

# Taxi demand prediction in New York City

□

In [1]:

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

import pandas as pd#pandas to create small dataframes

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

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

import time #Convert to unix time

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

# matplotlib: used to plot graphs
import matplotlib
# matplotlib.use('nbagg') : matplotlib uses this protocol which makes plots more user interactive
like zoom in and zoom out
matplotlib.use('nbagg')
import matplotlib.pyplot as plt
import seaborn as sns#Plots
from matplotlib import rcParams#Size of plots

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

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

# download mingwin: https://mingw-w64.org/doku.php/download/mingw-builds
# install it in your system and keep the path, mingw_path ='installed path'
mingw_path = 'C:\\Program Files\\mingw-w64\\x86_64-5.3.0-posix-seh-rt_v4-rev0\\mingw64\\bin'
os.environ['PATH'] = mingw_path + ';' + os.environ['PATH']

# to install xgboost: pip3 install xgboost
# if it didnt happen check install_xgboost.JPG
import xgboost as xgb

# to install sklearn: pip install -U scikit-learn
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error
from sklearn.metrics import mean_absolute_error
import warnings
warnings.filterwarnings("ignore")
```

```
C:\Users\sanjeev\Anaconda3\lib\site-packages\dask\dataframe\utils.py:15: FutureWarning:
pandas.util.testing is deprecated. Use the functions in the public API at pandas.testing instead.
import pandas.util.testing as tm
```

## Data Information

## Data Information

Get 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 (FHV)**

FHV transportation is accessed by a pre-arrangement with a dispatcher or limo company. These FHV's 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

In [2]:

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

```
# However unlike Pandas, operations on dask.dataframes don't trigger immediate computation,
# instead they add key-value pairs to an underlying Dask graph. Recall that in the diagram below,
# circles are operations and rectangles are results.
```

```
# to see the visulaization you need to install graphviz
# pip3 install graphviz if this doesnt work please check the install_graphviz.jpg in the drive
import os
os.environ["PATH"] += os.pathsep + 'C:/Users/sanjeev/Anaconda3/Lib/site-packages/graphviz-
2.38/release/bin'
month.visualize()
```

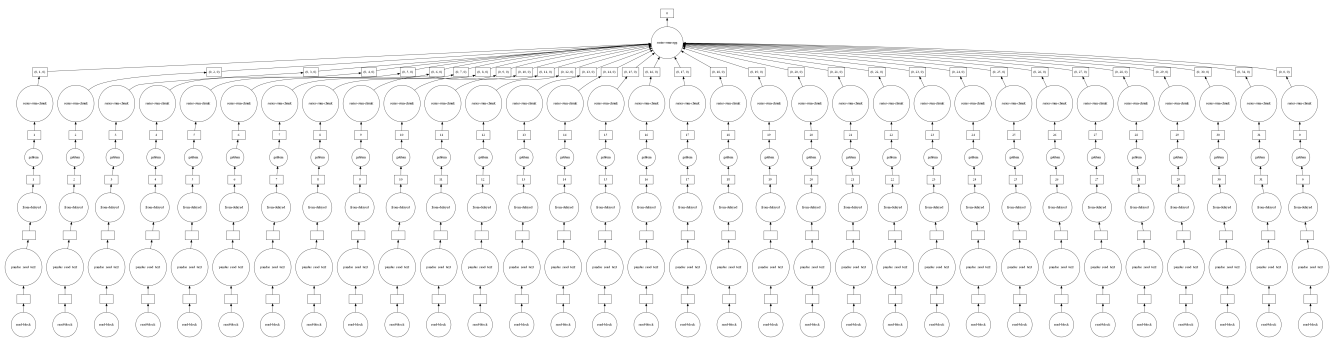
Out[3]:



In [4]:

```
month.fare_amount.sum().visualize()
```

Out[4]:



## Features in the dataset:

Field Name	Description
VendorID	A code indicating the TPEP provider that provided the record. 1. Creative Mobile Technologies 2. VeriFone Inc.
tpep_pickup_datetime	The date and time when the meter was engaged.
tpep_dropoff_datetime	The date and time when the meter was disengaged.
Passenger_count	The number of passengers in the vehicle. This is a driver-entered value.
Trip_distance	The elapsed trip distance in miles reported by the taximeter.

Pickup_longitude	Longitude where the meter was engaged.
Pickup_latitude	Latitude where the meter was engaged.
RateCodeID	The final rate code in effect at the end of the trip. 1. Standard rate 2. JFK 3. Newark 4. Nassau or Westchester 5. Negotiated fare 6. Group ride
Store_and_fwd_flag	This flag indicates whether the trip record was held in vehicle memory before sending to the vendor aka "store and forward," because the vehicle did not have a connection to the server. Y= store and forward trip N= not a store and forward trip
Dropoff_longitude	Longitude where the meter was disengaged.
Dropoff_latitude	Latitude where the meter was disengaged.
Payment_type	A numeric code signifying how the passenger paid for the trip. 1. Credit card 2. Cash 3. No charge 4. Dispute 5. Unknown 6. Voided trip
Fare_amount	The time-and-distance fare calculated by the meter.
Extra	Miscellaneous extras and surcharges. Currently, this only includes. the 0.50 and 1 rush hour and overnight charges.
MTA_tax	0.50 MTA tax that is automatically triggered based on the metered rate in use.
Improvement_surcharge	0.30 improvement surcharge assessed trips at the flag drop. the improvement surcharge began being levied in 2015.
Tip_amount	Tip amount – This field is automatically populated for credit card tips. Cash tips are not included.
Tolls_amount	Total amount of all tolls paid in trip.
Total_amount	The total amount charged to passengers. Does not include cash tips.

## ML Problem Formulation

### Time-series forecasting and Regression

- To find number of pickups, given location coordinates (latitude and longitude) and time, in the query region 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]:

```
#table below shows few datapoints along with all our features
month.head(5)
```

Out [5]:

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

## 1. Pickup Latitude and Pickup Longitude

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

In [6]:

```
# Plotting pickup coordinates which are outside the bounding box of New-York
# we will collect all the points outside the bounding box of newyork city to outlier_locations
outlier_locations = month[((month.pickup_longitude <= -74.15) | (month.pickup_latitude <= 40.5774) |
\
                        (month.pickup_longitude >= -73.7004) | (month.pickup_latitude >= 40.9176)))]

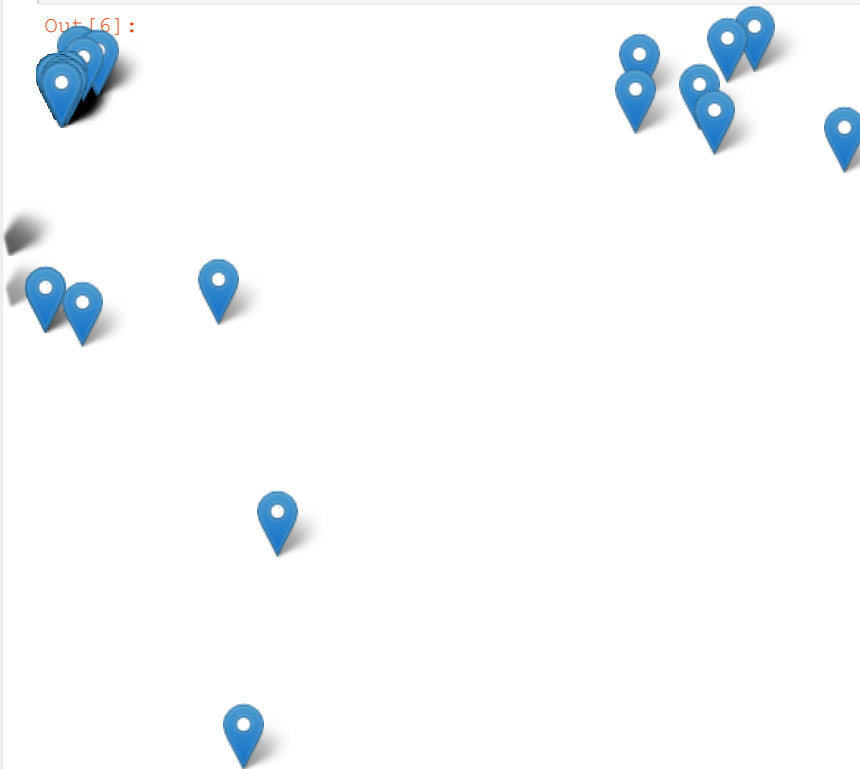
# creating a map with the a base location
# read more about the folium here: http://folium.readthedocs.io/en/latest/quickstart.html

# note: you dont need to remember any of these, you dont need indepth knowledge on these maps and
plots

map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')

# we will spot only first 100 outliers on the map, plotting all the outliers will take more time
sample_locations = outlier_locations.head(10000)
for i,j in sample_locations.iterrows():
    if int(j['pickup_latitude']) != 0:
        folium.Marker(list((j['pickup_latitude'],j['pickup_longitude'])).add_to(map_osm)
map_osm
```

Out[6]:



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

## 2. Dropoff Latitude & Dropoff Longitude

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

In [7]:

```

# Plotting dropoff coordinates which are outside the bounding box of New-York
# we will collect all the points outside the bounding box of newyork city to outlier_locations
outlier_locations = month[(month.dropoff_longitude <= -74.15) | (month.dropoff_latitude <= 40.5774) | \
                           (month.dropoff_longitude >= -73.7004) | (month.dropoff_latitude >= 40.9176))]

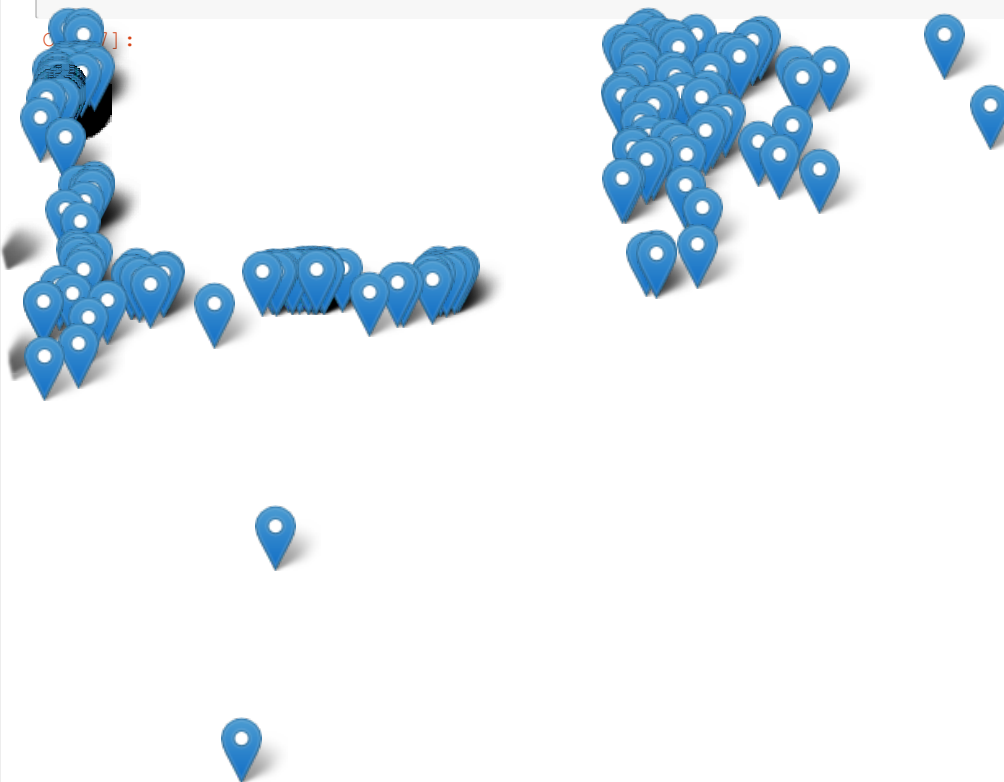
# creating a map with the a base location
# read more about the folium here: http://folium.readthedocs.io/en/latest/quickstart.html

# note: you dont need to remember any of these, you dont need indeepth knowledge on these maps and plots

map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')

# we will spot only first 100 outliers on the map, plotting all the outliers will take more time
sample_locations = outlier_locations.head(10000)
for i,j in sample_locations.iterrows():
    if int(j['pickup_latitude']) != 0:
        folium.Marker(list((j['dropoff_latitude'],j['dropoff_longitude'])).add_to(map_osm)
map_osm

```



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

### 3. Trip Durations:

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

In [8]:

```

#The timestamps are converted to unix so as to get duration(trip-time) & speed also pickup-times i
n unix are used while binning

# in out data we have time in the formate "YYYY-MM-DD HH:MM:SS" we convert thiss sting to python t
ime formate and then into unix time stamp
# https://stackoverflow.com/a/27914405
def convert_to_unix(s):
    return time.mktime(datetime.datetime.strptime(s, "%Y-%m-%d %H:%M:%S").timetuple())

```

```

# we return a data frame which contains the columns
# 1.'passenger_count' : self explanatory
# 2.'trip_distance' : self explanatory
# 3.'pickup_longitude' : self explanatory
# 4.'pickup_latitude' : self explanatory
# 5.'dropoff_longitude' : self explanatory
# 6.'dropoff_latitude' : self explanatory
# 7.'total_amount' : total fair that was paid
# 8.'trip_times' : duration of each trip
# 9.'pickup_times' : pickup time converted into unix time
# 10.'Speed' : velocity of each trip
def return_with_trip_times(month):
    duration = month[['tpep_pickup_datetime', 'tpep_dropoff_datetime']].compute()
    #pickups and dropoffs to unix time
    duration_pickup = [convert_to_unix(x) for x in duration['tpep_pickup_datetime'].values]
    duration_drop = [convert_to_unix(x) for x in duration['tpep_dropoff_datetime'].values]
    #calculate duration of trips
    durations = (np.array(duration_drop) - np.array(duration_pickup))/float(60)

    #append durations of trips and speed in miles/hr to a new dataframe
    new_frame =
month[['passenger_count', 'trip_distance', 'pickup_longitude', 'pickup_latitude', 'dropoff_longitude',
'dropoff_latitude', 'total_amount']].compute()

    new_frame['trip_times'] = durations
    new_frame['pickup_times'] = duration_pickup
    new_frame['Speed'] = 60*(new_frame['trip_distance']/new_frame['trip_times'])

    return new_frame

# print(frame_with_durations.head())
# passenger_count trip_distance pickup_longitude pickup_latitude dropoff_longitude
dropoff_latitude total_amount trip_times pickup_times Speed
# 1 1.59 -73.993896 40.750111 -73.974785 40.750618
17.05 18.050000 1.421329e+09 5.285319
# 1 3.30 -74.001648 40.724243 -73.994415 40.759109
.80 19.833333 1.420902e+09 9.983193
# 1 1.80 -73.963341 40.802788 -73.951820 40.824413
10.80 10.050000 1.420902e+09 10.746269
# 1 0.50 -74.009087 40.713818 -74.004326 40.719986
4.80 1.866667 1.420902e+09 16.071429
# 1 3.00 -73.971176 40.762428 -74.004181 40.742653
6.30 19.316667 1.420902e+09 9.318378
frame_with_durations = return_with_trip_times(month)

```

In [9]:

```

# the skewed box plot shows us the presence of outliers
sns.boxplot(y="trip_times", data =frame_with_durations)
plt.show()

```

In [10]:

```

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

```

```

0 percentile value is -1211.0166666666667
10 percentile value is 3.8333333333333335
20 percentile value is 5.383333333333334
30 percentile value is 6.816666666666666
40 percentile value is 8.3
50 percentile value is 9.95
60 percentile value is 11.866666666666667
70 percentile value is 14.283333333333333
80 percentile value is 17.633333333333333
90 percentile value is 23.45
100 percentile value is 548555.6333333333

```

In [11]:

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

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

In [13]:

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

In [14]:

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

In [15]:

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

In [16]:

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

In [17]:

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



## 4. Speed

In [18]:

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

In [19]:

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

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

```
#calculating speed values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100
for i in np.arange(0.0, 1.0, 0.1):
    var =frame_with_durations_modified["Speed"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
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]:

```
#removing further outliers based on the 99.9th percentile value
frame_with_durations_modified=frame_with_durations[(frame_with_durations.Speed>0) &
(frame_with_durations.Speed<45.31)]
```

In [23]:

```
#avg.speed of cabs in New-York
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]:

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

In [25]:

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

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

```
#calculating trip distance 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
for i in np.arange(0.0, 1.0, 0.1):
    var = frame_with_durations_modified["trip_distance"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
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]:

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

In [29]:

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

## 5. Total Fare

In [30]:

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

In [31]:

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

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

```
#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
for i in np.arange(0.0, 1.0, 0.1):
    var = frame_with_durations_modified["total_amount"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
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]:

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

In [35]:

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

In [36]:

```
#now looking at values not including the last two points we again find a drastic increase at around
```

```

d 1000 fare value
# we plot last 50 values excluding last two values
plt.plot(var[-50:-2])
plt.show()

```

## Remove all outliers/erronous points.

In [37]:

```

#removing all outliers based on our univariate analysis above
def remove_outliers(new_frame):

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

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

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

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

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

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

    new_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times < 720)]
    new_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_distance < 23)]
    new_frame = new_frame[(new_frame.Speed < 45.31) & (new_frame.Speed > 0)]
    new_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_amount >0)]

    print ("Total outliers removed",a - new_frame.shape[0])
    print ("----")
    return new_frame

```

In [38]:

```

print ("Removing outliers in the month of Jan-2015")
print ("----")
frame_with_durations_outliers_removed = remove_outliers(frame_with_durations)
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]:

```

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

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

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

# we need to choose number of clusters so that, there are more number of cluster regions
#that are close to any cluster center
# and make sure that the minimum inter cluster should not be very less
for increment in range(10, 100, 10):
    cluster_centers, cluster_len = find_clusters(increment)
    find_min_distance(cluster_centers, cluster_len)

```

```

On choosing a cluster size of 10
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 2.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 8.0
Min inter-cluster distance = 1.0945442325142543
---
On choosing a cluster size of 20
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 4.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 16.0
Min inter-cluster distance = 0.7131298007387813
---
On choosing a cluster size of 30

```

```

On choosing a cluster size of 30
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 22.0
Min inter-cluster distance = 0.5185088176172206
---
On choosing a cluster size of 40
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 32.0
Min inter-cluster distance = 0.5069768450363973
---
On choosing a cluster size of 50
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 12.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 38.0
Min inter-cluster distance = 0.365363025983595
---
On choosing a cluster size of 60
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 14.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 46.0
Min inter-cluster distance = 0.34704283494187155
---
On choosing a cluster size of 70
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 16.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 54.0
Min inter-cluster distance = 0.30502203163244707
---
On choosing a cluster size of 80
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 18.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 62.0
Min inter-cluster distance = 0.29220324531738534
---
On choosing a cluster size of 90
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 21.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 69.0
Min inter-cluster distance = 0.18257992857034985
---

```

## Inference:

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

In [40]:

```

# if check for the 50 clusters you can observe that there are two clusters with only 0.3 miles apart from each other
# so we choose 40 clusters for solve the further problem

# Getting 40 clusters using the kmeans
kmeans = MiniBatchKMeans(n_clusters=40, batch_size=10000, random_state=0).fit(coords)
frame_with_durations_outliers_removed['pickup_cluster'] =
kmeans.predict(frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']])

```

## Plotting the cluster centers:

In [41]:

```

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

```

Out[41]:

## Plotting the clusters:

In [42]:

```
#Visualising the clusters on a map
def plot_clusters(frame):
    city_long_border = (-74.03, -73.75)
    city_lat_border = (40.63, 40.85)
    fig, ax = plt.subplots(ncols=1, nrows=1)
    ax.scatter(frame.pickup_longitude.values[:100000], frame.pickup_latitude.values[:100000], s=10,
lw=0,
               c=frame.pickup_cluster.values[:100000], cmap='tab20', alpha=0.2)
    ax.set_xlim(city_long_border)
    ax.set_ylim(city_lat_border)
    ax.set_xlabel('Longitude')
    ax.set_ylabel('Latitude')
    plt.show()

plot_clusters(frame_with_durations_outliers_removed)
```

## Time-binning

In [43]:

```
#Refer:https://www.unixtimestamp.com/
# 1420070400 : 2015-01-01 00:00:00
# 1422748800 : 2015-02-01 00:00:00
# 1425168000 : 2015-03-01 00:00:00
# 1427846400 : 2015-04-01 00:00:00
# 1430438400 : 2015-05-01 00:00:00
# 1433116800 : 2015-06-01 00:00:00

# 1451606400 : 2016-01-01 00:00:00
# 1454284800 : 2016-02-01 00:00:00
# 1456790400 : 2016-03-01 00:00:00
# 1459468800 : 2016-04-01 00:00:00
# 1462060800 : 2016-05-01 00:00:00
# 1464739200 : 2016-06-01 00:00:00

def add_pickup_bins(frame, month, year):
    unix_pickup_times=[i for i in frame['pickup_times'].values]
    unix_times = [[1420070400,1422748800,1425168000,1427846400,1430438400,1433116800],\
                   [1451606400,1454284800,1456790400,1459468800,1462060800,1464739200]]

    start_pickup_unix=unix_times[year-2015][month-1]
    # https://www.timeanddate.com/time/zones/est
    # (int((i-start_pickup_unix)/6001.221)) : our unix time is in gmt to us are converting it to est
```



```

# (int((1-start_pickup_unix)/600)+33) : our unix time is in gmt so we are converting it to est
tenminutewise_binned_unix_pickup_times=[(int((i-start_pickup_unix)/600)+33) for i in unix_picku
p_times]
frame['pickup_bins'] = np.array(tenminutewise_binned_unix_pickup_times)
return frame

```

In [44]:

```

# clustering, making pickup bins and grouping by pickup cluster and pickup bins
frame_with_durations_outliers_removed['pickup_cluster'] =
kmeans.predict(frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']])
jan_2015_frame = add_pickup_bins(frame_with_durations_outliers_removed,1,2015)
jan_2015_groupby =
jan_2015_frame[['pickup_cluster','pickup_bins','trip_distance']].groupby(['pickup_cluster','pickup_
bins']).count()

```

In [45]:

```

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

```

Out[45]:

	passenger_count	trip_distance	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	total_amount	trip_times	picku
0	1	1.59	-73.993896	40.750111	-73.974785	40.750618	17.05	18.050000	1.42
1	1	3.30	-74.001648	40.724243	-73.994415	40.759109	17.80	19.833333	1.42
2	1	1.80	-73.963341	40.802788	-73.951820	40.824413	10.80	10.050000	1.42
3	1	0.50	-74.009087	40.713818	-74.004326	40.719986	4.80	1.866667	1.42
4	1	3.00	-73.971176	40.762428	-74.004181	40.742653	16.30	19.316667	1.42

In [46]:

```

# hear the trip_distance represents the number of pickups that are happend in that particular 10mi
n intravel
# this data frame has two indices
# primary index: pickup_cluster (cluster number)
# secondary index : pickup_bins (we devid whole months time into 10min intravels 24*31*60/10 =4464
bins)
jan_2015_groupby.head()

```

Out[46]:

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

In [47]:

```

# upto now we cleaned data and prepared data for the month 2015,
# now do the same operations for months Jan, Feb, March of 2016
# 1. get the dataframe which includes only required colums
# 2. adding trip times, speed, unix time stamp of pickup_time
# 4. remove the outliers based on trip_times, speed, trip_duration, total_amount
# 5. add pickup_cluster to each data point
# 6. add pickup_bin (index of 10min intravel to which that trip belongs to)
# 7. group by data, based on 'pickup_cluster' and 'pickuo_bin'

# Data Preparation for the months of Jan, Feb and March 2016
def datapreparation(month,kmeans,month no,vear no):

```

```

print ("Return with trip times..")

frame_with_durations = return_with_trip_times(month)

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

print ("Estimating clusters..")
frame_with_durations_outliers_removed['pickup_cluster'] =
kmeans.predict(frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']])
#frame_with_durations_outliers_removed_2016['pickup_cluster'] =
kmeans.predict(frame_with_durations_outliers_removed_2016[['pickup_latitude',
'pickup_longitude']])

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

return final_updated_frame,final_groupby_frame

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

jan_2016_frame,jan_2016_groupby = datapreparation(month_jan_2016,kmeans,1,2016)
feb_2016_frame,feb_2016_groupby = datapreparation(month_feb_2016,kmeans,2,2016)
mar_2016_frame,mar_2016_groupby = datapreparation(month_mar_2016,kmeans,3,2016)

```

```

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

```

# Gets the unique bins where pickup values are present for each each region
# for each cluster region we will collect all the indices of 10min intervals in which the pickups

```

```

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

```

In [49]:

```

# for every month we get all indices of 10min intravels in which atleast one pickup got happened

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

#feb
feb_2016_unique = return_unq_pickup_bins(feb_2016_frame)

#march
mar_2016_unique = return_unq_pickup_bins(mar_2016_frame)

```

In [50]:

```

# for each cluster number of 10min intravels with 0 pickups
for i in range(40):
    print("for the ",i,"th cluster number of 10min intavels with zero pickups: ",4464 -
len(set(jan_2015_unique[i])))
    print('-'*60)

```

```

for the 0 th cluster number of 10min intavels with zero pickups: 41
-----
for the 1 th cluster number of 10min intavels with zero pickups: 1986
-----
for the 2 th cluster number of 10min intavels with zero pickups: 30
-----
for the 3 th cluster number of 10min intavels with zero pickups: 355
-----
for the 4 th cluster number of 10min intavels with zero pickups: 38
-----
for the 5 th cluster number of 10min intavels with zero pickups: 154
-----
for the 6 th cluster number of 10min intavels with zero pickups: 35
-----
for the 7 th cluster number of 10min intavels with zero pickups: 34
-----
for the 8 th cluster number of 10min intavels with zero pickups: 118
-----
for the 9 th cluster number of 10min intavels with zero pickups: 41
-----
for the 10 th cluster number of 10min intavels with zero pickups: 26
-----
for the 11 th cluster number of 10min intavels with zero pickups: 45
-----
for the 12 th cluster number of 10min intavels with zero pickups: 43
-----
for the 13 th cluster number of 10min intavels with zero pickups: 29
-----
for the 14 th cluster number of 10min intavels with zero pickups: 27
-----
for the 15 th cluster number of 10min intavels with zero pickups: 32
-----
for the 16 th cluster number of 10min intavels with zero pickups: 41
-----
for the 17 th cluster number of 10min intavels with zero pickups: 59
-----
for the 18 th cluster number of 10min intavels with zero pickups: 1191
-----
for the 19 th cluster number of 10min intavels with zero pickups: 1358
-----
for the 20 th cluster number of 10min intavels with zero pickups: 54

```

```

-----
for the 21 th cluster number of 10min intavels with zero pickups: 30
-----
for the 22 th cluster number of 10min intavels with zero pickups: 30
-----
for the 23 th cluster number of 10min intavels with zero pickups: 164
-----
for the 24 th cluster number of 10min intavels with zero pickups: 36
-----
for the 25 th cluster number of 10min intavels with zero pickups: 42
-----
for the 26 th cluster number of 10min intavels with zero pickups: 32
-----
for the 27 th cluster number of 10min intavels with zero pickups: 215
-----
for the 28 th cluster number of 10min intavels with zero pickups: 37
-----
for the 29 th cluster number of 10min intavels with zero pickups: 42
-----
for the 30 th cluster number of 10min intavels with zero pickups: 1181
-----
for the 31 th cluster number of 10min intavels with zero pickups: 43
-----
for the 32 th cluster number of 10min intavels with zero pickups: 45
-----
for the 33 th cluster number of 10min intavels with zero pickups: 44
-----
for the 34 th cluster number of 10min intavels with zero pickups: 40
-----
for the 35 th cluster number of 10min intavels with zero pickups: 43
-----
for the 36 th cluster number of 10min intavels with zero pickups: 37
-----
for the 37 th cluster number of 10min intavels with zero pickups: 322
-----
for the 38 th cluster number of 10min intavels with zero pickups: 37
-----
for the 39 th cluster number of 10min intavels with zero pickups: 44
-----

```

there are two ways to fill up these values

- Fill the missing value with 0's
- Fill the missing values with the avg values
  - Case 1:(values missing at the start)
    - Ex1:  $\_ \_ \_ x \Rightarrow \text{ceil}(x/4), \text{ceil}(x/4), \text{ceil}(x/4), \text{ceil}(x/4)$
    - Ex2:  $\_ \_ \_ x \Rightarrow \text{ceil}(x/3), \text{ceil}(x/3), \text{ceil}(x/3)$
  - Case 2:(values missing in middle)
    - Ex1:  $x \_ \_ y \Rightarrow \text{ceil}((x+y)/4), \text{ceil}((x+y)/4), \text{ceil}((x+y)/4), \text{ceil}((x+y)/4)$
    - Ex2:  $x \_ \_ \_ y \Rightarrow \text{ceil}((x+y)/5), \text{ceil}((x+y)/5), \text{ceil}((x+y)/5), \text{ceil}((x+y)/5), \text{ceil}((x+y)/5)$
  - Case 3:(values missing at the end)
    - Ex1:  $x \_ \_ \_ \Rightarrow \text{ceil}(x/4), \text{ceil}(x/4), \text{ceil}(x/4), \text{ceil}(x/4)$
    - Ex2:  $x \_ \Rightarrow \text{ceil}(x/2), \text{ceil}(x/2)$

In [51]:

```

# Fills a value of zero for every bin where no pickup data is present
# the count_values: number pickups that are happened in each region for each 10min intravel
# there wont be any value if there are no pickups.
# values: number of unique bins

# for every 10min intravel(pickup_bin) we will check it is there in our unique bin,
# if it is there we will add the count_values[index] to smoothed data
# if not we add 0 to the smoothed data
# we finally return smoothed data
def fill_missing(count_values, values):
    smoothed_regions=[]
    ind=0
    for r in range(0,40):
        smoothed_bins=[]
        for i in range(4464):
            if i in values[r]:
                smoothed_bins.append(count_values[ind])

```

```

        smoothed_value.append(count_values[ind]),
        ind+=1
    else:
        smoothed_bins.append(0)
        smoothed_regions.extend(smoothed_bins)
    return smoothed_regions

```

In [52]:

```

# Fills a value of zero for every bin where no pickup data is present
# the count_values: number pickups that are happened in each region for each 10min intravel
# there wont be any value if there are no pickups.
# values: number of unique bins

# for every 10min intravel(pickup_bin) we will check it is there in our unique bin,
# if it is there we will add the count_values[index] to smoothed data
# if not we add smoothed data (which is calculated based on the methods that are discussed in the
above markdown cell)
# we finally return smoothed data
def smoothing(count_values, values):
    smoothed_regions=[] # stores list of final smoothed values of each region
    ind=0
    repeat=0
    smoothed_value=0
    for r in range(0,40):
        smoothed_bins=[] #stores the final smoothed values
        repeat=0
        for i in range(4464):
            if repeat!=0: # prevents iteration for a value which is already visited/resolved
                repeat-=1
                continue
            if i in values[r]: #checks if the pickup-bin exists
                smoothed_bins.append(count_values[ind]) # appends the value of the pickup bin if it
exists
            else:
                if i!=0:
                    right_hand_limit=0
                    for j in range(i,4464):
                        if j not in values[r]: #searches for the left-limit or the pickup-bin
value which has a pickup value
                            continue
                        else:
                            right_hand_limit=j
                            break
                    if right_hand_limit==0:
                        #Case 1: When we have the last/last few values are found to be missing,hence we
have no right-limit here
                        smoothed_value=count_values[ind-1]*1.0/((4463-i)+2)*1.0

                        for j in range(i,4464):
                            smoothed_bins.append(math.ceil(smoothed_value))
                            smoothed_bins[i-1] = math.ceil(smoothed_value)
                            repeat=(4463-i)
                            ind-=1
                        else:
                            #Case 2: When we have the missing values between two known values
                            smoothed_value=(count_values[ind-1]+count_values[ind])*1.0/((right_hand_lim
t-i)+2)*1.0

                            for j in range(i,right_hand_limit+1):
                                smoothed_bins.append(math.ceil(smoothed_value))
                                smoothed_bins[i-1] = math.ceil(smoothed_value)
                                repeat=(right_hand_limit-i)
                            else:
                                #Case 3: When we have the first/first few values are found to be missing,hence
we have no left-limit here
                                right_hand_limit=0
                                for j in range(i,4464):
                                    if j not in values[r]:
                                        continue
                                    else:
                                        right_hand_limit=j
                                        break
                                smoothed_value=count_values[ind]*1.0/((right_hand_limit-i)+1)*1.0
                                for j in range(i,right_hand_limit+1):
                                    smoothed_bins.append(math.ceil(smoothed_value))
                                    repeat=(right_hand_limit-i)
                                ind+=1

```

```
smoothed_regions.extend(smoothed_bins)
return smoothed_regions
```

In [53]:

```
#Filling Missing values of Jan-2015 with 0
# here in jan_2015_groupby dataframe the trip_distance represents the number of pickups that are h
appened
jan_2015_fill = fill_missing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)

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

In [54]:

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

number of 10min intravels among all the clusters 178560

In [55]:

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

In [56]:

```
# why we choose, these methods and which method is used for which data?

# 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
# 10st 10min intravel, 0 pickups happened in 2nd 10mins intravel, 0 pickups happened in 3rd 10min
intravel
# and 20 pickups happened in 4th 10min intravel.
# in fill_missing method we replace these values like 10, 0, 0, 20
# where as in smoothing method we replace these values as 6,6,6,6,6, if you can check the number o
f pickups
# that are happened in the first 40min are same in both cases, but if you can observe that we look
ing at the future values
# when you are using smoothing we are looking at the future number of pickups which might cause a
data leakage.

# so we use smoothing for jan 2015th data since it acts as our training data
# and we use simple fill_misssing method for 2016th data.
```

In [57]:

```
# Jan-2015 data is smoothed, Jan, Feb & March 2016 data missing values are filled with zero
jan_2015_smooth = smoothing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)
jan_2016_smooth = fill_missing(jan_2016_groupby['trip_distance'].values,jan_2016_unique)
feb_2016_smooth = fill_missing(feb_2016_groupby['trip_distance'].values,feb_2016_unique)
mar_2016_smooth = fill_missing(mar_2016_groupby['trip_distance'].values,mar_2016_unique)

# Making list of all the values of pickup data in every bin for a period of 3 months and storing t
hem region-wise
regions_cum = []

# a =[1,2,3]
# b = [2,3,4]
```

```

# a+b = [1, 2, 3, 2, 3, 4]

# number of 10min indices for jan 2015= 24*31*60/10 = 4464
# number of 10min indices for jan 2016 = 24*31*60/10 = 4464
# number of 10min indices for feb 2016 = 24*29*60/10 = 4176
# number of 10min indices for march 2016 = 24*31*60/10 = 4464
# regions_cum: it will contain 40 lists, each list will contain 4464+4176+4464 values which represents the number of pickups
# that are happened for three months in 2016 data

for i in range(0,40):
    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)])

# print(len(regions_cum))
# 40
# print(len(regions_cum[0]))
# 13104

```

## Time series and Fourier Transforms

In [58]:

```

def uniqueish_color():
    """There're better ways to generate unique colors, but this isn't awful."""
    return plt.cm.gist_ncar(np.random.random())

first_x = list(range(0,4464))
second_x = list(range(4464,8640))
third_x = list(range(8640,13104))
for i in range(40):
    plt.figure(figsize=(10,4))
    plt.plot(first_x,regions_cum[i][:4464], color=uniqueish_color(), label='2016 Jan month data')
    plt.plot(second_x,regions_cum[i][4464:8640], color=uniqueish_color(), label='2016 feb month data')
    plt.plot(third_x,regions_cum[i][8640:], color=uniqueish_color(), label='2016 march month data')
    plt.legend()
    plt.show()

```

In [59]:

```
# getting peaks: https://blog.ytotech.com/2015/11/01/findpeaks-in-python/
# read more about fft function :
https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fft.html
Y = np.fft.fft(np.array(jan_2016_smooth)[0:4460])
# read more about the fftfreq:
https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fftfreq.html
freq = np.fft.fftfreq(4460, 1)
n = len(freq)
plt.figure()
plt.plot( freq[:int(n/2)], np.abs(Y)[:int(n/2)] )
plt.xlabel("Frequency")
plt.ylabel("Amplitude")
plt.show()
```

In [70]:

```
def Get_FFT_Features(data):
    '''Gives top 5 amplitudes and their corresponding frequencies'''
    size = len(data)
    print(size)
    fft_size = 36
    frequencies = []
    amplitudes = []
    for i in range(0,size,fft_size):

        amp = np.fft.fft(np.array(data)[i:i+fft_size])
        freqs = np.fft.fftfreq(fft_size,1)
        top5_indices = np.argsort(amp)[::-1][:len(amp)][:5]
        top5_amps = amp[top5_indices]
        top5_amps_of_freqs = freqs[top5_indices]
        for j in range(0,fft_size):
            amplitudes.append(top5_amps)
            frequencies.append(top5_amps_of_freqs)
    amplitudes = np.asarray(amplitudes)
    frequencies = np.asarray(frequencies)
    print(amplitudes.shape)
    print(frequencies.shape)

    return amplitudes,frequencies
```

In [152]:

```
def normalize_complex_arr(a):
    a_oo = a - a.real.min() - 1j*a.imag.min() # origin offsetted
    return a_oo/np.abs(a_oo).max()
```

In [71]:

```
amps_jan_2016,freq_jan_2016 = Get_FFT_Features(jan_2016_smooth)
```

```
178560
(178560, 5)
(178560, 5)
```

In [75]:



```
features_jan_2016 = np.hstack((amps_jan_2016,freq_jan_2016))
print(features_jan_2016.shape)
```

```
(178560, 10)
```

In [76]:

```
amps_feb_2016,freq_feb_2016 = Get_FFT_Features(feb_2016_smooth)
```

```
178560
(178560, 5)
(178560, 5)
```

In [77]:

```
features_feb_2016 = np.hstack((amps_feb_2016,freq_feb_2016))
print(features_jan_2016.shape)
```

```
(178560, 10)
```

In [78]:

```
amps_mar_2016,freq_mar_2016 = Get_FFT_Features(mar_2016_smooth)
```

```
178560
(178560, 5)
(178560, 5)
```

In [79]:

```
features_mar_2016 = np.hstack((amps_mar_2016,freq_mar_2016))
print(features_mar_2016.shape)
```

```
(178560, 10)
```

In [80]:

```
print("FFT_Jan_2016 : ",features_jan_2016.shape)
print("FFT_Feb_2016 : ",features_feb_2016.shape)
print("FFT_Mar_2016 : ",features_mar_2016.shape)
```

```
FFT_Jan_2016 : (178560, 10)
FFT_Feb_2016 : (178560, 10)
FFT_Mar_2016 : (178560, 10)
```

In [81]:

```
FFT_Features = np.vstack((features_jan_2016,features_feb_2016,features_mar_2016))
```

In [157]:

```
FFT_Features = np.abs(FFT_Features)
```

In [158]:

```
columns = ["amp1","amp2","amp3","amp4","amp5","freq1","freq2","freq3","freq4","freq5"]
FFT_DataSet = pd.DataFrame(data=FFT_Features,columns = columns)
```

In [159]:

```
FFT_DataSet.head()
```

Out[159]:

```
Out[155]:
```

	amp1	amp2	amp3	amp4	amp5	freq1	freq2	freq3	freq4	freq5
0	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.472222	0.416667	0.416667
1	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.472222	0.416667	0.416667
2	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.472222	0.416667	0.416667
3	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.472222	0.416667	0.416667
4	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.472222	0.416667	0.416667

```
In [160]:
```

```
FFT_DataSet.shape
```

```
Out[160]:
```

```
(535680, 10)
```

```
In [84]:
```

```
#Preparing the Dataframe only with x(i) values as jan-2015 data and y(i) values as jan-2016
ratios_jan = pd.DataFrame()
ratios_jan['Given']=jan_2015_smooth
ratios_jan['Prediction']=jan_2016_smooth
ratios_jan['Ratios']=ratios_jan['Prediction']*1.0/ratios_jan['Given']*1.0
```

## Modelling: Baseline Models

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

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

## Simple Moving Averages

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

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

```
In [85]:
```

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

    ratios['MA_R_Predicted'] = predicted_values
    ratios['MA_R_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction']).v
```

```

alues))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err

```

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

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

In [86]:

```

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

    ratios['MA_P_Predicted'] = predicted_values
    ratios['MA_P_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].v
alues))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err

```

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

## Weighted Moving Averages

The Moving Avergaes Model used gave equal importance to all the values in the window used, but we know intuitively that the future is more likely to be similar to the latest values and less similar to the older values. Weighted Averages converts this analogy into a mathematical relationship giving the highest weight while computing the averages to the latest previous value and decreasing weights to the subsequent older ones

Weighted Moving Averages using Ratio Values -  $R_t = (N * R_{t-1} + (N-1) * R_{t-2} + (N-2) * R_{t-3} \dots 1 * R_{t-n}) / (N * (N+1) / 2)$

In [87]:

```

def WA_R_Predictions(ratios,month):
    predicted_ratio=(ratios['Ratios'].values)[0]
    alpha=0.5
    error=[]
    predicted_values=[]
    window_size=5
    predicted_ratio_values=[]
    for i in range(0,4464*40):
        if i%4464==0:
            predicted_ratio_values.append(0)
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_ratio_values.append(predicted_ratio)
        predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
        error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(ratios['Pred
iction'].values)[i],1))))
        if i+1>=window_size:
            sum_values=0
            sum_of_coeff=0
            for j in range(window_size,0,-1):
                sum_values += j*(ratios['Ratios'].values)[i-window_size+j]
                sum_of_coeff+=j

```

```

        sum_of_coeff+=j
        predicted_ratio=sum_values/sum_of_coeff
    else:
        sum_values=0
        sum_of_coeff=0
        for j in range(i+1,0,-1):
            sum_values += j*(ratios['Ratios'].values)[j-1]
            sum_of_coeff+=j
        predicted_ratio=sum_values/sum_of_coeff

    ratios['WA_R_Predicted'] = predicted_values
    ratios['WA_R_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err

```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 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 [88]:

```

def WA_P_Predictions(ratios,month):
    predicted_value=(ratios['Prediction'].values)[0]
    error=[]
    predicted_values=[]
    window_size=2
    for i in range(0,4464*40):
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
        if i+1>=window_size:
            sum_values=0
            sum_of_coeff=0
            for j in range(window_size,0,-1):
                sum_values += j*(ratios['Prediction'].values)[i-window_size+j]
                sum_of_coeff+=j
            predicted_value=int(sum_values/sum_of_coeff)

        else:
            sum_values=0
            sum_of_coeff=0
            for j in range(i+1,0,-1):
                sum_values += j*(ratios['Prediction'].values)[j-1]
                sum_of_coeff+=j
            predicted_value=int(sum_values/sum_of_coeff)

    ratios['WA_P_Predicted'] = predicted_values
    ratios['WA_P_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err

```

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

## Exponential Weighted Moving Averages

[https://en.wikipedia.org/wiki/Moving\\_average#Exponential\\_moving\\_average](https://en.wikipedia.org/wiki/Moving_average#Exponential_moving_average) Through weighted averaged we have satisfied the analogy of giving higher weights to the latest value and decreasing weights to the subsequent ones but we still do not know which is the correct weighting scheme as there are 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 [89]:

```
def EA_R1_Predictions(ratios,month):
    predicted_ratio=(ratios['Ratios'].values)[0]
    alpha=0.6
    error=[]
    predicted_values=[]
    predicted_ratio_values=[]
    for i in range(0,4464*40):
        if i%4464==0:
            predicted_ratio_values.append(0)
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_ratio_values.append(predicted_ratio)
        predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
        error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(ratios['Prediction'].values)[i],1))))
        predicted_ratio = (alpha*predicted_ratio) + (1-alpha)*((ratios['Ratios'].values)[i])

    ratios['EA_R1_Predicted'] = predicted_values
    ratios['EA_R1_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err
```

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

In [90]:

```
def EA_P1_Predictions(ratios,month):
    predicted_value= (ratios['Prediction'].values)[0]
    alpha=0.3
    error=[]
    predicted_values=[]
    for i in range(0,4464*40):
        if i%4464==0:
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
        predicted_value =int((alpha*predicted_value) + (1-alpha)*((ratios['Prediction'].values)[i])
    )

    ratios['EA_P1_Predicted'] = predicted_values
    ratios['EA_P1_Error'] = error
    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err
```

In [91]:

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

## Comparison between baseline models

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

In [92]:

```
print ("Error Metric Matrix (Forecasting Methods) - MAPE & MSE")
print ("-----")
print ("Moving Averages (Ratios) - MAPE: ",mean_err[0]," MSE: ",me
ian_err[0])
print ("Moving Averages (2016 Values) - MAPE: ",mean_err[1]," MSE: ",m
edian_err[1])
print ("-----")
print ("Weighted Moving Averages (Ratios) - MAPE: ",mean_err[2]," MSE: ",me
dian_err[2])
print ("Weighted Moving Averages (2016 Values) - MAPE: ",mean_err[3]," MSE: ",me
dian_err[3])
print ("-----")
print ("Exponential Moving Averages (Ratios) - MAPE: ",mean_err[4]," MSE: ",media
n_err[4])
print ("Exponential Moving Averages (2016 Values) - MAPE: ",mean_err[5]," MSE: ",media
n_err[5])
```

Error Metric Matrix (Forecasting Methods) - MAPE & MSE

```
-----
-
Moving Averages (Ratios) - MAPE: 0.1821155173392136 MSE: 400.06
5504032258
Moving Averages (2016 Values) - MAPE: 0.14292849686975506 MSE: 174.
4901993727598
-----
-
Weighted Moving Averages (Ratios) - MAPE: 0.1784869254376018 MSE:
384.01578741039424
Weighted Moving Averages (2016 Values) - MAPE: 0.13551088436182082 MSE:
162.46707549283155
-----
-
Exponential Moving Averages (Ratios) - MAPE: 0.17783550194861494 MSE:
378.34610215053766
Exponential Moving Averages (2016 Values) - MAPE: 0.1350915263669572 MSE:
159.73614471326164
-----
```

**Please 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 [223]:

```
# Preparing data to be split into train and test, The below prepares data in cumulative form which
will be later split into test and train
# number of 10min indices for jan 2015= 24*31*60/10 = 4464
# number of 10min indices for jan 2016 = 24*31*60/10 = 4464
# number of 10min indices for feb 2016 = 24*29*60/10 = 4176
# number of 10min indices for march 2016 = 24*31*60/10 = 4464
```

```

# regions_cum: it will contain 40 lists, each list will contain 4464+4176+4464 values which represents the number of pickups
# that are happened for three months in 2016 data
number_of_pickups = 4464+4176+4464
# print(len(regions_cum))
# 40
# print(len(regions_cum[0]))
# 12960

# we take number of pickups that are happened in last 30 10min intravels
number_of_time_stamps = 40
total_timestamps = number_of_pickups - number_of_time_stamps
# output variable
# it is list of lists
# it will contain number of pickups 13094 for each cluster
output = []

# tsne_lat will contain 13104-10=13094 times latitude of cluster center for every cluster
# Ex: [[cent_lat 13099times],[cent_lat 13099times], [cent_lat 13094times].... 40 lists]
# it is list of lists
tsne_lat = []

# tsne_lon will contain 13104-5=13099 times logitude of cluster center for every cluster
# Ex: [[cent_long 13099times],[cent_long 13099times], [cent_long 13099times].... 40 lists]
# it is list of lists
tsne_lon = []

# we will code each day
# sunday = 0, monday=1, tue = 2, wed=3, thur=4, fri=5,sat=6
# for every cluster we will be adding 13094 values, each value represent to which day of the week that pickup bin belongs to
# it is list of lists
tsne_weekday = []

# its an numpy array, of shape (523760, 5)
# each row corresponds to an entry in out data
# for the first row we will have [f0,f1,f2,f3,f4] fi=number of pickups happened in i+1th 10min interval(bin)
# the second row will have [f1,f2,f3,f4,f5]
# the third row will have [f2,f3,f4,f5,f6]
# and so on...
tsne_feature = []

tsne_feature = [0]*number_of_time_stamps
for i in range(0,40):
    tsne_lat.append([kmeans.cluster_centers_[i][0]]*total_timestamps)
    tsne_lon.append([kmeans.cluster_centers_[i][1]]*total_timestamps)
    # jan 1st 2016 is thursday, so we start our day from 4: "(int(k/144))%7+4"
    # our prediction start from 5th 10min intravel since we need to have number of pickups that are happened in last 15 pickup bins
    tsne_weekday.append([int(((int(k/144))%7+4)%7) for k in range(number_of_time_stamps,number_of_pickups)])
    # regions_cum is a list of lists [[x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3..x13104], .. 40 lists]
    tsne_feature = np.vstack((tsne_feature, [regions_cum[i][r:r+number_of_time_stamps] for r in range(0,len(regions_cum[i])-number_of_time_stamps)]))
    output.append(regions_cum[i][number_of_time_stamps:])
tsne_feature = tsne_feature[1:]

```

In [224]:

```
tsne_feature.shape
```

Out[224]:

```
(522560, 40)
```

In [225]:

```
FFT_Final = FFT_DataSet[:tsne_feature.shape[0]]
print(FFT_Final.shape)
```

```
(522560, 10)
```

In [226]:

```
len(tsne_lat[0])*len(tsne_lat) == tsne_feature.shape[0] == len(tsne_weekday)*len(tsne_weekday[0]) =  
= 40*total_timestamps == len(output)*len(output[0])
```

Out[226]:

True

In [227]:

```
# Getting the predictions of exponential moving averages to be used as a feature in cumulative for  
m  
  
# upto now we computed 8 features for every data point that starts from 50th min of the day  
# 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-1th 10min intravel  
# 5. f_t_2: number of pickups that are happened previous t-2th 10min intravel  
# 6. f_t_3: number of pickups that are happened previous t-3th 10min intravel  
# 7. f_t_4: number of pickups that are happened previous t-4th 10min intravel  
# 8. f_t_5: number of pickups that are happened previous t-5th 10min intravel  
  
# 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 avarage =>  $p'(t) = \alpha * p'(t-1) + (1-\alpha) * P(t-1)$   
alpha=0.3  
  
# it is a temporary array that store exponential weighted moving avarage for each 10min intravel,  
# for each cluster it will get reset  
# for every cluster it contains 13104 values  
predicted_values=[]  
  
# it is similar like tsne_lat  
# it is list of lists  
# predict_list is a list of lists [[x5,x6,x7..x13104], [x5,x6,x7..x13104], [x5,x6,x7..x13104], [x5  
,x6,x7..x13104], [x5,x6,x7..x13104], .. 40 lsits]  
predict_list = []  
tsne_flat_exp_avg = []  
for r in range(0,40):  
    for i in range(0,number_of_pickups):  
        if i==0:  
            predicted_value= regions_cum[r][0]  
            predicted_values.append(0)  
            continue  
            predicted_values.append(predicted_value)  
            predicted_value =int((alpha*predicted_value) + (1-alpha)*(regions_cum[r][i]))  
        predict_list.append(predicted_values[number_of_time_stamps:])  
        predicted_values=[]
```

In [228]:

```
# train, test split : 70% 30% split  
# Before we start predictions using the tree based regression models we take 3 months of 2016 pick  
up 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  
train_data_size = int(total_timestamps*0.7)  
test_data_size = int(total_timestamps*0.3)  
print("size of train data :",train_data_size)  
print("size of test data :", test_data_size)
```

size of train data : 9144

size of test data : 3919

In [229]:

```
# extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our training data
```



```

train_features = [tsne_feature[i*total_timestamps:(total_timestamps*i+train_data_size)] for i in range(0,40)]
# temp = [0]*(12955 - 9068)
test_features = [tsne_feature[(total_timestamps*(i))+train_data_size:total_timestamps*(i+1)] for i in range(0,40)]

```

In [230]:

```

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]),"features")
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]),"features")

```

Number of data clusters 40 Number of data points in trian data 9144 Each data point contains 40 features  
 Number of data clusters 40 Number of data points in test data 3920 Each data point contains 40 features

In [231]:

```

# extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our training data
tsne_train_flat_lat = [i[:train_data_size] for i in tsne_lat]
tsne_train_flat_lon = [i[:train_data_size] for i in tsne_lon]
tsne_train_flat_weekday = [i[:train_data_size] for i in tsne_weekday]
tsne_train_flat_output = [i[:train_data_size] for i in output]
tsne_train_flat_exp_avg = [i[:train_data_size] for i in predict_list]

```

In [232]:

```

# extracting the rest of the timestamp values i.e 30% of 12956 (total timestamps) for our test data
tsne_test_flat_lat = [i[train_data_size:] for i in tsne_lat]
tsne_test_flat_lon = [i[train_data_size:] for i in tsne_lon]
tsne_test_flat_weekday = [i[train_data_size:] for i in tsne_weekday]
tsne_test_flat_output = [i[train_data_size:] for i in output]
tsne_test_flat_exp_avg = [i[train_data_size:] for i in predict_list]

```

In [233]:

```

# 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
train_new_features = []
for i in range(0,40):
    train_new_features.extend(train_features[i])
test_new_features = []
for i in range(0,40):
    test_new_features.extend(test_features[i])

```

In [234]:

```

# converting lists of lists into sinle list i.e flatten
# a = [[1,2,3,4],[4,6,7,8]]
# print(sum(a,[]))
# [1, 2, 3, 4, 4, 6, 7, 8]

tsne_train_lat = sum(tsne_train_flat_lat, [])
tsne_train_lon = sum(tsne_train_flat_lon, [])
tsne_train_weekday = sum(tsne_train_flat_weekday, [])
tsne_train_output = sum(tsne_train_flat_output, [])
tsne_train_exp_avg = sum(tsne_train_flat_exp_avg, [])

```

In [235]:

```

# converting lists of lists into sinle list i.e flatten
# a = [[1,2,3,4],[4,6,7,8]]
# print(sum(a,[]))
# [1, 2, 3, 4, 4, 6, 7, 8]

tsne_test_lat = sum(tsne_test_flat_lat, [])
tsne_test_lon = sum(tsne_test_flat_lon, [])
tsne_test_weekday = sum(tsne_test_flat_weekday, [])

```

```
tsne_test_weekday = sum(tsne_test_flat_weekday, [])
tsne_test_output = sum(tsne_test_flat_output, [])
tsne_test_exp_avg = sum(tsne_test_flat_exp_avg, [])
```

In [240]:

```
#columns
columns = []

for i in range(0,number_of_time_stamps):
    columns.append('ft_'+str(i+1))

print(columns)
```

```
['ft_1', 'ft_2', 'ft_3', 'ft_4', 'ft_5', 'ft_6', 'ft_7', 'ft_8', 'ft_9', 'ft_10', 'ft_11',
'ft_12', 'ft_13', 'ft_14', 'ft_15', 'ft_16', 'ft_17', 'ft_18', 'ft_19', 'ft_20', 'ft_21', 'ft_22',
'ft_23', 'ft_24', 'ft_25', 'ft_26', 'ft_27', 'ft_28', 'ft_29', 'ft_30', 'ft_31', 'ft_32', 'ft_33',
'ft_34', 'ft_35', 'ft_36', 'ft_37', 'ft_38', 'ft_39', 'ft_40']
```

In [241]:

```
# Preparing the data frame for our train data
#columns =
['ft_20','ft_19','ft_18','ft_17','ft_16','ft_15','ft_14','ft_13','ft_12','ft_11','ft_10','ft_9','ft_8',
'ft_7','ft_6','ft_5','ft_4','ft_3','ft_2','ft_1']
df_train = pd.DataFrame(data=train_new_features, columns=columns)
df_train['lat'] = tsne_train_lat
df_train['lon'] = tsne_train_lon
df_train['weekday'] = tsne_train_weekday
df_train['exp_avg'] = tsne_train_exp_avg

print(df_train.shape)
```

(365760, 44)

In [243]:

```
# Preparing the data frame for our train data
df_test = pd.DataFrame(data=test_new_features, columns=columns)
df_test['lat'] = tsne_test_lat
df_test['lon'] = tsne_test_lon
df_test['weekday'] = tsne_test_weekday
df_test['exp_avg'] = tsne_test_exp_avg
print(df_test.shape)
```

(156800, 44)

In [244]:

```
FFT_train = FFT_Final[:df_train.shape[0]]
FFT_test = FFT_Final[df_train.shape[0]:]
```

In [245]:

```
FFT_train.shape
```

Out[245]:

(365760, 10)

In [246]:

```
FFT_test = FFT_test.reset_index(drop=True)
FFT_test.shape
```

Out[246]:

(156800, 10)

In [247]:

```
FFT_test.head()
```

Out[247]:

	amp1	amp2	amp3	amp4	amp5	freq1	freq2	freq3	freq4	freq5
0	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778	0.333333	0.333333
1	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778	0.333333	0.333333
2	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778	0.333333	0.333333
3	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778	0.333333	0.333333
4	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778	0.333333	0.333333

In [248]:

```
df_test.head()
```

Out[248]:

	ft_1	ft_2	ft_3	ft_4	ft_5	ft_6	ft_7	ft_8	ft_9	ft_10	...	ft_35	ft_36	ft_37	ft_38	ft_39	ft_40	lat	lon	weekday	e
0	116	140	127	123	137	128	111	109	129	132	...	154	133	148	172	153	162	40.776228	73.982119	-	4
1	140	127	123	137	128	111	109	129	132	127	...	133	148	172	153	162	200	40.776228	73.982119	-	4
2	127	123	137	128	111	109	129	132	127	106	...	148	172	153	162	200	169	40.776228	73.982119	-	4
3	123	137	128	111	109	129	132	127	106	115	...	172	153	162	200	169	149	40.776228	73.982119	-	4
4	137	128	111	109	129	132	127	106	115	133	...	153	162	200	169	149	168	40.776228	73.982119	-	4

5 rows × 44 columns

In [249]:

```
df_train = pd.concat([df_train,FFT_train],axis=1)
```

In [250]:

```
df_train.head()
```

Out[250]:

	ft_1	ft_2	ft_3	ft_4	ft_5	ft_6	ft_7	ft_8	ft_9	ft_10	...	amp1	amp2	amp3	amp4	amp5	freq1	freq2	
0	0	63	217	189	137	135	129	150	164	152	...	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.47
1	63	217	189	137	135	129	150	164	152	131	...	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.47
2	217	189	137	135	129	150	164	152	131	138	...	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.47
3	189	137	135	129	150	164	152	131	138	147	...	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.47
4	137	135	129	150	164	152	131	138	147	127	...	3242.0	76.100043	76.100043	91.706957	91.706957	0.0	0.472222	0.47

5 rows × 54 columns

In [251]:

```
df_test = pd.concat([df_test,FFT_test],axis=1)
```

In [252]:

```
df_test.shape
```

Out[252]:

(156800, 54)

In [253]:

```
df_test.head()
```

Out[253]:

	ft_1	ft_2	ft_3	ft_4	ft_5	ft_6	ft_7	ft_8	ft_9	ft_10	...	amp1	amp2	amp3	amp4	amp5	freq1	freq2	freq3
0	116	140	127	123	137	128	111	109	129	132	...	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778
1	140	127	123	137	128	111	109	129	132	127	...	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778
2	127	123	137	128	111	109	129	132	127	106	...	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778
3	123	137	128	111	109	129	132	127	106	115	...	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778
4	137	128	111	109	129	132	127	106	115	133	...	14.0	7.130255	7.130255	3.605551	3.605551	0.0	0.027778	0.027778

5 rows × 54 columns

## Using Linear Regression

In [254]:

```
# find more about LinearRegression function here http://scikit-learn.org/stable/modules/generated/sklearn.linear\_model.LinearRegression.html
# -----
# default paramters
# sklearn.linear_model.LinearRegression(fit_intercept=True, normalize=False, copy_X=True, n_jobs=1
# )

# some of methods of LinearRegression()
# fit(X, y[, sample_weight]) Fit linear model.
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict using the linear model
# score(X, y[, sample_weight]) Returns the coefficient of determination R^2 of the prediction.
# set_params(**params) Set the parameters of this estimator.
# -----
# video link: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/geometric-in-tuition-1-2-copy-8/
# -----

from sklearn.linear_model import LinearRegression
from sklearn.model_selection import GridSearchCV

# Create linear regression
lr_reg=LinearRegression()

# Create hyperparameter options
parameters = {'fit_intercept':[True,False], 'normalize':[True,False], 'copy_X':[True, False]}

# Create grid search using 10-fold cross validation
clf = GridSearchCV(lr_reg, parameters, cv=10, verbose=1)

best_model = clf.fit(df_train, tsne_train_output)

y_pred = best_model.predict(df_test)
lr_test_predictions = [round(value) for value in y_pred]
y_pred = best_model.predict(df_train)
lr_train_predictions = [round(value) for value in y_pred]
```

Fitting 10 folds for each of 8 candidates, totalling 80 fits

[Parallel(n\_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.  
[Parallel(n\_jobs=1)]: Done 80 out of 80 | elapsed: 1.4min finished

## Using Random Forest Regressor

## Using Random Forest Regressor

In [255]:

```
# Training a hyper-parameter tuned random forest regressor on our train data
# find more about LinearRegression function here http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html
# -----
# default paramters
# sklearn.ensemble.RandomForestRegressor(n_estimators=10, criterion='mse', max_depth=None, min_samples_split=2,
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False)

# some of methods of RandomForestRegressor()
# apply(X) Apply trees in the forest to X, return leaf indices.
# decision_path(X) Return the decision path in the forest
# fit(X, y[, sample_weight]) Build a forest of trees from the training set (X, y).
# get_params([deep]) Get parameters for this estimator.
# predict(X) Predict regression target for X.
# score(X, y[, sample_weight]) Returns the coefficient of determination R^2 of the prediction.
# -----
# video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-using-decision-trees-2/
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
# -----

from sklearn.model_selection import RandomizedSearchCV

# create RandomForestRegressor model
reg1 = RandomForestRegressor(max_features='sqrt',min_samples_leaf=4,min_samples_split=3,n_jobs=-1)

# number of estimators
estimators = [10,20,30,40,50,60,70,80,90,100]
# max depths
max_depth = [100,200,300,400,500,600,700,800,900,1000]
#parameters
parameters = dict(n_estimators = estimators,max_depth = max_depth)

# Create randomized search using 10-fold cross validation
clf = RandomizedSearchCV(reg1, parameters, cv=10, verbose=1)

best_model = clf.fit(df_train, tsne_train_output)

#reg1.fit(df_train, tsne_train_output)
```

Fitting 10 folds for each of 10 candidates, totalling 100 fits

```
[Parallel(n_jobs=1)]: Using backend SequentialBackend with 1 concurrent workers.
[Parallel(n_jobs=1)]: Done 100 out of 100 | elapsed: 76.5min finished
```

In [256]:

```
# Predicting on test data using our trained random forest model

# the models reg1 is already hyper parameter tuned
# the parameters that we got above are found using grid search

y_pred = best_model.predict(df_test)
rndf_test_predictions = [round(value) for value in y_pred]
y_pred = best_model.predict(df_train)
rndf_train_predictions = [round(value) for value in y_pred]
```

In [257]:

```
print(best_model.best_params_)
```

```
{'n_estimators': 90, 'max_depth': 500}
```

In [258]:

```
#feature importances based on analysis using random forest
print (df_train.columns)
print (best_model.best_estimator_.feature_importances_)
```

```
Index(['ft_1', 'ft_2', 'ft_3', 'ft_4', 'ft_5', 'ft_6', 'ft_7', 'ft_8', 'ft_9',
      'ft_10', 'ft_11', 'ft_12', 'ft_13', 'ft_14', 'ft_15', 'ft_16', 'ft_17',
      'ft_18', 'ft_19', 'ft_20', 'ft_21', 'ft_22', 'ft_23', 'ft_24', 'ft_25',
      'ft_26', 'ft_27', 'ft_28', 'ft_29', 'ft_30', 'ft_31', 'ft_32', 'ft_33',
      'ft_34', 'ft_35', 'ft_36', 'ft_37', 'ft_38', 'ft_39', 'ft_40', 'lat',
      'lon', 'weekday', 'exp_avg', 'amp1', 'amp2', 'amp3', 'amp4', 'amp5',
      'freq1', 'freq2', 'freq3', 'freq4', 'freq5'],
      dtype='object')
[8.50266737e-04  9.12709688e-04  8.14274737e-04  1.82558274e-03
 8.75523311e-04  7.56928075e-04  7.84805610e-04  1.14429368e-03
 6.76901376e-04  8.41612697e-04  3.08516274e-03  6.99333766e-04
 9.05372480e-04  6.95647507e-04  7.80830693e-04  7.87235581e-04
 7.28211497e-04  1.26637804e-03  6.86521054e-04  1.24503118e-03
 1.22916856e-03  7.76750387e-04  1.99261998e-03  5.48900507e-03
 1.10257003e-02  1.27622678e-02  7.71542571e-03  2.14747695e-02
 1.52015795e-02  2.30992450e-02  3.03613157e-02  5.06710345e-02
 3.53863019e-02  7.63124569e-02  6.19823512e-02  4.72189174e-02
 9.34888631e-02  7.91933897e-02  8.75505067e-02  1.54544275e-01
 6.94798004e-04  7.10703739e-04  3.18979461e-04  1.56600511e-01
 5.80030677e-04  5.08651026e-04  5.03588396e-04  5.02264280e-04
 4.94847433e-04  4.78035918e-06  2.91580514e-04  2.79286346e-04
 3.40035410e-04  3.31376312e-04]
```

## Using XgBoost Regressor

In [259]:

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

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

# some of methods of RandomForestRegressor()
# fit(X, y, sample_weight=None, eval_set=None, eval_metric=None, early_stopping_rounds=None, verbo
se=True, xgb_model=None)
# get_params([deep]) Get parameters for this estimator.
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE: This function is no
t thread safe.
# get_score(importance_type='weight') -> get the feature importance
# -----
# video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-using-decision-trees-2/
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-en-sembles/
# -----

x_model = xgb.XGBRegressor(
    min_child_weight=3,
    learning_rate=0.1,
    gamma=0,
    subsample=0.8,
    reg_alpha=200, reg_lambda=200,
    colsample_bytree=0.8, nthread=4)

parameters = dict(n_estimators = estimators, max_depth = max_depth)
# Create randomized search using 10-fold cross validation
clf = RandomizedSearchCV(x_model, parameters, cv=10, verbose=0)
```



1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65 66 67 68 69 70 71 72 73 74 75 76 77 78 79 80 81 82 83 84 85 86 87 88 89 90 91 92 93 94 95 96 97 98 99 100 101 102 103 104 105 106 107 108 109 110 111 112 113 114 115 116 117 118 119 120 121 122 123 124 125 126 127 128 129 130 131 132 133 134 135 136 137 138 139 140 141 142 143 144 145 146 147 148 149 150 151 152 153 154 155 156 157 158 159 160 161 162 163 164 165 166 167 168 169 170 171 172 173 174 175 176 177 178 179 180 181 182 183 184 185 186 187 188 189 190 191 192 193 194 195 196 197 198 199 200 201 202 203 204 205 206 207 208 209 210 211 212 213 214 215 216 217 218 219 220 221 222 223 224 225 226 227 228 229 230 231 232 233 234 235 236 237 238 239 240 241 242 243 244 245 246 247 248 249 250 251 252 253 254 255 256 257 258 259 260 261 262 263 264 265 266 267 268 269 270 271 272 273 274 275 276 277 278 279 280 281 282 283 284 285 286 287 288 289 290 291 292 293 294 295 296 297 298 299 300 301 302 303 304 305 306 307 308 309 310 311 312 313 314 315 316 317 318 319 320 321 322 323 324 325 326 327 328 329 330 331 332 333 334 335 336 337 338 339 340 341 342 343 344 345 346 347 348 349 350 351 352 353 354 355 356 357 358 359 360 361 362 363 364 365 366 367 368 369 370 371 372 373 374 375 376 377 378 379 380 381 382 383 384 385 386 387 388 389 390 391 392 393 394 395 396 397 398 399 400 401 402 403 404 405 406 407 408 409 410 411 412 413 414 415 416 417 418 419 420 421 422 423 424 425 426 427 428 429 430 431 432 433 434 435 436 437 438 439 440 441 442 443 444 445 446 447 448 449 450 451 452 453 454 455 456 457 458 459 460 461 462 463 464 465 466 467 468 469 470 471 472 473 474 475 476 477 478 479 480 481 482 483 484 485 486 487 488 489 490 491 492 493 494 495 496 497 498 499 500 501 502 503 504 505 506 507 508 509 510 511 512 513 514 515 516 517 518 519 520 521 522 523 524 525 526 527 528 529 530 531 532 533 534 535 536 537 538 539 540 541 542 543 544 545 546 547 548 549 550 551 552 553 554 555 556 557 558 559 560 561 562 563 564 565 566 567 568 569 570 571 572 573 574 575 576 577 578 579 580 581 582 583 584 585 586 587 588 589 590 591 592 593 594 595 596 597 598 599 600 601 602 603 604 605 606 607 608 609 610 611 612 613 614 615 616 617 618 619 620 621 622 623 624 625 626 627 628 629 630 631 632 633 634 635 636 637 638 639 640 641 642 643 644 645 646 647 648 649 650 651 652 653 654 655 656 657 658 659 660 661 662 663 664 665 666 667 668 669 670 671 672 673 674 675 676 677 678 679 680 681 682 683 684 685 686 687 688 689 690 691 692 693 694 695 696 697 698 699 700 701 702 703 704 705 706 707 708 709 710 711 712 713 714 715 716 717 718 719 720 721 722 723 724 725 726 727 728 729 730 731 732 733 734 735 736 737 738 739 740 741 742 743 744 745 746 747 748 749 750 751 752 753 754 755 756 757 758 759 760 761 762 763 764 765 766 767 768 769 770 771 772 773 774 775 776 777 778 779 780 781 782 783 784 785 786 787 788 789 790 791 792 793 794 795 796 797 798 799 800 801 802 803 804 805 806 807 808 809 810 811 812 813 814 815 816 817 818 819 820 821 822 823 824 825 826 827 828 829 830 831 832 833 834 835 836 837 838 839 840 841 842 843 844 845 846 847 848 849 850 851 852 853 854 855 856 857 858 859 860 861 862 863 864 865 866 867 868 869 870 871 872 873 874 875 876 877 878 879 880 881 882 883 884 885 886 887 888 889 890 891 892 893 894 895 896 897 898 899 900 901 902 903 904 905 906 907 908 909 910 911 912 913 914 915 916 917 918 919 920 921 922 923 924 925 926 927 928 929 930 931 932 933 934 935 936 937 938 939 940 941 942 943 944 945 946 947 948 949 950 951 952 953 954 955 956 957 958 959 960 961 962 963 964 965 966 967 968 969 970 971 972 973 974 975 976 977 978 979 980 981 982 983 984 985 986 987 988 989 990 991 992 993 994 995 996 997 998 999 1000 1001 1002 1003 1004 1005 1006 1007 1008 1009 1010 1011 1012 1013 1014 1015 1016 1017 1018 1019 1020 1021 1022 1023 1024 1025 1026 1027 1028 1029 1030 1031 1032 1033 1034 1035 1036 1037 1038 1039 1040 1



[illegible]

[illegible]

In [260]:

```
# best params
print(best_model.best_params_)

{'n_estimators': 80, 'max_depth': 1000}
```

In [261]:

```
#predicting with our trained Xg-Boost regressor
# the models x_model is already hyper parameter tuned
# the parameters that we got above are found using grid search

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

In [262]:

```
#feature importances based on analysis using Xg-Boost regressor
print (df_train.columns)
print (best_model.best_estimator_.feature_importances_)
```

```
Index(['ft_1', 'ft_2', 'ft_3', 'ft_4', 'ft_5', 'ft_6', 'ft_7', 'ft_8', 'ft_9',
      'ft_10', 'ft_11', 'ft_12', 'ft_13', 'ft_14', 'ft_15', 'ft_16', 'ft_17',
      'ft_18', 'ft_19', 'ft_20', 'ft_21', 'ft_22', 'ft_23', 'ft_24', 'ft_25',
      'ft_26', 'ft_27', 'ft_28', 'ft_29', 'ft_30', 'ft_31', 'ft_32', 'ft_33',
      'ft_34', 'ft_35', 'ft_36', 'ft_37', 'ft_38', 'ft_39', 'ft_40', 'lat',
      'lon', 'weekday', 'exp_avg', 'amp1', 'amp2', 'amp3', 'amp4', 'amp5',
      'freq1', 'freq2', 'freq3', 'freq4', 'freq5'],
      dtype='object')
[3.5418334e-04 3.1704977e-04 2.8725908e-04 2.7568502e-04 2.9790969e-04
 2.8436180e-04 2.4933799e-04 2.8915118e-04 2.8017498e-04 2.7868417e-04
 3.0663831e-04 2.9693861e-04 2.8172130e-04 2.8742096e-04 2.8720437e-04
 3.0155331e-04 3.0219546e-04 2.5430281e-04 2.4903662e-04 2.8536329e-04
 2.8616394e-04 2.7684221e-04 3.2212940e-04 3.2533391e-04 3.2010081e-04
 3.5176121e-04 5.7714048e-04 3.9326455e-04 2.6812218e-04 2.6801552e-04
 3.3313027e-04 3.2873705e-04 3.9856075e-04 2.8351642e-04 6.1904563e-04
 7.5649086e-04 7.3094165e-04 6.3385270e-03 1.4791957e-02 2.5580955e-01
 5.4774532e-04 6.1345805e-04 4.2051674e-04 7.0675659e-01 2.2230060e-04
 1.9099837e-04 2.3598480e-04 1.8621419e-04 2.2578031e-04 1.9322192e-04
 2.2689860e-04 2.1455281e-04 2.1147827e-04 2.0863676e-04]
```

## Calculating the error metric values for various models

In [263]:

```
train_mape=[]
test_mape=[]

train_mape.append((mean_absolute_error(tsne_train_output,df_train['ft_1'].values))/(sum(tsne_train_output)/len(tsne_train_output)))
train_mape.append((mean_absolute_error(tsne_train_output,df_train['exp_avg'].values))/(sum(tsne_train_output)/len(tsne_train_output)))
train_mape.append((mean_absolute_error(tsne_train_output,rndf_train_predictions))/(sum(tsne_train_output)/len(tsne_train_output)))
train_mape.append((mean_absolute_error(tsne_train_output,xgb_train_predictions))/(sum(tsne_train_output)/len(tsne_train_output)))
train_mape.append((mean_absolute_error(tsne_train_output,lr_train_predictions))/(sum(tsne_train_output)/len(tsne_train_output)))

test_mape.append((mean_absolute_error(tsne_test_output, df_test['ft_1'].values))/(sum(tsne_test_output)/len(tsne_test_output)))
test_mape.append((mean_absolute_error(tsne_test_output, df_test['exp_avg'].values))/(sum(tsne_test_output)/len(tsne_test_output)))
test_mape.append((mean_absolute_error(tsne_test_output, rndf_test_predictions))/(sum(tsne_test_output)/len(tsne_test_output)))
test_mape.append((mean_absolute_error(tsne_test_output, xgb_test_predictions))/(sum(tsne_test_output)/len(tsne_test_output)))
test_mape.append((mean_absolute_error(tsne_test_output, lr_test_predictions))/(sum(tsne_test_output)/len(tsne_test_output)))
```

## Error Metric Matrix

In [264]:

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

Error Metric Matrix (Tree Based Regression Methods) - MAPE

```
-----
-
Baseline Model - Train: 0.7682905677289276 Test: 0.7500048682479439
Exponential Averages Forecasting - Train: 0.13278679231447216 Test: 0.12945347561727225
Linear Regression - Train: 0.1314256977361752 Test: 0.12730665021854357
Random Forest Regression - Train: 0.07265114294745745 Test: 0.12076343720367838
XgBoost Regression - Train: 0.10318715675623338 Test: 0.1198830277585033
-----
-
```

## Assignments

In [265]:

```
'''
Task 1: Incorporate Fourier features as features into Regression models and measure MAPE. <br>

Task 2: Perform hyper-parameter tuning for Regression models.
        2a. Linear Regression: Grid Search
        2b. Random Forest: Random Search
        2c. Xgboost: Random Search
Task 3: Explore more time-series features using Google search/Quora/Stackoverflow
to reduce the MAPE to < 12%
'''
```

Out[265]:

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

In [ ]: