # colsample_bylevel=1, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, base_score=0.5, random_state=0, seed=None,
9 # missing=None, **kwargs)
10
11 # some of methods of RandomForestRegressor()
12 # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbose=True, xgb_model=None)
13 # get_params([deep]) Get parameters for this estimator.
14 # predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is not thread safe.
15 # get_score(importance_type='weight') -> get the feature importance
16 # -----
17 # video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-using-decision-trees-2/
18 # video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
19 # -----
20 xg = xgb.XGBRegressor(n_thread=8)
21 params = dict()
22 params['learning_rate'] = [0.1,0.01,0.001,0.2]
23 params['max_depth'] = [3,4,5,10]
24 params['min_child_weight'] = [3,5,10]
25 params['gamma'] = [0]
26 params['subsample'] = [0.8,0.5,0.6]
27 params['reg_alpha'] = [200]
28 params['reg_lambda'] = [200]
29 params['colsample_bytree'] = [0.8,0.5,0.6]
30 params['n_estimators'] = [int(x) for x in np.linspace(start = 40, stop = 1000, num = 10)]
31
32 rs_xgb = RandomizedSearchCV(xg,params,cv=2,verbose=3,scoring='neg_mean_absolute_error')
33 rs_xgb.fit(df_train,train_output)

```

Fitting 2 folds for each of 10 candidates, totalling 20 fits

[CV] subsample=0.8, reg\_lambda=200, reg\_alpha=200, n\_estimators=893, min\_child\_weight=10, max\_depth=10, learning\_rate=0.01, gamma=0, colsample\_bytree=0.5

[21:57:20] WARNING: src/objective/regression\_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

[Parallel(n\_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.

[CV] subsample=0.8, reg\_lambda=200, reg\_alpha=200, n\_estimators=893, min\_child\_weight=10, max\_depth=10, learning\_rate=0.01, gamma=0, colsample\_bytree=0.5, score=-8.790, total=1.3min

[CV] subsample=0.8, reg\_lambda=200, reg\_alpha=200, n\_estimators=893, min\_child\_weight=10, max\_depth=10, learning\_rate=0.01, gamma=0, colsample\_bytree=0.5

[21:58:38] WARNING: src/objective/regression\_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

[Parallel(n\_jobs=1)]: Done 1 out of 1 | elapsed: 1.3min remaining: 0.0s

[CV] subsample=0.8, reg\_lambda=200, reg\_alpha=200, n\_estimators=893, min\_child\_weight=10, max\_depth=10, learning\_rate=0.01, gamma=0, colsample\_bytree=0.5, score=-10.188, total=1.3min

[CV] subsample=0.6, reg\_lambda=200, reg\_alpha=200, n\_estimators=146, min\_child\_weight=10, max\_depth=3, learning\_rate=0.2, gamma=0, colsample\_bytree=0.5

[21:59:56] WARNING: src/objective/regression\_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

[Parallel(n\_jobs=1)]: Done 2 out of 2 | elapsed: 2.6min remaining: 0.0s

[CV] subsample=0.6, reg\_lambda=200, reg\_alpha=200, n\_estimators=146, min\_child\_weight=10, max\_depth=3, learning\_rate=0.2, gamma=0, colsample\_bytree=0.5, score=-9.285, total=5.0s

[CV] subsample=0.6, reg\_lambda=200, reg\_alpha=200, n\_estimators=146, min\_child\_weight=10, max\_depth=3, learning\_rate=0.2, gamma=0, colsample\_bytree=0.5

[22:00:01] WARNING: src/objective/regression\_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

[CV] subsample=0.6, reg\_lambda=200, reg\_alpha=200, n\_estimators=146, min\_child\_weight=10, max\_depth=3, learning\_rate=0.2, gamma=0, colsample\_bytree=0.5, score=-10.758, total=4.9s

```
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=146, min_child_weight=5, max_depth=10, learning_rate=0.2, gamma=0, colsample_bytree=0.5
[22:00:06] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=146, min_child_weight=5, max_depth=10, learning_rate=0.2, gamma=0, colsample_bytree=0.5, score=-8.900, total
= 15.2s
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=146, min_child_weight=5, max_depth=10, learning_rate=0.2, gamma=0, colsample_bytree=0.5
[22:00:21] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=146, min_child_weight=5, max_depth=10, learning_rate=0.2, gamma=0, colsample_bytree=0.5, score=-10.020, tota
l= 14.8s
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=1000, min_child_weight=10, max_depth=3, learning_rate=0.001, gamma=0, colsample_bytree=0.8
[22:00:36] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=1000, min_child_weight=10, max_depth=3, learning_rate=0.001, gamma=0, colsample_bytree=0.8, score=-46.243, t
otal= 38.0s
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=1000, min_child_weight=10, max_depth=3, learning_rate=0.001, gamma=0, colsample_bytree=0.8
[22:01:14] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=1000, min_child_weight=10, max_depth=3, learning_rate=0.001, gamma=0, colsample_bytree=0.8, score=-49.349, t
otal= 37.3s
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=893, min_child_weight=3, max_depth=10, learning_rate=0.001, gamma=0, colsample_bytree=0.8
[22:01:51] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=893, min_child_weight=3, max_depth=10, learning_rate=0.001, gamma=0, colsample_bytree=0.8, score=-51.023, to
tal= 1.4min
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=893, min_child_weight=3, max_depth=10, learning_rate=0.001, gamma=0, colsample_bytree=0.8
[22:03:13] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=893, min_child_weight=3, max_depth=10, learning_rate=0.001, gamma=0, colsample_bytree=0.8, score=-54.902, to
tal= 1.4min
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=573, min_child_weight=10, max_depth=3, learning_rate=0.01, gamma=0, colsample_bytree=0.8
[22:04:36] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=573, min_child_weight=10, max_depth=3, learning_rate=0.01, gamma=0, colsample_bytree=0.8, score=-9.511, tota
l= 21.6s
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=573, min_child_weight=10, max_depth=3, learning_rate=0.01, gamma=0, colsample_bytree=0.8
[22:04:58] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=573, min_child_weight=10, max_depth=3, learning_rate=0.01, gamma=0, colsample_bytree=0.8, score=-10.631, tot
al= 21.6s
[CV] subsample=0.6, reg_lambda=200, reg_alpha=200, n_estimators=40, min_child_weight=10, max_depth=4, learning_rate=0.1, gamma=0, colsample_bytree=0.6
[22:05:19] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.6, reg_lambda=200, reg_alpha=200, n_estimators=40, min_child_weight=10, max_depth=4, learning_rate=0.1, gamma=0, colsample_bytree=0.6, score=-10.182, total
= 2.1s
[CV] subsample=0.6, reg_lambda=200, reg_alpha=200, n_estimators=40, min_child_weight=10, max_depth=4, learning_rate=0.1, gamma=0, colsample_bytree=0.6
[22:05:21] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.6, reg_lambda=200, reg_alpha=200, n_estimators=40, min_child_weight=10, max_depth=4, learning_rate=0.1, gamma=0, colsample_bytree=0.6, score=-11.508, total
= 2.0s
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=786, min_child_weight=10, max_depth=10, learning_rate=0.1, gamma=0, colsample_bytree=0.6
[22:05:23] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=786, min_child_weight=10, max_depth=10, learning_rate=0.1, gamma=0, colsample_bytree=0.6, score=-8.747, tota
l= 1.4min
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=786, min_child_weight=10, max_depth=10, learning_rate=0.1, gamma=0, colsample_bytree=0.6
[22:06:48] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=786, min_child_weight=10, max_depth=10, learning_rate=0.1, gamma=0, colsample_bytree=0.6, score=-9.979, tota
l= 1.4min
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=253, min_child_weight=3, max_depth=5, learning_rate=0.1, gamma=0, colsample_bytree=0.5
[22:08:10] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=253, min_child_weight=3, max_depth=5, learning_rate=0.1, gamma=0, colsample_bytree=0.5, score=-8.840, total=
13.2s
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=253, min_child_weight=3, max_depth=5, learning_rate=0.1, gamma=0, colsample_bytree=0.5
[22:08:23] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.5, reg_lambda=200, reg_alpha=200, n_estimators=253, min_child_weight=3, max_depth=5, learning_rate=0.1, gamma=0, colsample_bytree=0.5, score=-10.303, total
= 12.9s
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=360, min_child_weight=5, max_depth=3, learning_rate=0.2, gamma=0, colsample_bytree=0.5
[22:08:36] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=360, min_child_weight=5, max_depth=3, learning_rate=0.2, gamma=0, colsample_bytree=0.5, score=-8.873, total=
10.5s
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=360, min_child_weight=5, max_depth=3, learning_rate=0.2, gamma=0, colsample_bytree=0.5
```



```
[22:08:47] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
[CV] subsample=0.8, reg_lambda=200, reg_alpha=200, n_estimators=360, min_child_weight=5, max_depth=3, learning_rate=0.2, gamma=0, colsample_bytree=0.5, score=-10.162, total
= 10.3s
[22:08:57] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.

[Parallel(n_jobs=1)]: Done 20 out of 20 | elapsed: 11.6min finished
```

```
Out[259]: RandomizedSearchCV(cv=2, estimator=XGBRegressor(n_thread=8),
                          param_distributions={'colsample_bytree': [0.8, 0.5, 0.6],
                                              'gamma': [0],
                                              'learning_rate': [0.1, 0.01, 0.001,
                                                                0.2],
                                              'max_depth': [3, 4, 5, 10],
                                              'min_child_weight': [3, 5, 10],
                                              'n_estimators': [40, 146, 253, 360, 466,
                                                             573, 680, 786, 893,
                                                             1000],
                                              'reg_alpha': [200], 'reg_lambda': [200],
                                              'subsample': [0.8, 0.5, 0.6]},
                          scoring='neg_mean_absolute_error', verbose=3)
```

```
In [261]: 1 print(rs_xgb.best_estimator_)
```

```
XGBRegressor(colsample_bytree=0.6, max_depth=10, min_child_weight=10,
             n_estimators=786, n_thread=8, reg_alpha=200, reg_lambda=200,
             subsample=0.8)
```

```
In [262]: 1 xg = xgb.XGBRFRegressor(colsample_bytree=0.6, max_depth=10, min_child_weight=10,
2           n_estimators=786, n_thread=8, reg_alpha=200, reg_lambda=200,
3           subsample=0.8)
4 xg.fit(df_train, train_output)
```

```
[22:33:34] WARNING: src/objective/regression_obj.cu:152: reg:linear is now deprecated in favor of reg:squarederror.
```

```
Out[262]: XGBRFRegressor(colsample_bytree=0.6, max_depth=10, min_child_weight=10,
                        n_estimators=786, n_thread=8, reg_alpha=200, reg_lambda=200)
```

```
In [264]: 1 train_pred = xg.predict(df_train)
2 xg_train_pred = [round(value) for value in train_pred]
3 test_pred = xg.predict(df_test)
4 xg_test_pred = [round(value) for value in test_pred]
```

```
In [268]: 1 #feature importances
2 #xg.booster().get_score(importance_type='weight')
```

```
In [290]: 1 print(df_train.columns)
2 print(xg.feature_importances_)
```

```
Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'lat', 'lon', 'weekday',
      'exp_avg', 'A1', 'A2', 'A3', 'A4', 'A5', 'F1', 'F2', 'F3', 'F4', 'F5',
      'exp_avg_2'],
      dtype='object')
[3.7699025e-02 4.8193872e-02 9.0655975e-02 1.3181305e-01 1.8441990e-01
 3.6479064e-04 1.9213855e-06 1.8774941e-07 2.2546878e-01 1.4097838e-03
 3.0596522e-05 2.6543017e-07 3.8726605e-03 4.7213740e-07 3.0464292e-07
 3.4868592e-07 3.2071367e-07 3.9105612e-07 3.9132601e-07 2.7606696e-01]
```

Calculating the error metric values for various models

In [272]: ▶

```
1 train_mape=[]
2 test_mape=[]
3
4 train_mape.append((mean_absolute_error(train_output,df_train['ft_1'].values))/(sum(train_output)/len(train_output)))
5 train_mape.append((mean_absolute_error(train_output,df_train['exp_avg'].values))/(sum(train_output)/len(train_output)))
6 train_mape.append((mean_absolute_error(train_output,rf_train_pred))/(sum(train_output)/len(train_output)))
7 train_mape.append((mean_absolute_error(train_output, xg_train_pred))/(sum(train_output)/len(train_output)))
8 train_mape.append((mean_absolute_error(train_output, lr_train_pred))/(sum(train_output)/len(train_output)))
9
10 test_mape.append((mean_absolute_error(test_output, df_test['ft_1'].values))/(sum(test_output)/len(test_output)))
11 test_mape.append((mean_absolute_error(test_output, df_test['exp_avg'].values))/(sum(test_output)/len(test_output)))
12 test_mape.append((mean_absolute_error(test_output, rf_test_pred))/(sum(test_output)/len(test_output)))
13 test_mape.append((mean_absolute_error(test_output, xg_test_pred))/(sum(test_output)/len(test_output)))
14 test_mape.append((mean_absolute_error(test_output, lr_test_pred))/(sum(test_output)/len(test_output)))
```

In [273]: ▶

```
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[3]," Test: ",test_mape[3])
6 print ("Random Forest Regression - Train: ",train_mape[2]," Test: ",test_mape[2])
```

Error Metric Matrix (Tree Based Regression Methods) - MAPE

-----

Baseline Model -	Train: 0.1066779557743926	Test: 0.10382670907415832
Exponential Averages Forecasting -	Train: 0.10438967160385766	Test: 0.10147760418060392
Linear Regression -	Train: 0.1025037303227335	Test: 0.10003708106992266
Random Forest Regression -	Train: 0.04122459403921813	Test: 0.06612825845123457

Error Metric Matrix

In [274]: ▶

```
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: 0.1066779557743926	Test: 0.10382670907415832
Exponential Averages Forecasting -	Train: 0.10438967160385766	Test: 0.10147760418060392
Linear Regression -	Train: 0.08005411702162843	Test: 0.07367970365726387
Random Forest Regression -	Train: 0.04122459403921813	Test: 0.06612825845123457
XgBoost Regression -	Train: 0.1025037303227335	Test: 0.10003708106992266

-----

Results

```
In [288]: 1 from prettytable import PrettyTable as pt
2 x = pt()
3 x.field_names = ["Model Name","Train MAPE","Test MAPE" ]
4 x.add_row(['Baseline Model',train_mape[0]*100,test_mape[0]*100])
5 x.add_row(['Exponential Averages Forecasting',train_mape[1]*100,test_mape[1]*100])
6 x.add_row(['Linear Regression',train_mape[2]*100,test_mape[2]*100])
7 x.add_row(['Random Forest Regression',train_mape[3]*100,test_mape[3]*100])
8 x.add_row(['XgBoost Regression',train_mape[4]*100,test_mape[4]*100])
9 print(x)
```

Model Name	Train MAPE	Test MAPE
Baseline Model	10.667795577439259	10.382670907415832
Exponential Averages Forecasting	10.438967160385767	10.147760418060392
Linear Regression	4.122459403921813	6.612825845123457
Random Forest Regression	10.25037303227335	10.003708106992267
XgBoost Regression	8.005411702162842	7.367970365726387

Conclusions :

- Linear Regression despite being the simpler model gives the least MAPE among others
- We used Holt Winter's method and FFT features for Linear Regression,Random Forest and XGB models .
- Holt's Winter's forecasteting with mixing parameter values (0.25,0.1,0.2) has a prominent impact on the dataset as it is able to reduce MAPE below 10.

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/)
- <https://blog.ytotech.com/2015/11/01/findpeaks-in-python/> (https://blog.ytotech.com/2015/11/01/findpeaks-in-python/)
- <https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fft.html> (https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fft.html)