

Taxi demand prediction in New York City



In [1]:

```
#Importing Libraries
!pip3 install graphviz
!pip3 install dask
!pip install "dask[complete]"
!pip3 install toolz
!pip3 install云dcloudpickle
# 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 doesn't work refer to 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 straight line distance
between two (lat,lon) pairs in miles
!pip install gpxpy
```

```
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")
```

```
Requirement already satisfied: folium in c:\\users\\sys.ai\\anaconda3\\lib\\site-packages (0.12.1)
Requirement already satisfied: numpy in c:\\users\\sys.ai\\appdata\\roaming\\python38\\site-packages (from folium) (1.19.5)
Requirement already satisfied: requests in c:\\users\\sys.ai\\anaconda3\\lib\\site-packages (from folium) (2.24.0)
Requirement already satisfied: branca>=0.3.0 in c:\\users\\sys.ai\\anaconda3\\lib\\site-packages (from folium) (0.4.2)
Requirement already satisfied: jinja2>=2.9 in c:\\users\\sys.ai\\anaconda3\\lib\\site-packages (from folium) (2.11.2)
Requirement already satisfied: MarkupSafe>=0.23 in c:\\users\\sys.ai\\anaconda3\\lib\\site-packages (from jinja2>=2.9->folium) (1.1.1)
Requirement already satisfied: chardet<4,>=3.0.2 in c:\\users\\sys.ai\\anaconda3\\lib\\site-packages (from requests->folium) (3.0.4)
Requirement already satisfied: urllib3!=1.25.0,!=1.25.1,<1.26,>=1.21.1 in c:\\users\\sys.ai\\anaconda3\\lib\\site-packages (from requests->folium) (1.25.11)
Requirement already satisfied: certifi>=2017.4.17 in c:\\users\\sys.ai\\anaconda3\\lib\\site-packages (from requests->folium) (2020.6.20)
Requirement already satisfied: idna<3,>=2.5 in c:\\users\\sys.ai\\anaconda3\\lib\\site-packages (from requests->folium) (2.10)
WARNING: You are using pip version 21.1.3; however, version 21.2.2 is available
```

e.
You should consider upgrading via the 'c:\users\sys.ai\anaconda3\python.exe -m
Requirement already satisfied: gpxpy in c:\users\sys.ai\anaconda3\lib\site-pac
kages (1.4.2)
WARNING: You are using pip version 21.1.3; however, version 21.2.2 is availabl
e.
You should consider upgrading via the 'c:\users\sys.ai\anaconda3\python.exe -m
pip install --upgrade pip' command.

Data Information

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

Information on taxis:

Yellow Taxi: Yellow Medallion Taxicabs

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

For Hire Vehicles (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
yellow_tripdata_2015-12	1.67Gb	11460573	19

In [7]:

```
!wget wget --header="Host: doc-04-bk-docs.googleusercontent.com"
--header="User-Agent: Mozilla/5.0 (Windows NT 10.0; Win64; x64)
AppleWebKit/537.36 (KHTML, like Gecko) Chrome/88.0.4324.190
Safari/537.36" --header="Accept: text/html,application/*
/xhtml+xml,application/xml;q=0.9,image/avif,image/webp,image/apng,
/*;q=0.8,application/signed-exchange;v=b3;q=0.9" --header="Accept-
Language: en-US,en;q=0.9" --header="Cookie:
AUTH_nso6dcn1mbidkt5qr539a2jiefc09pqv_nonce=e4vb7isdofdpc"
--header="Connection: keep-alive" "https://doc-04-bk-
docs.googleusercontent.com/docs/securesc
/nss2f5s2soorprev6d4t4qp3n5ekp9nh/mhg8tsfscqgl381loj75cigdj523bpec
/1615629600000/06629147635963609455/13017565264516993811/1kcIZlf-
LQiQhqfSCZb719Nh6Rqkp2zKK?e=download&authuser=0&
nonce=e4vb7isdofdpc&user=13017565264516993811&
hash=vuqdmlld4i50g7m2f36ji4ttkokjgnu3i" -c -O
'yellow_tripdata_2015-01.csv'
```

```
--2021-03-13 10:01:44-- http://wget/
Resolving wget (wget)... failed: Name or service not known.
wget: unable to resolve host address 'wget'
--2021-03-13 10:01:44-- https://doc-04-bk-docs.googleusercontent.com/docs/sec
uresc/nss2f5s2soorprev6d4t4qp3n5ekp9nh/mhg8tsfscqgl381loj75cigdj523bpec/161562
9600000/06629147635963609455/13017565264516993811/1kcIZlf-LQiQhqfSCZb719Nh6Rqk
p2zKK?e=download&authuser=0&nonce=e4vb7isdofdpc&user=13017565264516993811&hash
=vuqdmlld4i50g7m2f36ji4ttkokjgnu3i
Resolving doc-04-bk-docs.googleusercontent.com (doc-04-bk-docs.googleusercontent.com)...
74.125.195.132, 2607:f8b0:400e:c09::84
Connecting to doc-04-bk-docs.googleusercontent.com (doc-04-bk-docs.googleusercontent.com)|74.125.195.132|:443... connected.
HTTP request sent, awaiting response... 200 OK
Length: unspecified [text/csv]
Saving to: 'yellow_tripdata_2015-01.csv'

yellow_tripdata_201      [          =>      ] 1.85G 148MB/s   in 21s

2021-03-13 10:02:06 (88.4 MB/s) - 'yellow_tripdata_2015-01.csv' saved [1985964
692]

FINISHED --2021-03-13 10:02:06--
Total wall clock time: 22s
Downloaded: 1 files, 1.8G in 21s (88.4 MB/s)
```

In [2]:

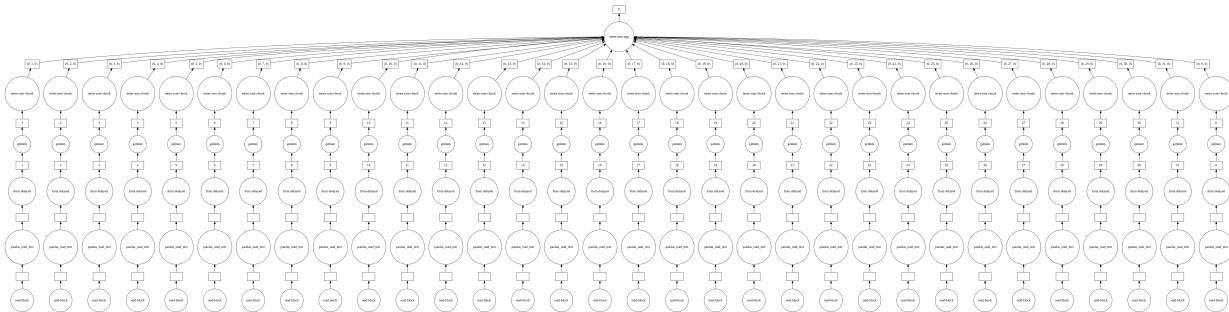
```
#Looking at the features
# dask dataframe : # https://github.com/dask/dask-tutorial
/blob/master/07_dataframe.ipynb
month = dd.read_csv('yellow_tripdata_2015-01.csv')
print(month.columns)
```

```
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'],
```

In [10]: `month.fare_amount.sum().visualize()`

Out[10]:



Features in the dataset:

Field Name	Description
	A code indicating the TPEP provider that provided the record.
VendorID	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.
	The final rate code in effect at the end of the trip.
RateCodeID	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 [12] :

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

Out[12] :

	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance	pickup_lo
0	2	2015-01-15 19:05:39	2015-01-15 19:23:42	1	1.59	-73
1	1	2015-01-10 20:33:38	2015-01-10 20:53:28	1	3.30	-74

VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance	pickup_longitude	pickup_latitude
2	1 2015-01-10 20:33:38	2015-01-10 20:43:41	1	1.80	-73	40.7744
3	1 2015-01-10 20:33:39	2015-01-10 20:35:31	1	0.50	-74	40.7744

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 [3]:

```
# 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
indeephth 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)
```

Out[3] : Make this Notebook Trusted to load map: File -> Trust Notebook

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 [4]:

```
# 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)
```

Out [4] : Make this Notebook Trusted to load map: File -> Trust Notebook

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

3. Trip Durations:

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

In [5]:

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

# in out data we have time in the formate "YYYY-MM-DD HH:MM:SS" we
convert this sting to python time formate and then into unix time
stamp
# 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

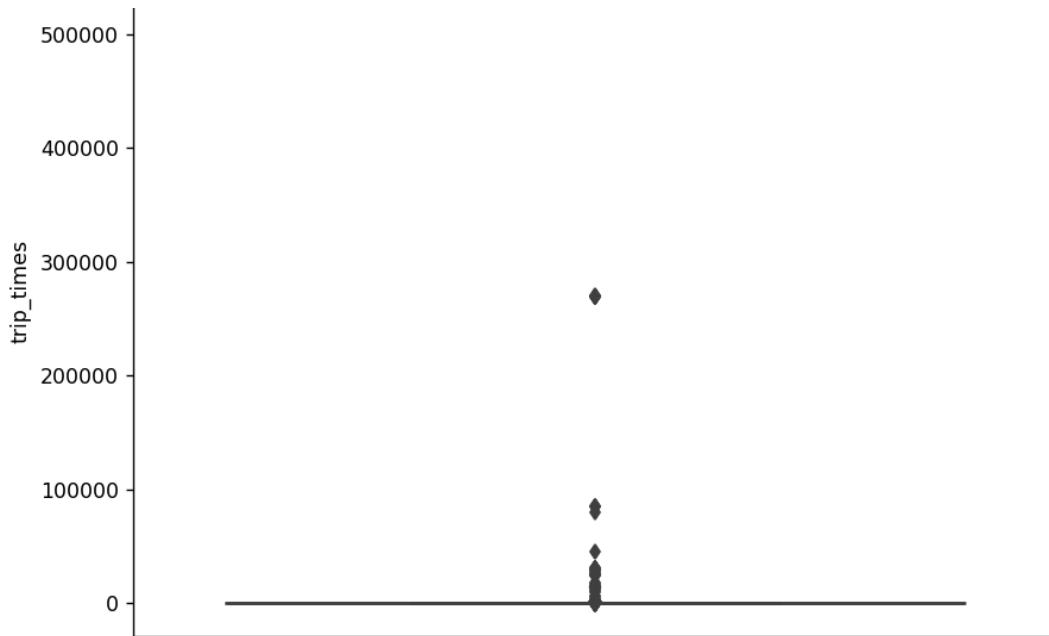
def return_with_trip_times(month):
    duration =
month[['tpep_pickup_datetime', 'tpep_dropoff_datetime']].compute()
    #pickups and dropoffs to unix time
    duration_pickup = [convert_to_unix(x) for x in
duration['tpep_pickup_datetime'].values]
    duration_drop = [convert_to_unix(x) for x in
duration['tpep_dropoff_datetime'].values]
    #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['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  
17.80            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  
16.30            19.316667      1.420902e+09      9.318378  
frame_with_durations = return_with_trip_times(month)
```

In [7]:

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



```
In [6]: #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 -1211.016666666667
10 percentile value is 3.833333333333335
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 14.28333333333333
80 percentile value is 17.63333333333333
90 percentile value is 23.45
100 percentile value is 548555.6333333333
```

In [7]:

```
#looking further from the 99th percentile
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])
```

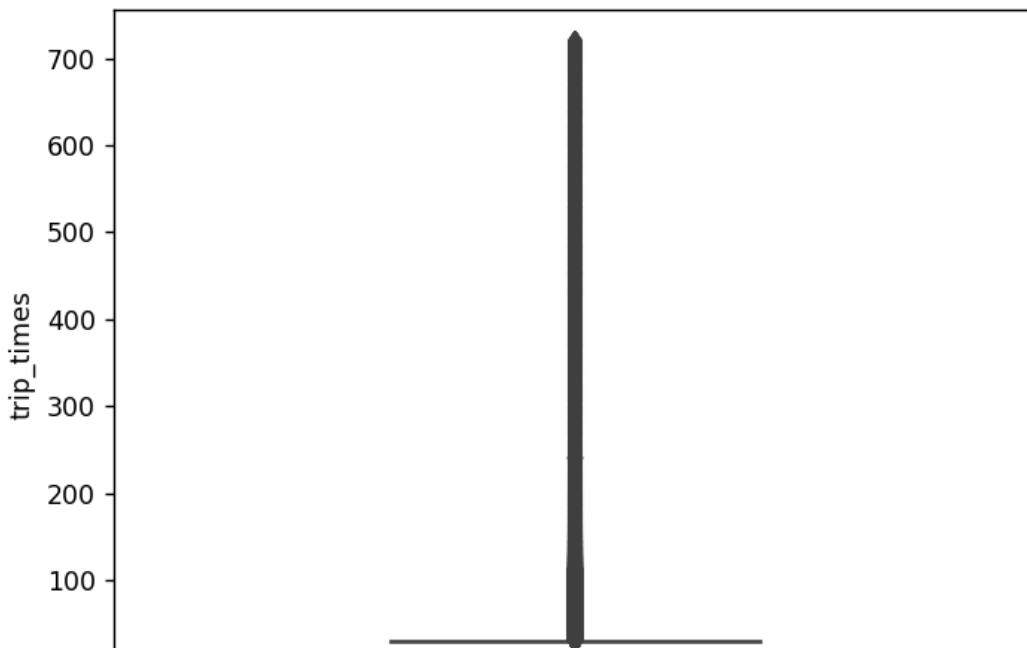
```
90 percentile value is 23.45
91 percentile value is 24.35
92 percentile value is 25.383333333333333
93 percentile value is 26.55
94 percentile value is 27.93333333333334
95 percentile value is 29.58333333333332
96 percentile value is 31.68333333333334
97 percentile value is 34.46666666666667
98 percentile value is 38.71666666666667
99 percentile value is 46.75
100 percentile value is 548555.6333333333
```

In [8]:

```
#removing data based on our analysis and TLC regulations
frame_with_durations_modified=frame_with_durations[(frame_with_durations['trip_time'] &lt; 720)]
```

In [10]:

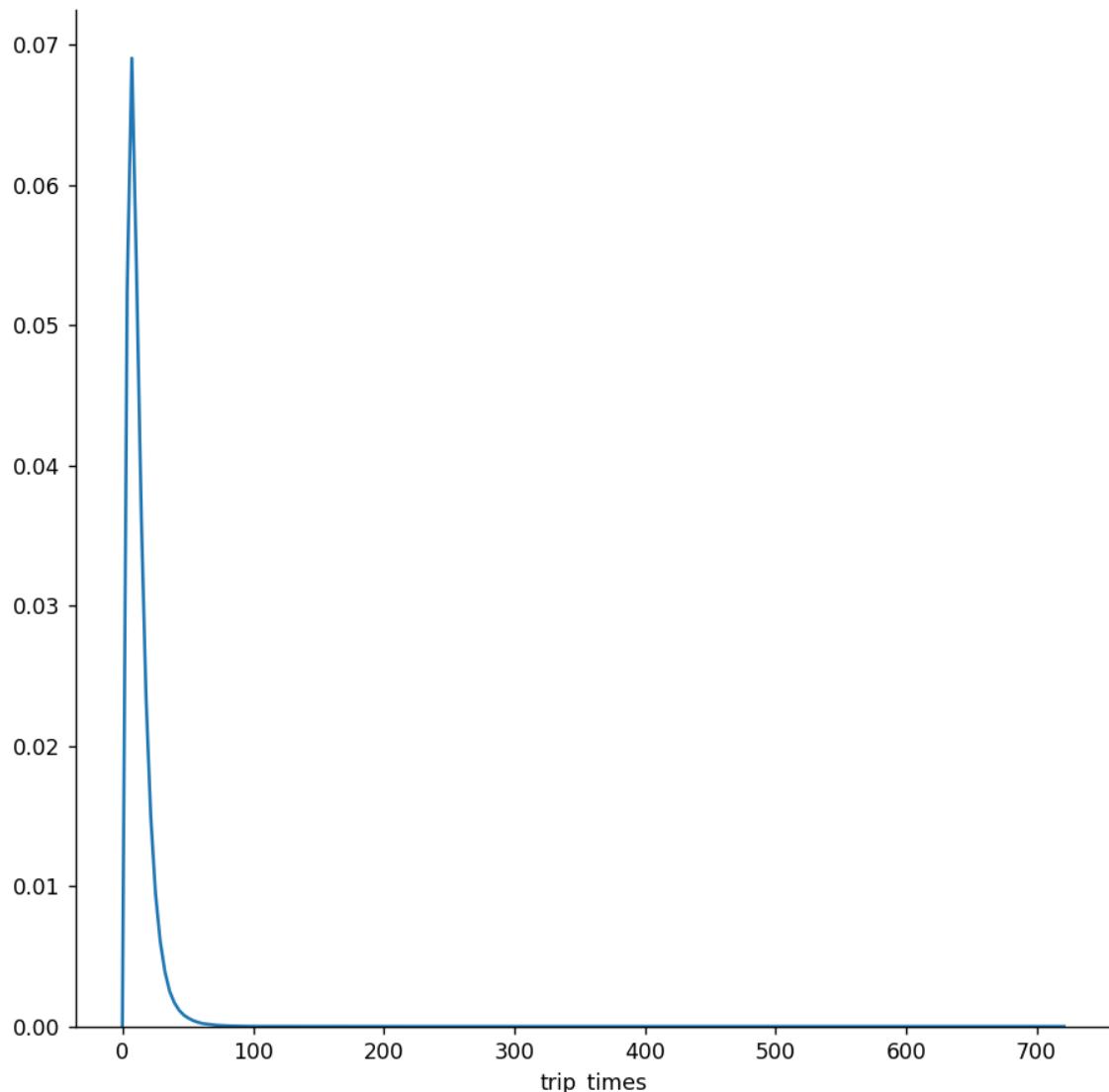
```
#box-plot after removal of outliers
sns.boxplot(y="trip_times", data =frame_with_durations_modified)
plt.show()
```





In [11]:

```
#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();
```

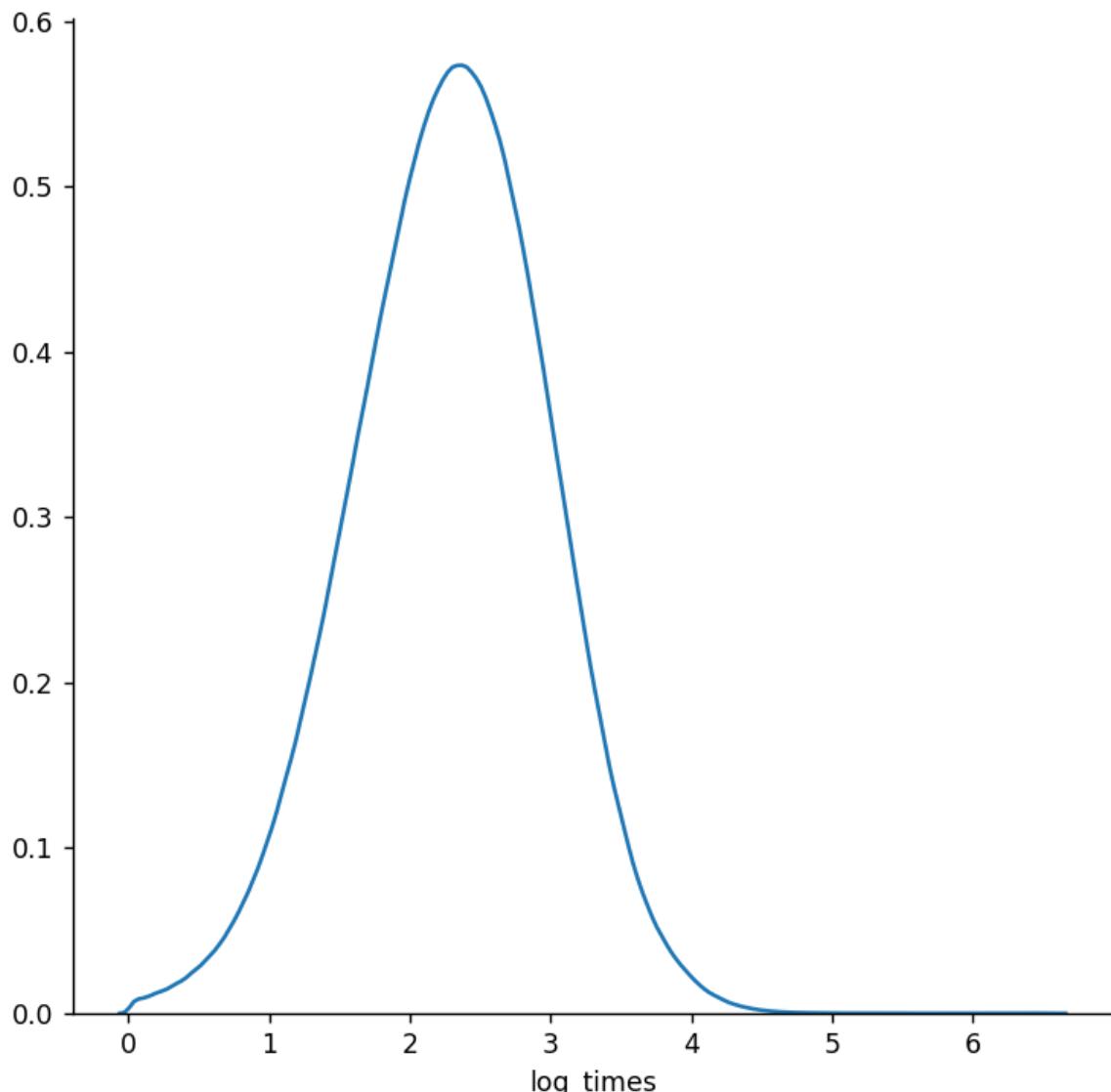


In [9]:

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

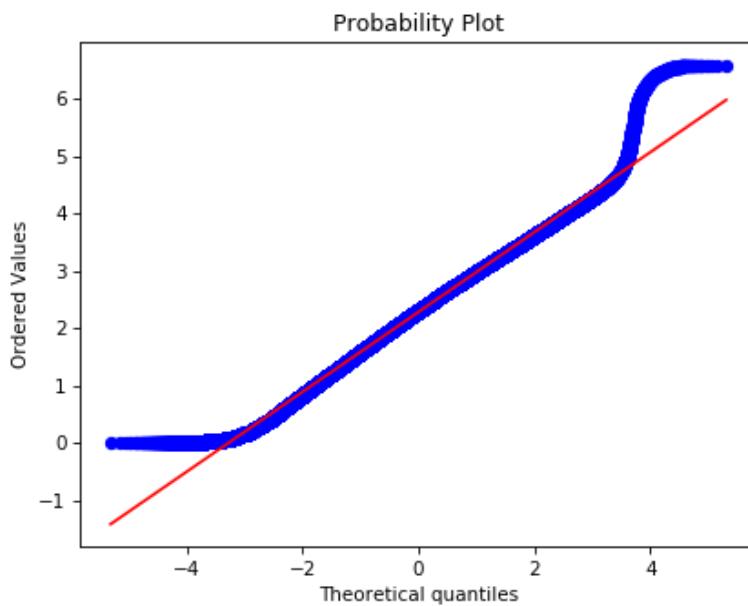
In [18]:

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



In []:

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

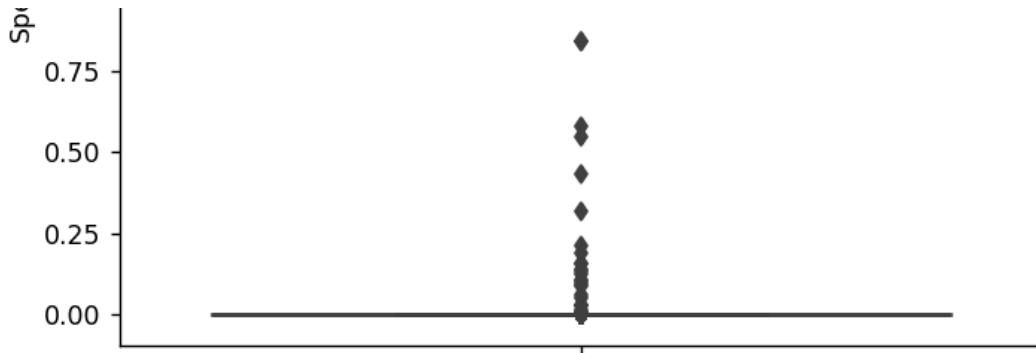


4. Speed

In [10]:

```
# 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['fare_amount'])
sns.boxplot(y="Speed", data =frame_with_durations_modified)
plt.show()
```





```
In [11]: #calculating speed 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["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 [12]: #calculating speed 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["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 [13]:

```

#calculating speed values at each percentile
99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100

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 [14]:

```

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

```

In [15]:

```

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

```

Out[15]: 12.450173996027528

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

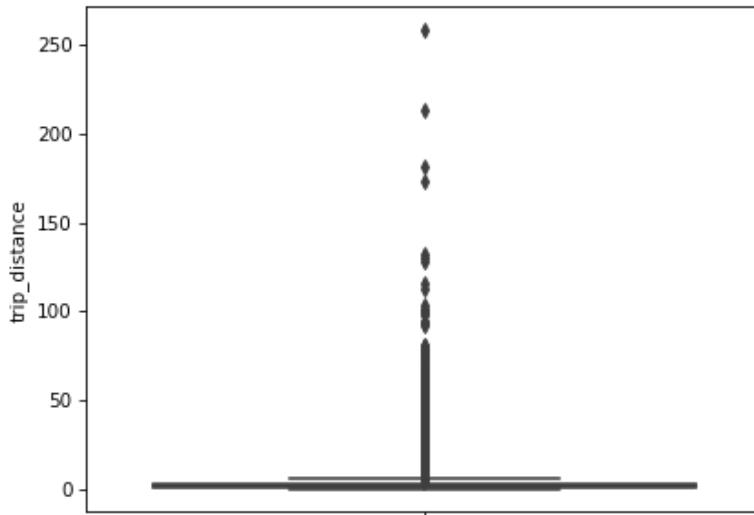
4. Trip Distance

In []:

```

# 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 [16]:

```
#calculating trip distance 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["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 [17]:

```
#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 [18]:

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

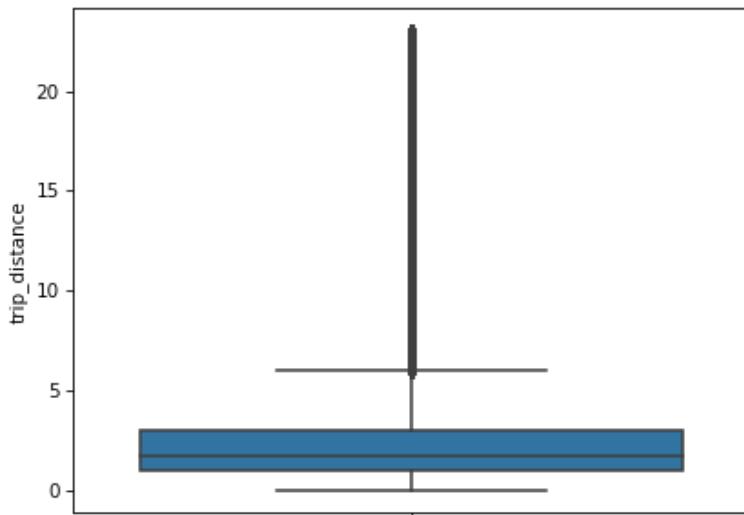
```
99.0 percentile value is 18.17
99.1 percentile value is 18.37
99.2 percentile value is 18.6
99.3 percentile value is 18.83
99.4 percentile value is 19.13
99.5 percentile value is 19.5
99.6 percentile value is 19.96
99.7 percentile value is 20.5
99.8 percentile value is 21.22
99.9 percentile value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

In [19]:

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

In []:

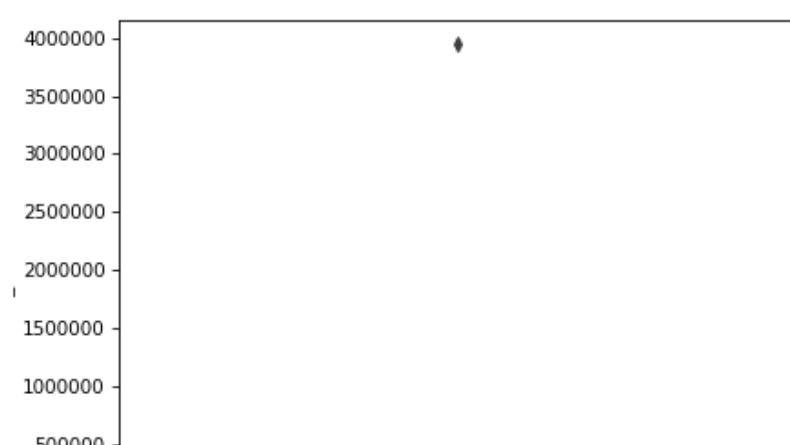
```
#box-plot after removal of outliers
sns.boxplot(y="trip_distance", data =
frame_with_durations_modified)
plt.show()
```

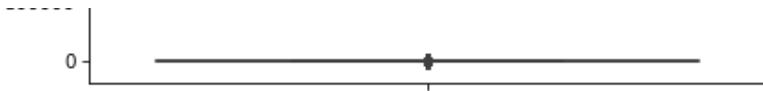


5. Total Fare

In []:

```
# 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 [20]:

```
#calculating total fare amount 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["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 [21]:

```
#calculating total fare amount 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["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 [22]:

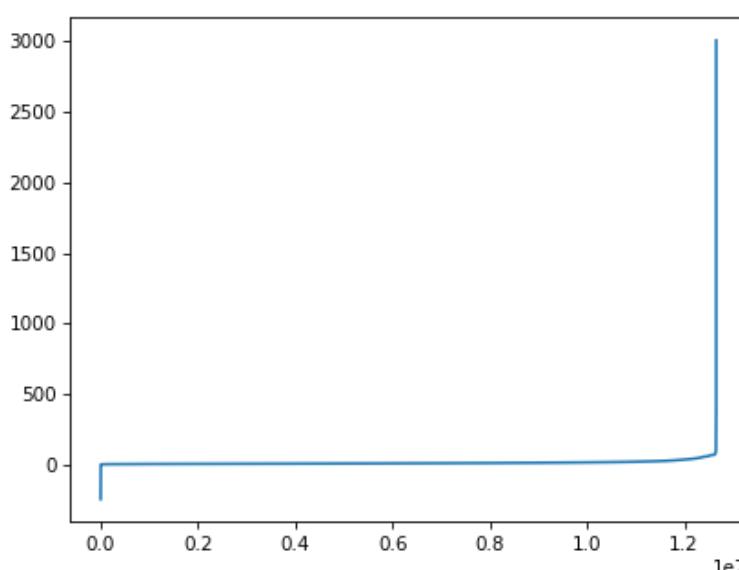
```
#calculating total fare amount values at each percentile
99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100
for i in np.arange(0.0, 1.0, 0.1):
    var = frame_with_durations_modified["total_amount"].values
    var = np.sort(var, axis = None)
    print("{} percentile value is {}".format(99+i, var[int(len(var))*float(99+i)/100]))])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

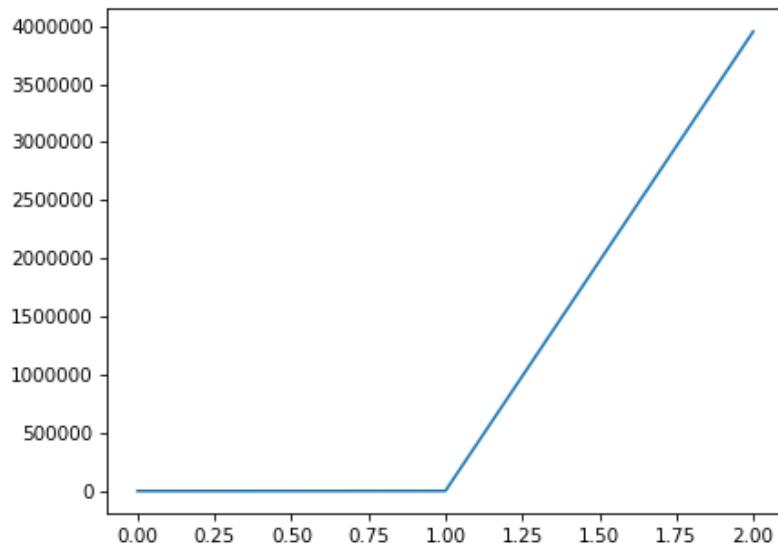
In []:

```
#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()
```



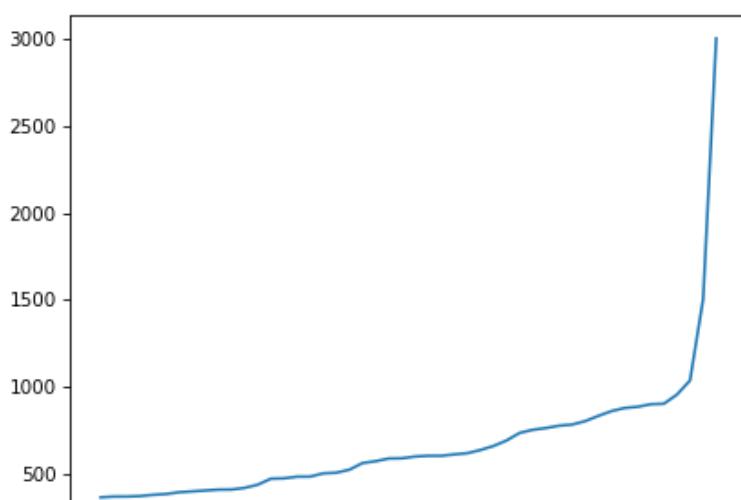
In []:

```
# 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 []:

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





Remove all outliers/errorous points.

In [23]:

```
#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_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]
```

```
        (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 [24]:

```
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_dura
```

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

Data-preperation

Clustering/Segmentation

In [25]:

```
#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 =
gpxy.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,"Avg.
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,"---")

def find_clusters(increment):
    kmeans = MiniBatchKMeans(n_clusters=increment,
batch_size=10000,random_state=42).fit(coords)
    frame_with_durations_outliers_removed['pickup_cluster'] =
```

```
# 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 =  1.0945442325142662  
---  
On choosing a cluster size of  20  
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2):  
4.0  
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):  
16.0  
Min inter-cluster distance =  0.7131298007388065  
---  
On choosing a cluster size of  30  
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2):  
8.0  
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):  
22.0  
Min inter-cluster distance =  0.5185088176172186  
---  
On choosing a cluster size of  40  
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2):  
8.0  
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):  
32.0  
Min inter-cluster distance =  0.5069768450365043  
---  
On choosing a cluster size of  50  
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2):  
12.0  
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):  
38.0  
Min inter-cluster distance =  0.36536302598358383  
---  
On choosing a cluster size of  60  
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2):  
14.0  
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):  
46.0
```

```
Min inter-cluster distance =  0.34704283494173577
---
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 =  0.30502203163245994
---
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 =  0.292203245317388
---
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 =  0.18257992857033273
```

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 [26]: # 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 [27]:

```
# 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 [27]: Make this Notebook Trusted to load map: File -> Trust Notebook

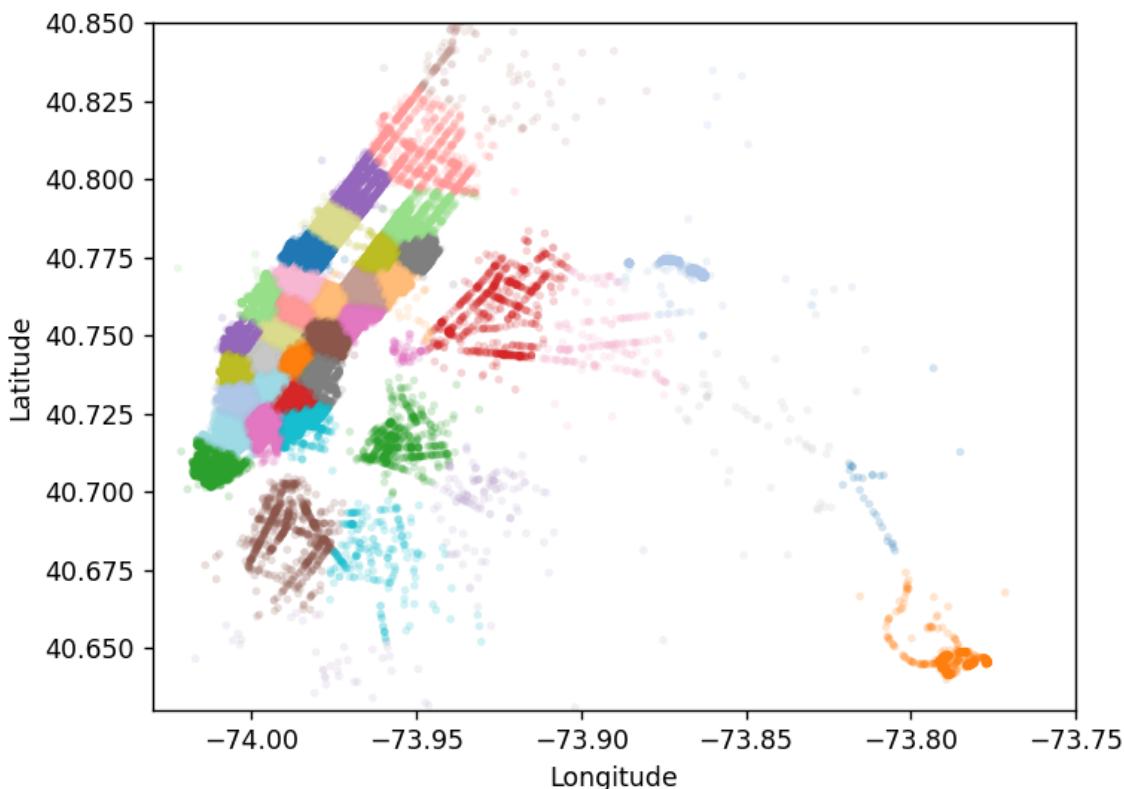
Plotting the clusters:

In [28]:

```
#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 [29]:

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

In [30]:

```
# 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']).sum().reset_index()
```

In [29]:

```
# 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[29]:

	passenger_count	trip_distance	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
0	1	1.59	-73.993896	40.750111	-73.974785	40.7506
1	1	3.30	-74.001648	40.724243	-73.994415	40.7591
2	1	1.80	-73.963341	40.802788	-73.951820	40.8244
3	1	0.50	-74.009087	40.713818	-74.004326	40.7195
4	1	3.00	-73.971176	40.762428	-74.004181	40.7426

In [30]:

```
# hear the trip_distance represents the number of pickups that are
happened in that particular 10min intravel

# this data frame has two indices

# primary index: pickup_cluster (cluster number)

# secondary index : pickup_bins (we divided whole months time into
10min intravels 24*31*60/10 =4464bins)

jan_2015_groupby.head()
```

Out[30]:

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

~~trip distance~~

In [31]:

```
# 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 columns  
# 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 'pickup_bin'  
  
# Data Preparation for the months of Jan, Feb and March 2016  
def datapreparation(month, kmeans, month_no, year_no):  
  
    print ("Return with trip times..")  
  
    frame_with_durations = return_with_trip_times(month)  
  
    print ("Remove outliers..")  
    frame_with_durations_outliers_removed =  
remove_outliers(frame_with_durations)  
  
    print ("Estimating clusters..")  
    frame_with_durations_outliers_removed['pickup_cluster'] =  
kmeans.predict(frame_with_durations_outliers_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']]  
  
return final_updated_frame, final_groupby_frame
```

```
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 [32]:

```
# Gets the unique bins where pickup values are present for each
# each reigion

# for each cluster region we will collect all the indices of 10min
intravels in which the pickups are 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 [33]:

```
# 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 []:

```
# for each cluster number of 10min intavels 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 => ceil(x/4), ceil(x/4), ceil(x/4), ceil(x/4)

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

- Case 2:(values missing in middle)

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

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

- Case 3:(values missing at the end)

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

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

In [34]:

```
# 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])
                ind+=1
            else:
                smoothed_bins.append(0)
        smoothed_regions.extend(smoothed_bins)
    return smoothed_regions
```

In [35]:

```
# Fills a value of zero for every bin where no pickup data is
# present

# the count_values: number picks 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 reigion

    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
```

```
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_limit-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 i in range(i,right_hand_limit+1):
```

In [36]:

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

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

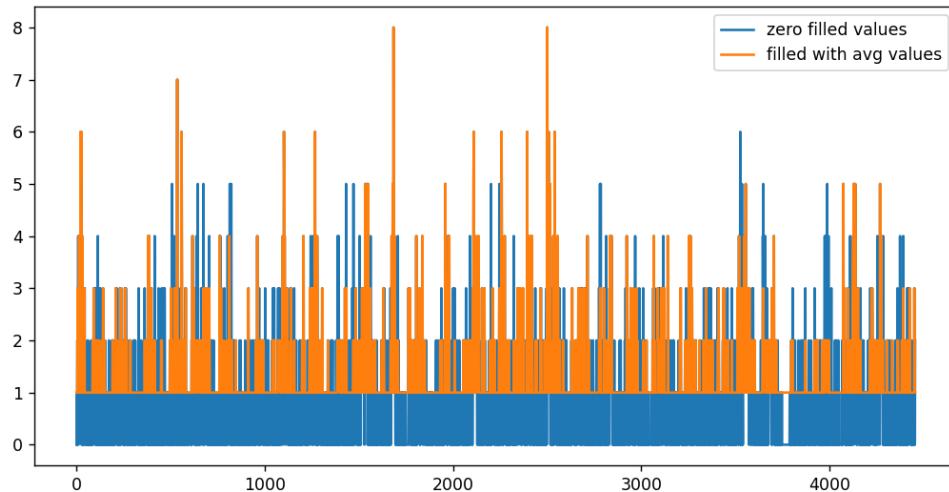
In [37]:

```
# 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 [25]:

```
# 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 []:

```
# 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 of pickups
# that are happened in the first 40min are same in both cases, but
# if you can observe that we looking at the future values
# wheen 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 [38]:

```
# 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 them 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
```

In [27]:

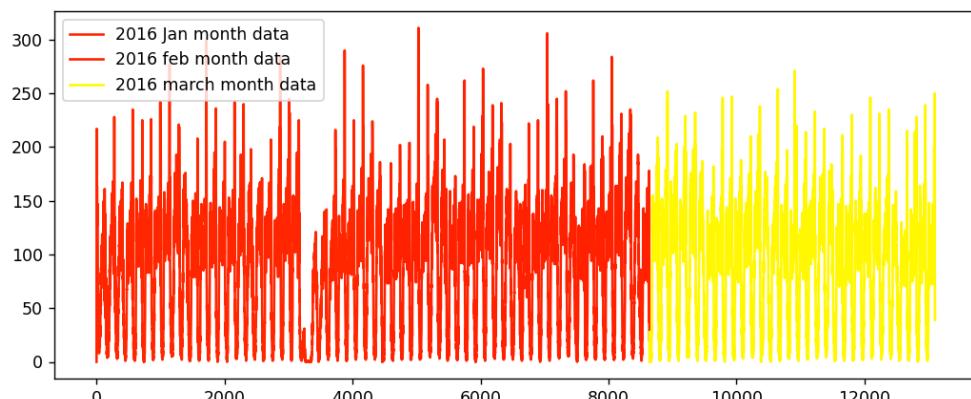
```
import joblib
joblib.dump(jan_2015_smooth, 'jan_2015_smooth.pkl')
joblib.dump(jan_2016_smooth, 'jan_2016_smooth.pkl')
joblib.dump(feb_2016_smooth, 'feb_2016_smooth.pkl')
joblib.dump(mar_2016_smooth, 'mar_2016_smooth.pkl')
```

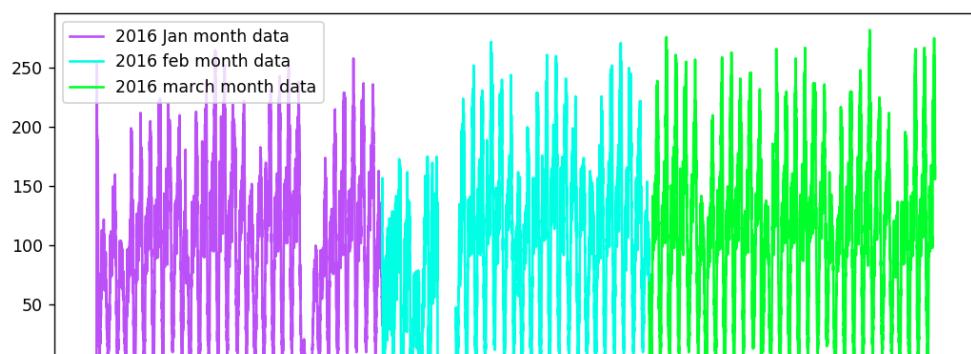
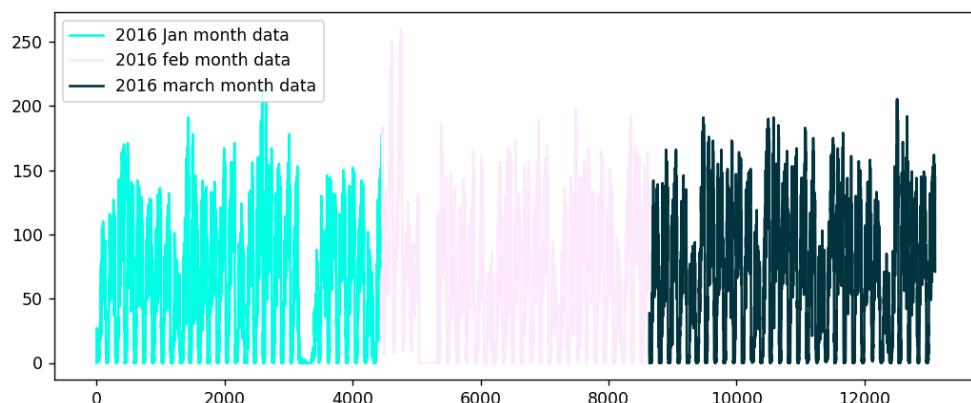
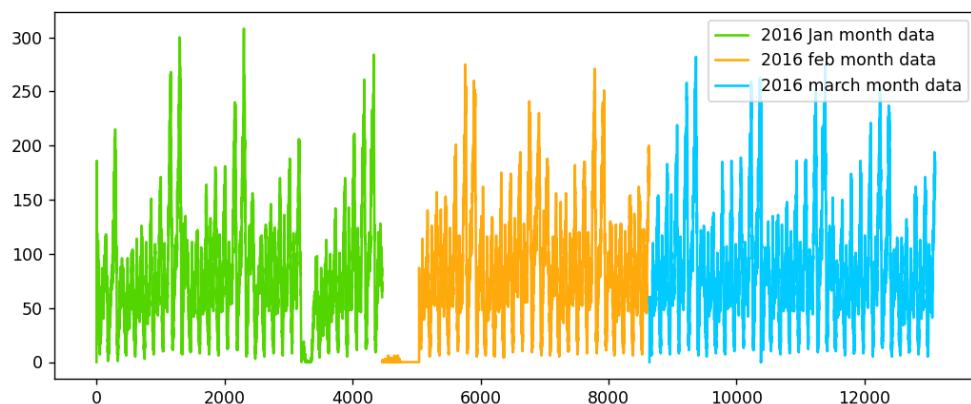
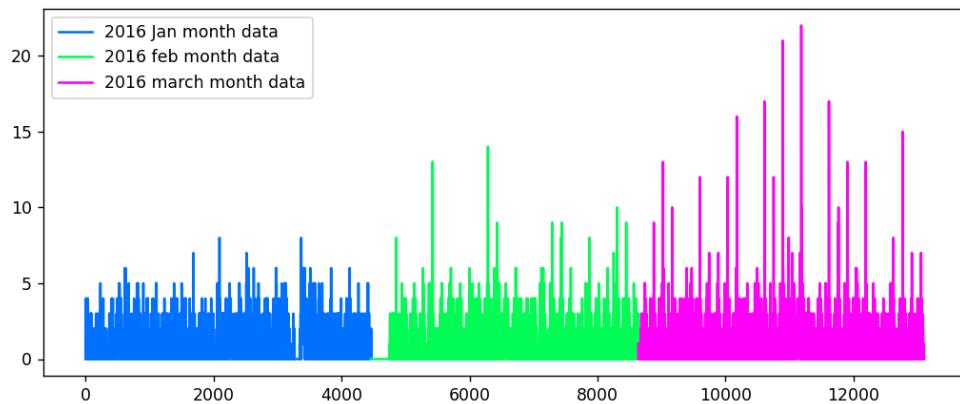
Out[27]: ['mar_2016_smooth.pkl']

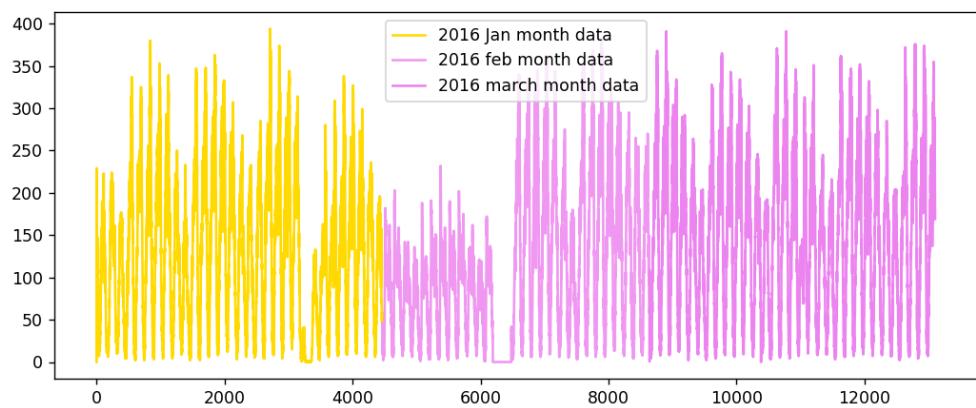
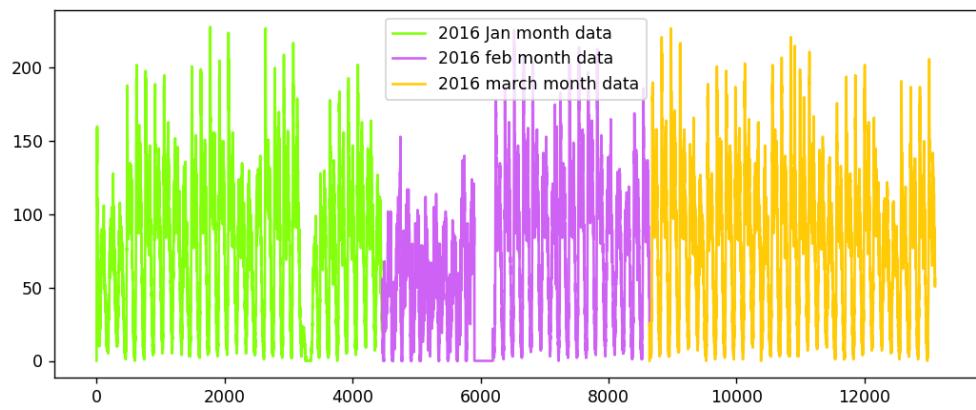
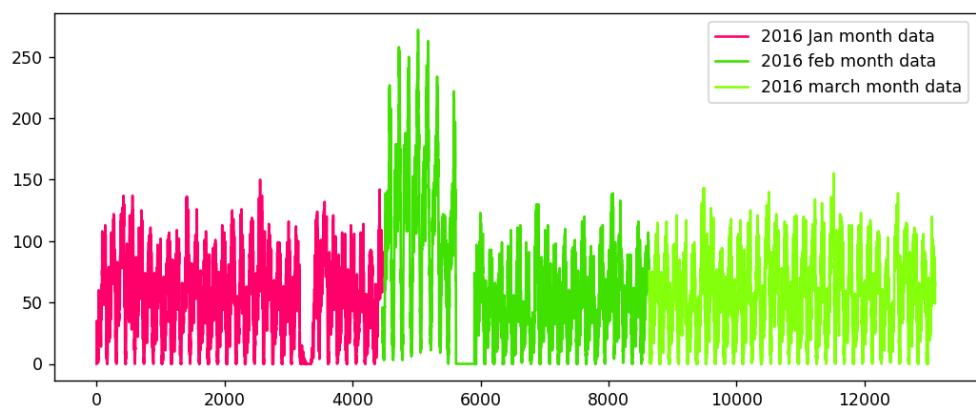
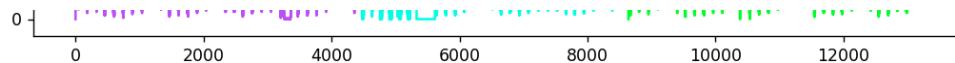
Time series and Fourier Transforms

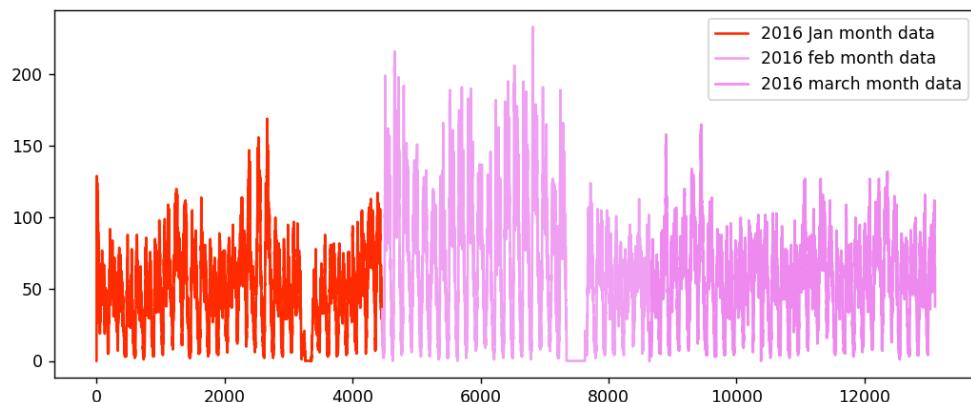
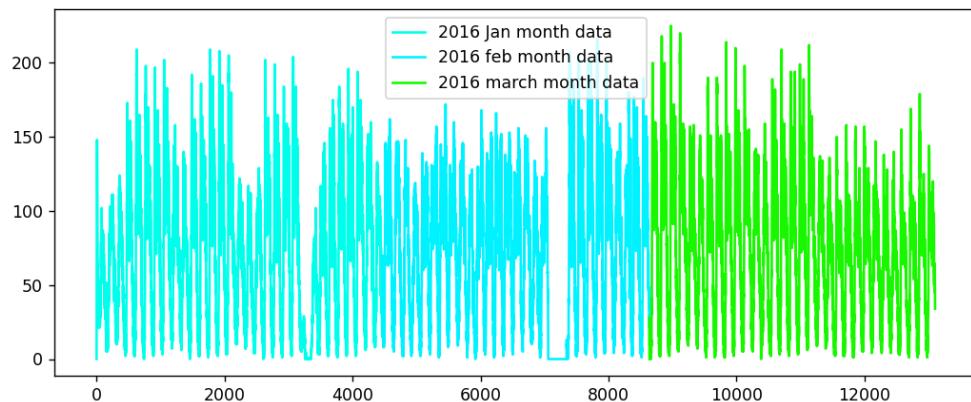
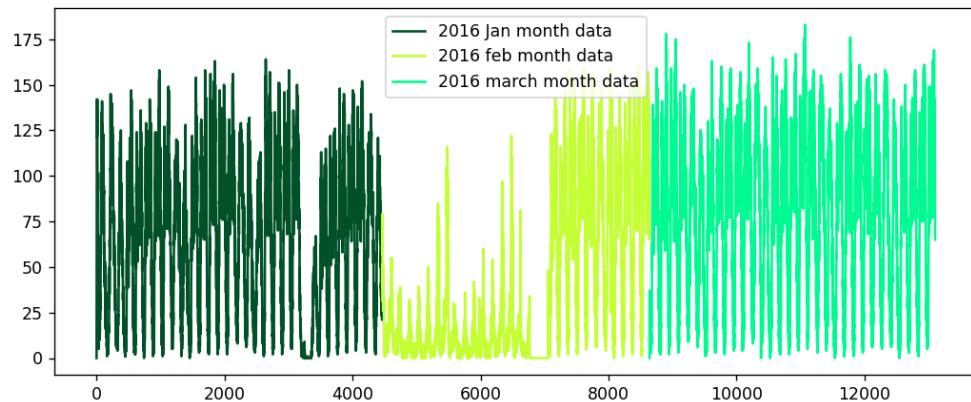
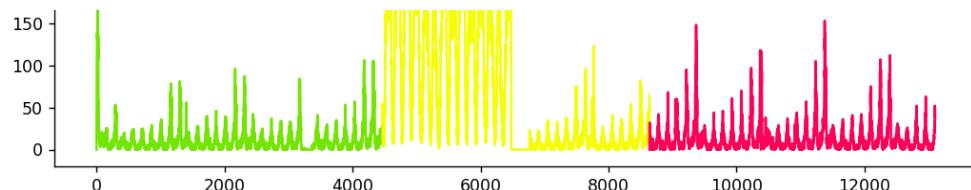
In [39]:

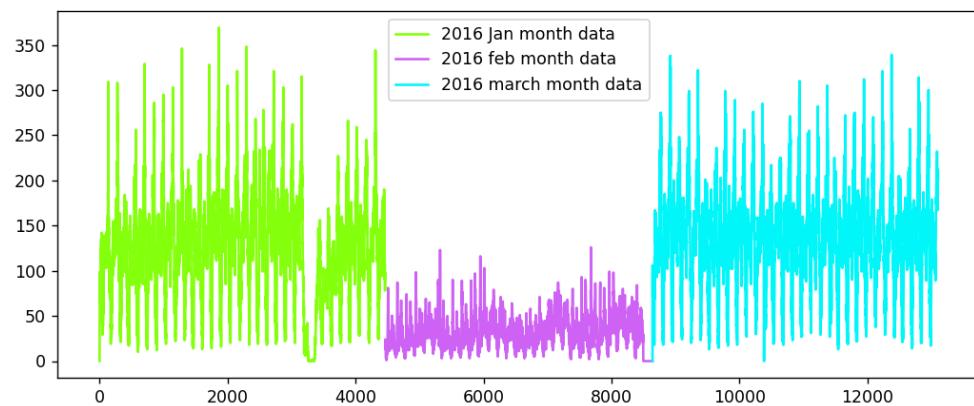
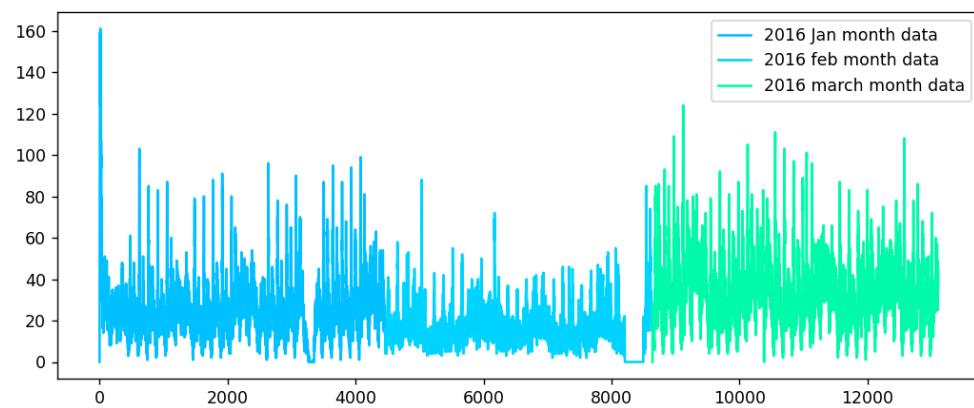
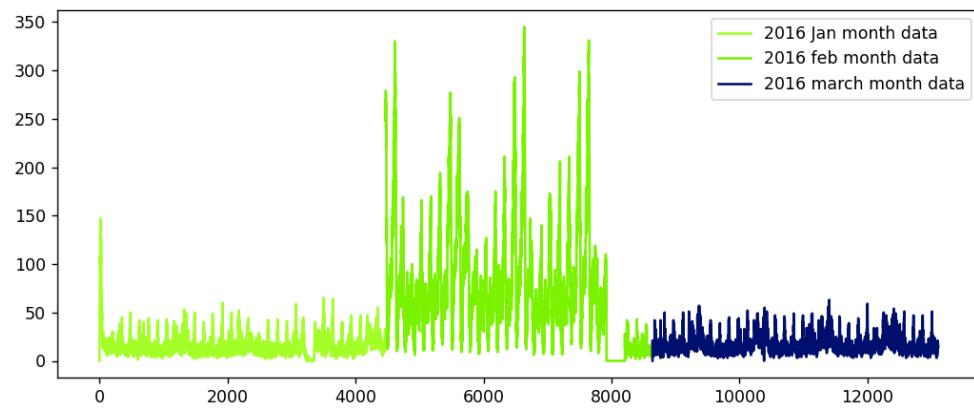
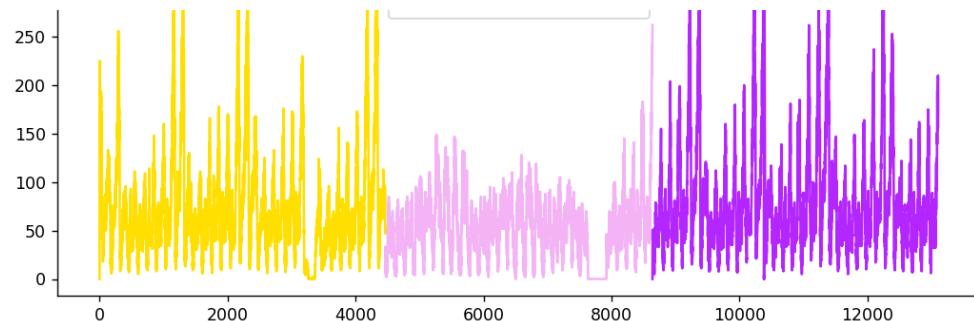
```
def uniqueish_color():
    """There're better ways to generate unique colors, but this
isn't awful."""
    return plt.cm.gist_ncar(np.random.random())
first_x = list(range(0,4464))
second_x = list(range(4464,8640))
third_x = list(range(8640,13104))
for i in range(40):
    plt.figure(figsize=(10,4))
    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()
```

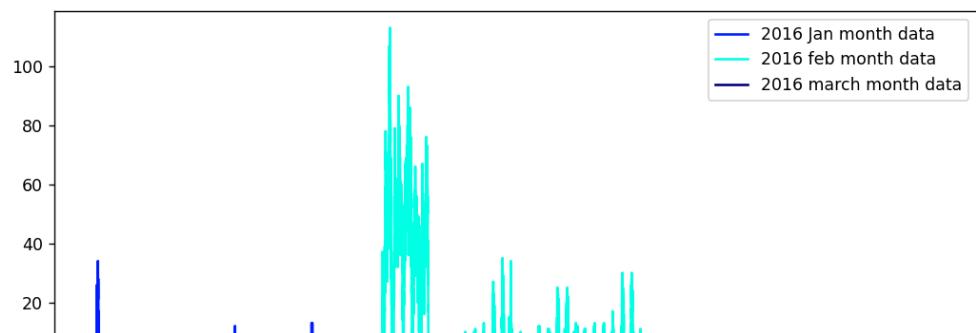
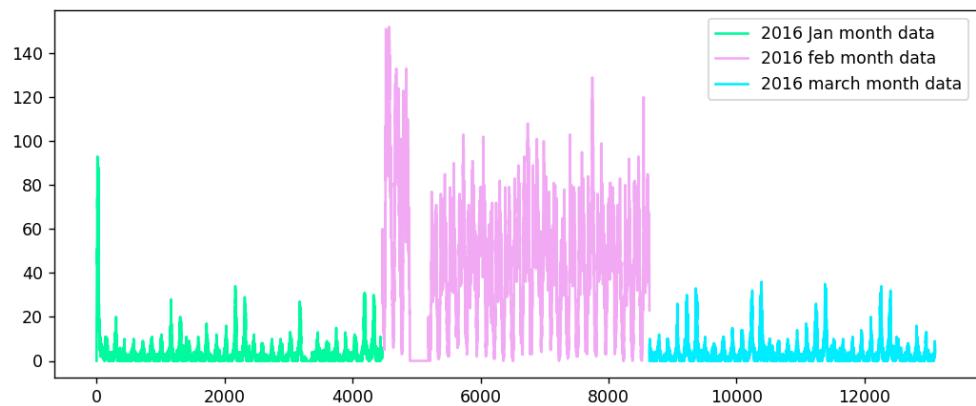
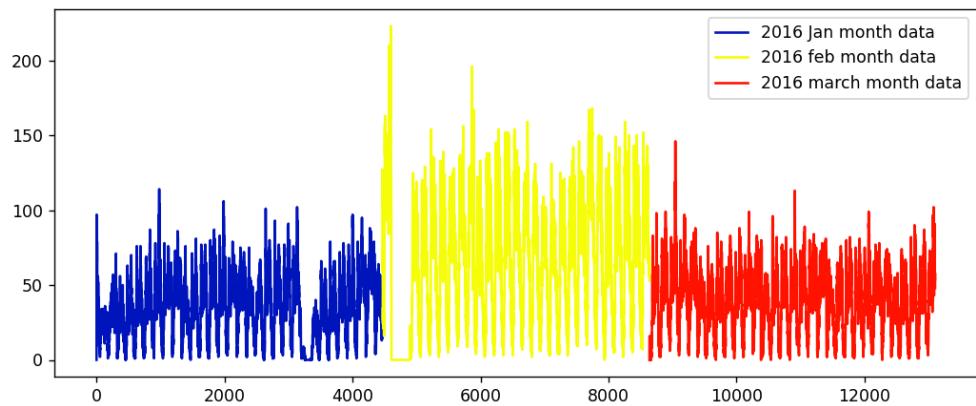
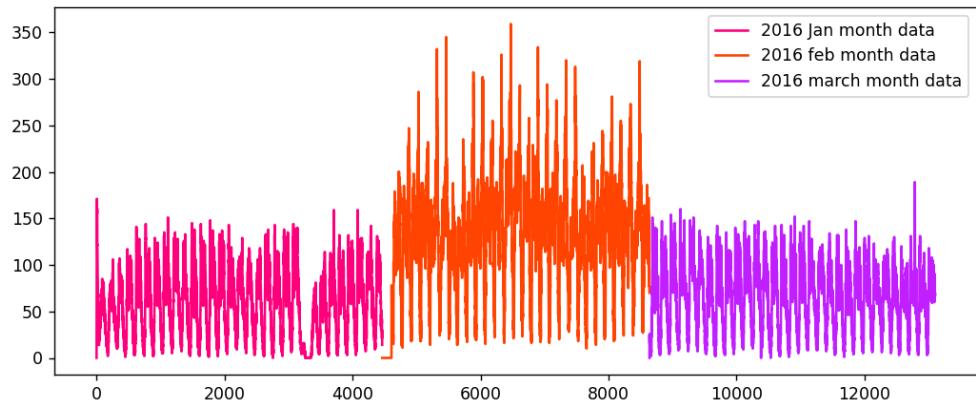


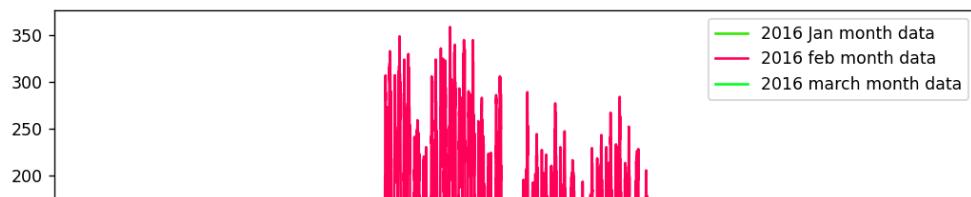
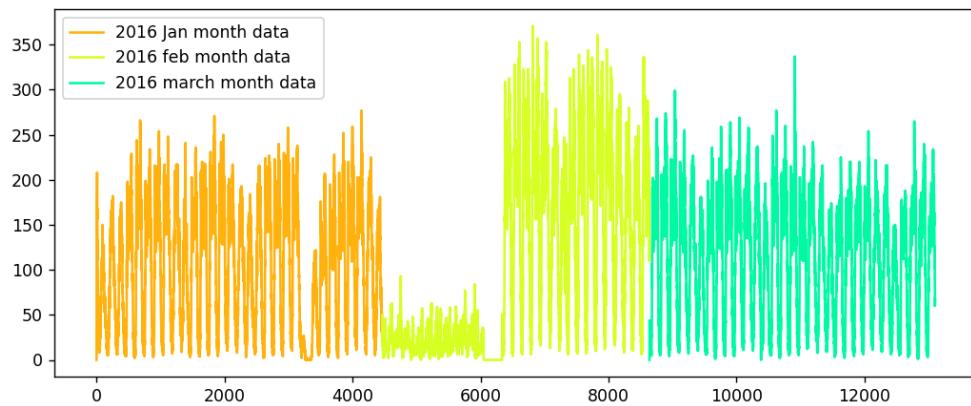
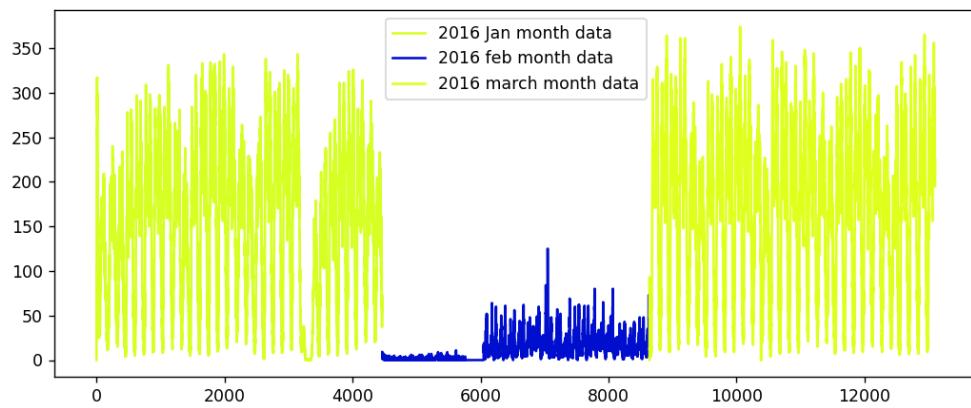
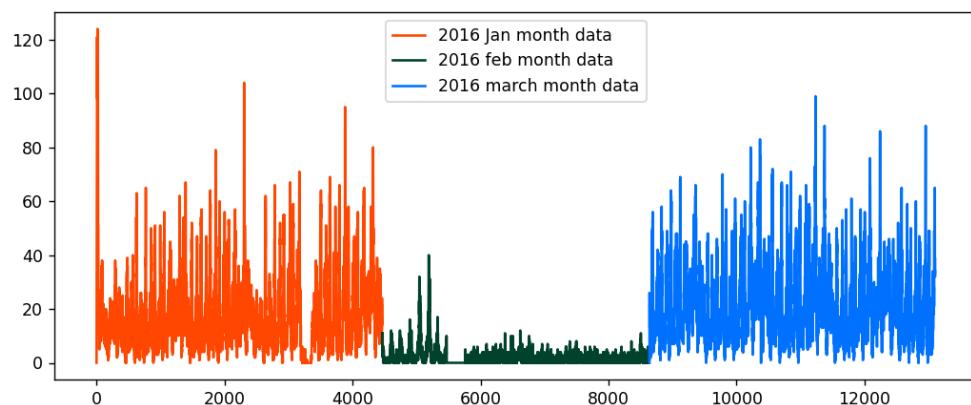
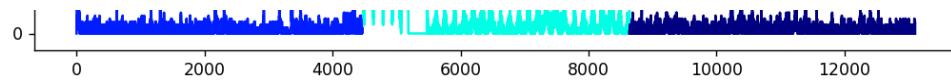


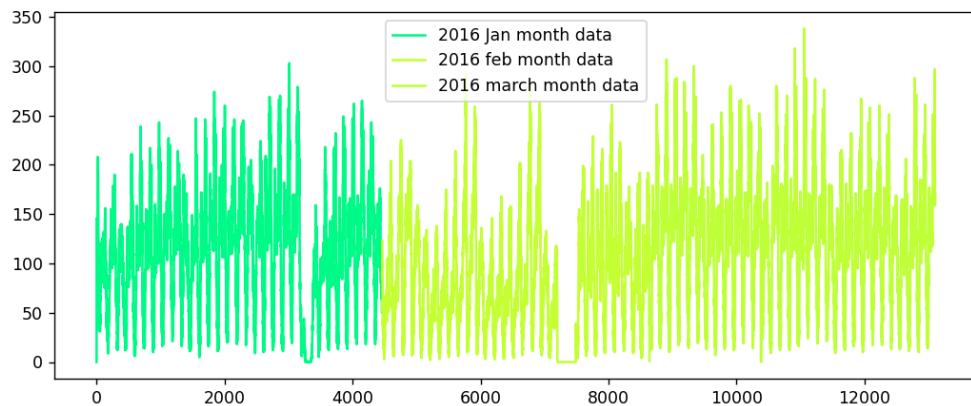
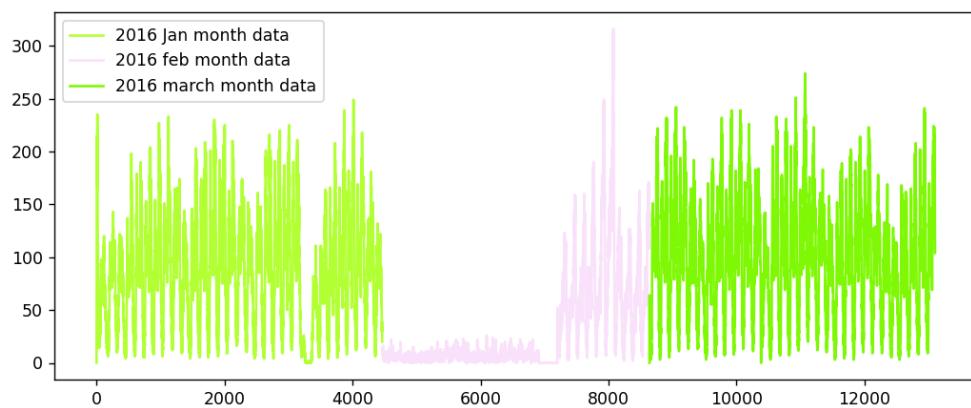
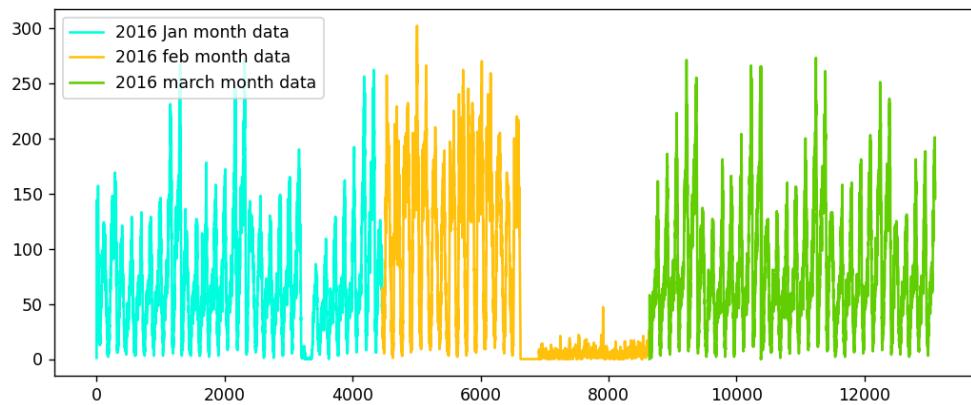
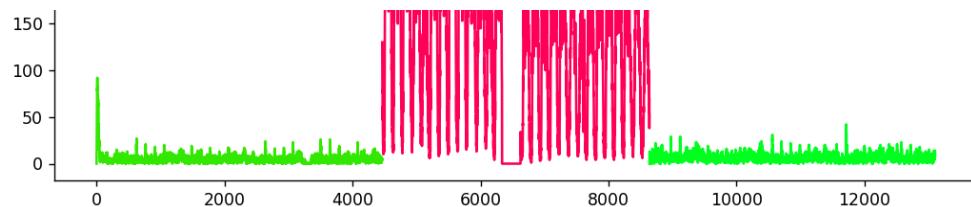


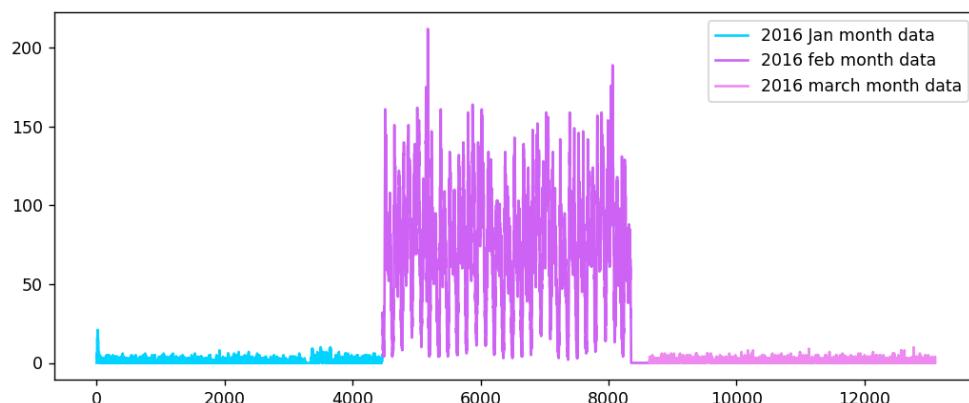
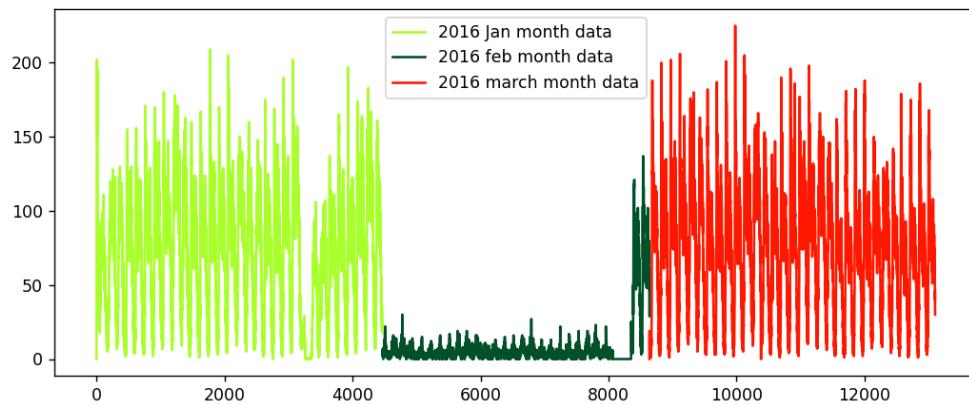
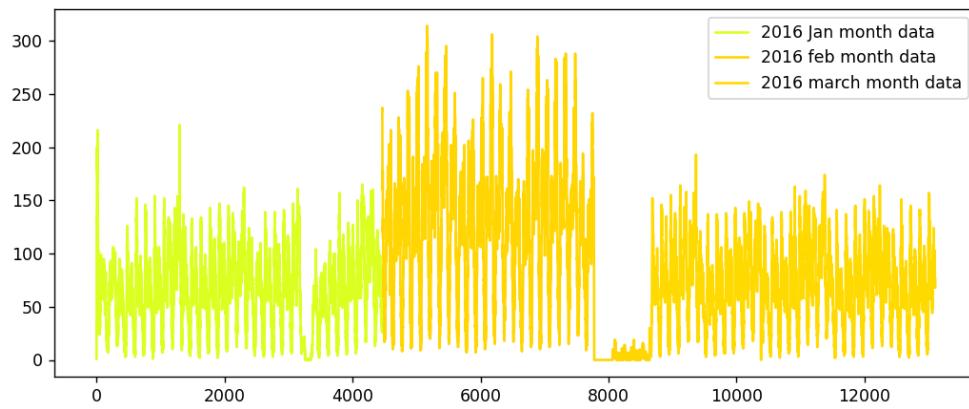
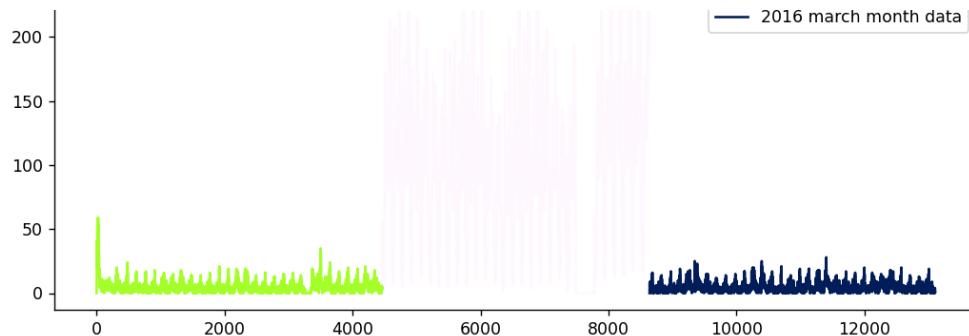


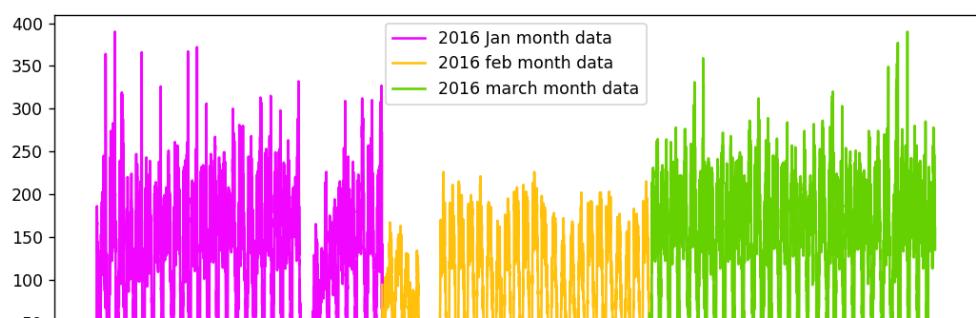
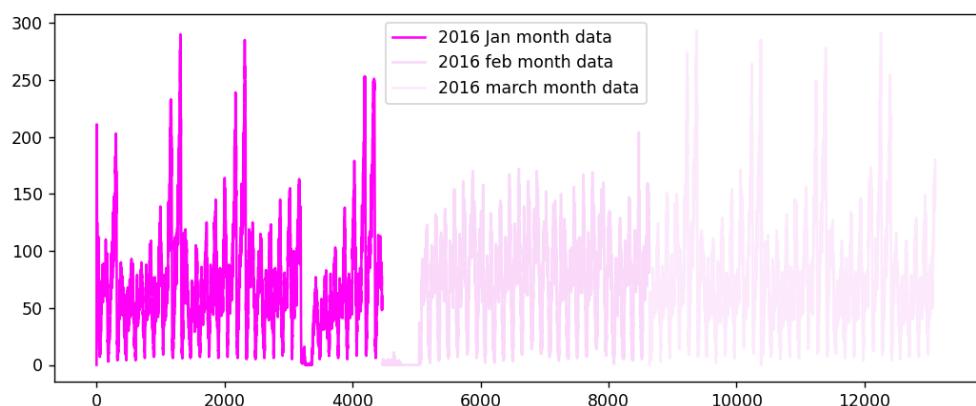
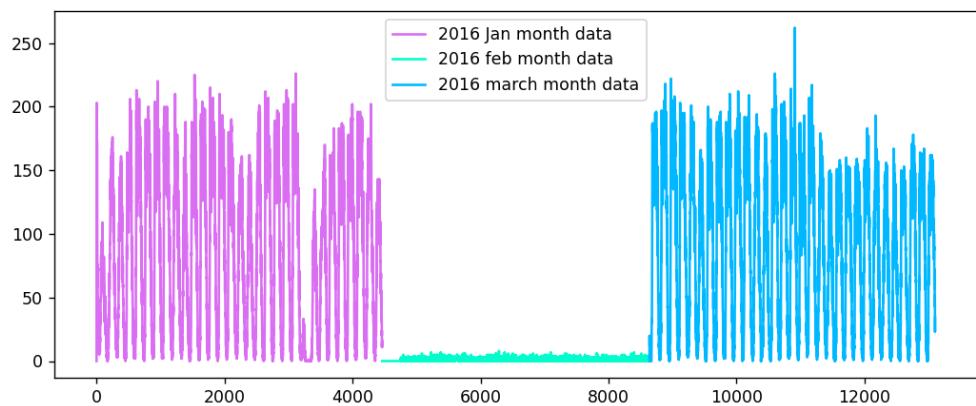
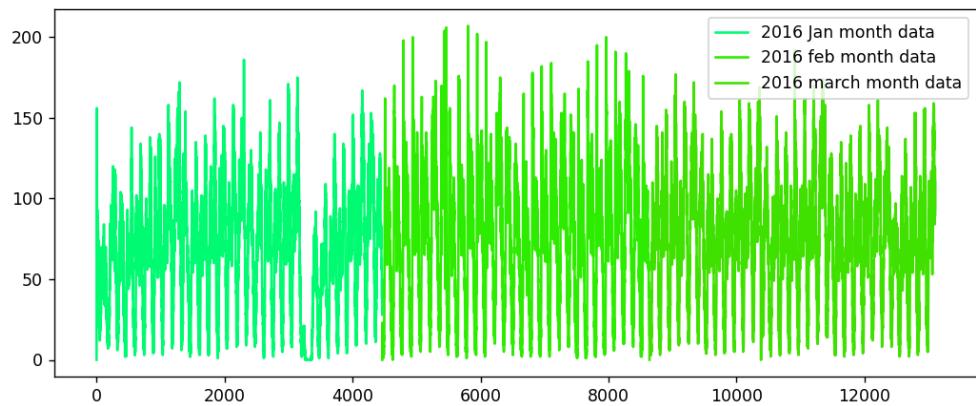


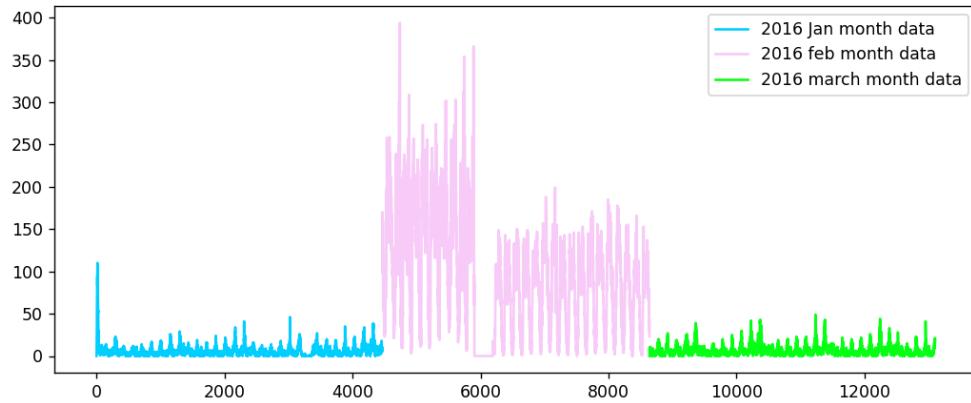
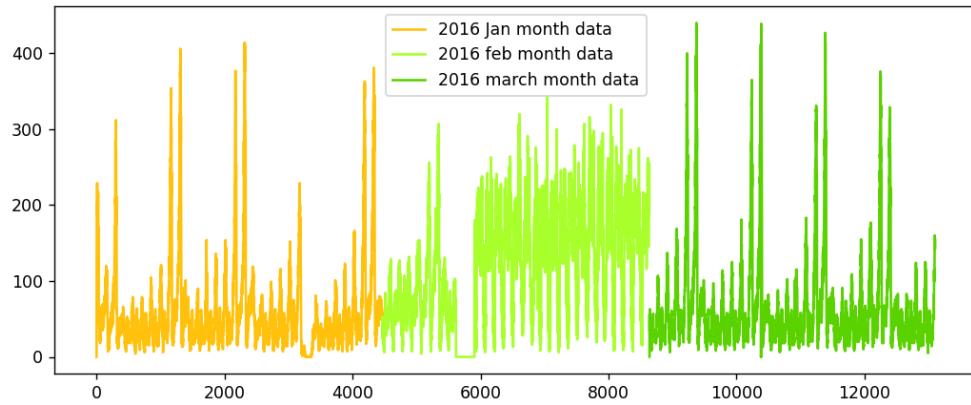
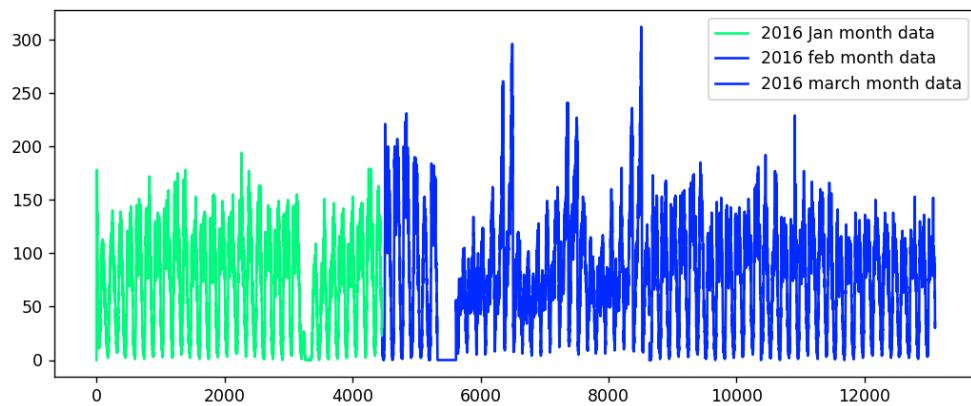
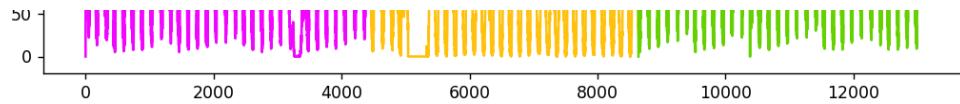


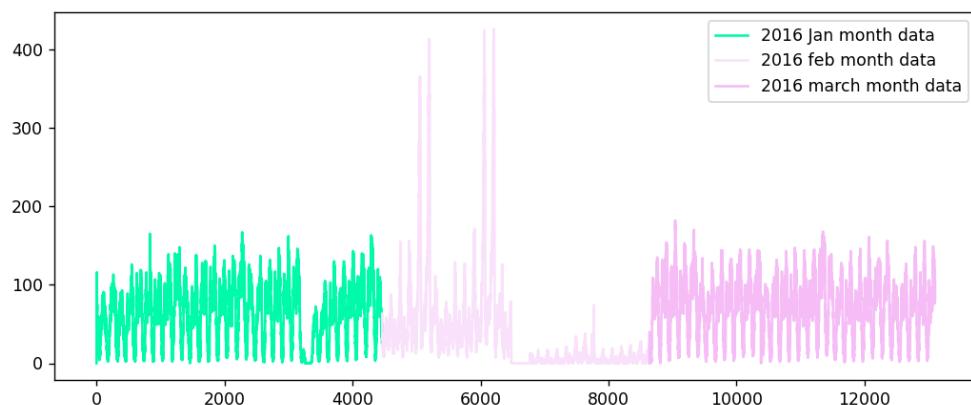
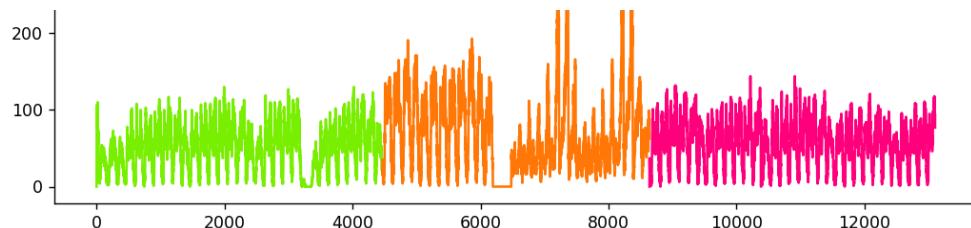






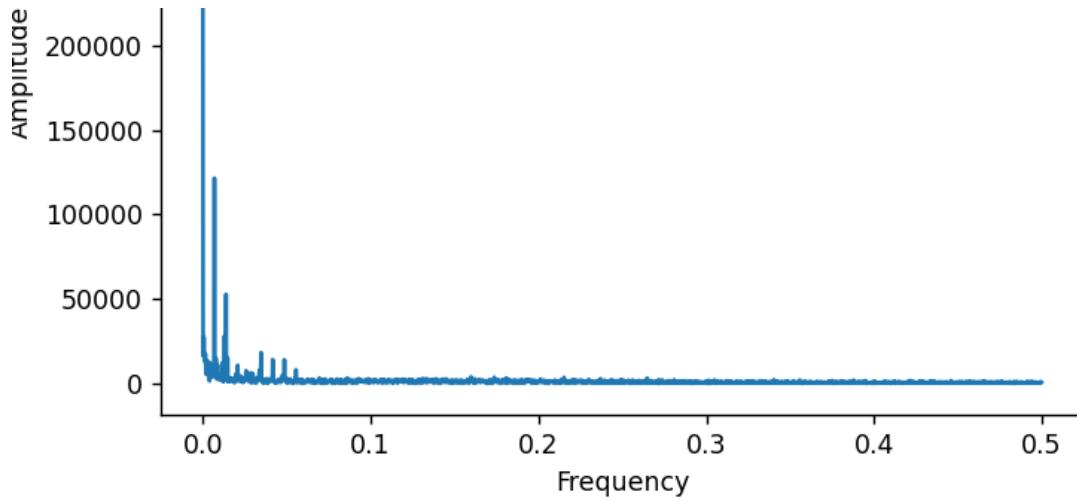






```
In [41]: # getting peaks: https://blog.ytotech.com/2015/11/01/findpeaks-in-python/
# read more about fft function : https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fft.html
Y      = np.fft.fft(np.array(jan_2016_smooth)[0:4460])
# read more about the fftfreq: https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fftfreq.html
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 [37]:

```
#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']
```

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

$$R_t = P_t^{2016}/P_t^{2015} \quad (1)$$

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 -

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

In [38]:

```

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

```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 3 is optimal for getting the best results using Moving Averages using previous Ratio values therefore we get

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

Next we use the Moving averages of the 2016 values itself to predict the future value using

$$P_t = (P_{t-1} + P_{t-2} + P_{t-3} \dots P_{t-n})/n \quad (4)$$

In [39]:

```

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

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

```

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

$$P_t = P_{t-1} \quad (5)$$

Weighted Moving Averages

The Moving Average 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 -

$$R_t = (N * R_{t-1} + (N - 1) * R_{t-2} + (N - 2) * R_{t-3} \dots 1 * R_{t-n}) / (N * (N + 1) / 2) \quad (6)$$

In [40]:

```
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['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['Ratios'].values)[i-window_size+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))
```

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

$$R_t = (5 * R_{t-1} + 4 * R_{t-2} + 3 * R_{t-3} + 2 * R_{t-4} + R_{t-5})/15 \quad (7)$$

Weighted Moving Averages using Previous 2016 Values -

$$P_t = (N * P_{t-1} + (N - 1) * P_{t-2} + (N - 2) * P_{t-3} \dots 1 * P_{t-n})/(N * (N + 1)/2) \quad (8)$$

In [41]:

```

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

$$P_t = (2 * P_{t-1} + P_{t-2})/3 \quad (9)$$

Exponential Weighted Moving Averages

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

$$(\alpha) \quad (10)$$

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

$$\alpha = 0.9 \quad (11)$$

then the number of days on which the value of the current iteration is based is~

$$1/(1 - \alpha) = 10 \quad (12)$$

i.e. we consider values 10 days prior before we predict the value for the current iteration. Also the weights are assigned using

$$2/(N + 1) = 0.18 \quad (13)$$

,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.

$$R'_t = \alpha * R_{t-1} + (1 - \alpha) * R'_{t-1} \quad (14)$$

In [42]:

```

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

```

$$P'_t = \alpha * P_{t-1} + (1 - \alpha) * P'_{t-1} \quad (15)$$

In [43]:

```
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 [44]:

```
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 []:

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

```
Error Metric Matrix (Forecasting Methods) - MAPE & MSE
-----
-----
Moving Averages (Ratios) - MAPE: 0.182115517339
MSE: 400.0625504032258
Moving Averages (2016 Values) - MAPE: 0.14292849687
MSE: 174.84901993727598
-----
-----
Weighted Moving Averages (Ratios) - MAPE: 0.178486925438
MSE: 384.01578741039424
Weighted Moving Averages (2016 Values) - MAPE: 0.135510884362
MSE: 162.46707549283155
-----
-----
Exponential Moving Averages (Ratios) - MAPE: 0.177835501949
MSE: 378.34610215053766
Exponential Moving Averages (2016 Values) - MAPE: 0.135091526367
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:-

$$P'_t = \alpha * P_{t-1} + (1 - \alpha) * P'_{t-1} \quad (16)$$

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 [47]:

```
# 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

# print(len(regions_cum))
# 40
# print(len(regions_cum[0]))
# 12960

# we take number of pickups that are happened in last 5 10min
intravels
number_of_time_stamps = 5

# output varaiable
# it is list of lists
# it will contain number of pickups 13099 for each cluster
output = []

# tsne_lat will contain 13104-5=13099 times lattitude of cluster
center for every cluster
# Ex: [[cent_lat 13099times],[cent_lat 13099times], [cent_lat
13099times].... 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]
```

```
# for every cluster we will be adding 13099 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 (523960, 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 intravel(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]]*13099)
    tsne_lon.append([kmeans.cluster_centers_[i][1]]*13099)
    # 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 5 pickup bins
    tsne_weekday.append([int(((int(k/144))%7+4)%7) for k in
range(5,4464+4176+4464)])
    # 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 lsits]
    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][5:])
tsne_feature = tsne_feature[1:]
```

```
In [48]: len(tsne_lat[0])*len(tsne_lat) == tsne_feature.shape[0] ==
len(tsne_weekday)*len(tsne_weekday[0]) == 40*13099 ==
len(output)*len(output[0])
```

```
Out[48]: True
```

In [49]:

```
# Getting the predictions of exponential moving averages to be used
# as a feature in cumulative form

# upto now we computed 8 features for every data point that starts
# from 50th min of the day
# 1. cluster center lattitude
# 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
# avarage gives us the best error
# we will try to add the same exponential weighted moving avarage
# 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]
```

```
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[5:])
predicted_values=[]
```

In []:

In [50]:

```
# train, test split : 70% 30% 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
print("size of train data :", int(13099*0.7))
print("size of test data :", int(13099*0.3))
```

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

In [51]:

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

```
In [52]:  
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 9169 Each data point contains 5 features
Number of data clusters 40 Number of data points in test data 3930 Each data p oint contains 5 features

```
In [53]:  
# extracting first 9169 timestamp values i.e 70% of 13099 (total  
timestamps) for our training data  
tsne_train_flat_lat = [i[:9169] for i in tsne_lat]  
tsne_train_flat_lon = [i[:9169] for i in tsne_lon]  
tsne_train_flat_weekday = [i[:9169] for i in tsne_weekday]  
tsne_train_flat_output = [i[:9169] for i in output]  
tsne_train_flat_exp_avg = [i[:9169] for i in predict_list]
```

```
In [54]:  
# extracting the rest of the timestamp values i.e 30% of 12956  
(total timestamps) for our test data  
tsne_test_flat_lat = [i[9169:] for i in tsne_lat]  
tsne_test_flat_lon = [i[9169:] for i in tsne_lon]  
tsne_test_flat_weekday = [i[9169:] for i in tsne_weekday]  
tsne_test_flat_output = [i[9169:] for i in output]  
tsne_test_flat_exp_avg = [i[9169:] for i in predict_list]
```

```
In [55]:  
# 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 [56]:

```
# 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 [57]:

```
# 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_output = sum(tsne_test_flat_output, [])
tsne_test_exp_avg = sum(tsne_test_flat_exp_avg, [])
```

In [58]:

```
# Preparing the data frame for our train data
columns = ['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)
```

(366760, 9)

In [59]:

```
# 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)
```

(157200, 9)

In [60]:

```
df_test.head()
```

Out[60]:

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

Using Linear Regression

In []:

```
# find more about LinearRegression function here http://scikit-  
learn.org/stable/modules/generated  
/sklearn.linear_model.LinearRegression.html  
# -----  
# default parameters  
# 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-intuition-1-2-copy-8/  
# -----  
  
from sklearn.linear_model import LinearRegression  
lr_reg=LinearRegression().fit(df_train, tsne_train_output)  
  
y_pred = lr_reg.predict(df_test)  
lr_test_predictions = [round(value) for value in y_pred]  
y_pred = lr_reg.predict(df_train)  
lr_train_predictions = [round(value) for value in y_pred]
```

Using Random Forest Regressor

In []:

```
# 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 parameters

# 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/
# -----


regrl =
RandomForestRegressor(max_features='sqrt', min_samples_leaf=4, min_samples_
n_jobs=-1)
regrl.fit(df_train, tsne_train_output)
```

Out[]: RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=None,
max_features='sqrt', max_leaf_nodes=None,

```
min_impurity_decrease=0.0, min_impurity_split=None,
min_samples_leaf=4, min_samples_split=3,
min_weight_fraction_leaf=0.0, n_estimators=40, n_jobs=-1,
oob_score=False, random_state=None, verbose=0, warm_start=False)
```

```
In [ ]: # Predicting on test data using our trained random forest model
```

```
# the models regr1 is already hyper parameter tuned
# the parameters that we got above are found using grid search

y_pred = regr1.predict(df_test)
rndf_test_predictions = [round(value) for value in y_pred]
y_pred = regr1.predict(df_train)
rndf_train_predictions = [round(value) for value in y_pred]
```

```
In [ ]: #feature importances based on analysis using random forest
```

```
print (df_train.columns)
print (regr1.feature_importances_)
```

```
Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'lat', 'lon', 'weekday',
       'exp_avg'],
      dtype='object')
[ 0.00477243  0.07614745  0.14289548  0.1857027   0.23859285  0.00227886
  0.00261956  0.00162121  0.34536947]
```

Using XgBoost Regressor

In []:

```
# 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\_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, subsample=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, verbose=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 not 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-ensembles/  
# -----  
  
x_model = xgb.XGBRegressor(  
    learning_rate =0.1,  
    n_estimators=1000,  
    max_depth=3,  
    min_child_weight=3,  
    gamma=0,  
    subsample=0.8,
```

```
Out[ ]: XGBRegressor(base_score=0.5, colsample_bylevel=1, colsample_bytree=0.8,
                     gamma=0, learning_rate=0.1, max_delta_step=0, max_depth=3,
                     min_child_weight=3, missing=None, n_estimators=1000, nthread=4,
                     objective='reg:linear', reg_alpha=200, reg_lambda=200,
                     scale_pos_weight=1, seed=0, silent=True, subsample=0.8)
```

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

```
In [ ]: #feature importances
x_model.booster().get_score(importance_type='weight')
```

```
Out[ ]: {'exp_avg': 806,
         'ft_1': 1008,
         'ft_2': 1016,
         'ft_3': 863,
         'ft_4': 746,
         'ft_5': 1053,
         'lat': 602,
         'lon': 612,
         'weekday': 195}
```

Calculating the error metric values for various models

In []:

```
train_mape = []
test_mape = []

train_mape.append((mean_absolute_error(tsne_train_output, df_train['ft_1'].values) / 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 []:

```
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: 0.140052758787	Test: 0.136531257048
Exponential Averages Forecasting -	Train: 0.13289968436	Test: 0.129361804204
Linear Regression -	Train: 0.13331572016	Test: 0.129120299401
Random Forest Regression -	Train: 0.0917619544199	Test: 0.127244647137
XgBoost Regression -	Train: 0.129387355679	Test: 0.126861699078

Assignments

In []:

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

Incorporate Fourier features as features into Regression models
and measure MAPE.

In [66]:

```
# https://stackoverflow.com/questions/3694918/how-to-extract-
# frequency-associated-with-fft-values-in-python
# https://github.com/jinalsalvi/NYC-Taxi-Demand-Prediction
# blob/master/NYC%20Final.ipynb
# Code Ref.: https://github.com/pranaysawant/New-York-Taxi-Demand-
# Prediction/blob/master/NYC_Final.ipynb

fourier_features = pd.DataFrame(columns=
['f_1','a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5'])

for region in range(0,40):
    df_jan = pd.DataFrame()
    df_feb = pd.DataFrame()
    df_mar = pd.DataFrame()

    ampJan = np.fft.fft(np.array(regions_cum[region][0:4464]))
    freqJan = np.fft.fftfreq((4464), 1)
    df_jan['Frequency'] = freqJan
    df_jan['Amplitude'] = ampJan

    ampFeb = np.fft.fft(np.array(regions_cum[region])[4464:(4176+4464)])
    freqFeb = np.fft.fftfreq((4176), 1)
    df_feb['Frequency'] = freqFeb
    df_feb['Amplitude'] = ampFeb

    ampMar = np.fft.fft(np.array(regions_cum[region])[(4176+4464):(4176+4464+4464)])
    freqMar = np.fft.fftfreq((4464), 1)
    df_mar['Frequency'] = freqMar
    df_mar['Amplitude'] = ampMar

    list_jan = []
    list_feb = []
    list_mar = []
```

```
mar_sorted = df_mar.sort_values(by=["Amplitude"], ascending=False) [:5].reset_index(drop=True).T

for i in range(0,5):
    list_jan.append(float(jan_sorted[i]['Frequency']))
    list_jan.append(float(jan_sorted[i]['Amplitude']))

    list_feb.append(float(feb_sorted[i]['Frequency']))
    list_feb.append(float(feb_sorted[i]['Amplitude']))

    list_mar.append(float(mar_sorted[i]['Frequency']))
    list_mar.append(float(mar_sorted[i]['Amplitude']))

frame_jan = pd.DataFrame([list_jan]*4464)
frame_feb = pd.DataFrame([list_feb]*4176)
frame_mar = pd.DataFrame([list_mar]*4464)

frame_jan.columns =
['f_1','a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5',]
frame_feb.columns =
['f_1','a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5',]
frame_mar.columns =
['f_1','a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5',]

fourier_features = fourier_features.append(frame_jan,
ignore_index=True)
fourier_features = fourier_features.append(frame_feb,
ignore_index=True)
fourier_features = fourier_features.append(frame_mar,
ignore_index=True)

for i in range(0,13104):
    if i==0:
        predicted_value= reasions_cum[reasions][0]
```

```
fourier_features = fourier_features.fillna(0)
```

In [78]:

```
# Code Ref.: https://github.com/pranaysawant/New-York-Taxi-Demand-Prediction/blob/master/NYC_Final.ipynb

final_fourier_train = pd.DataFrame(columns=['a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5'])
final_fourier_test = pd.DataFrame(columns=['a_1','f_2','a_2','f_3','a_3','f_4','a_4','f_5','a_5'])

# extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our training data
for i in range(0,40):
    final_fourier_train =
final_fourier_train.append(fourier_features[i*13099:(13099*i+9169)])
    final_fourier_test =
final_fourier_test.append(fourier_features[(13099*(i))+9169:13099*(i+1)])
final_fourier_train.reset_index(inplace=True)
final_fourier_test.reset_index(inplace=True)
```

In [79]:

```
print("fourier_train shape",final_fourier_train.shape)
print("fourier_test shape",final_fourier_test.shape)
final_fourier_train.head()
```

```
fourier_train shape (366760, 10)
fourier_test shape (157200, 10)
```

Out[79]:

index	a_1	f_2	a_2	f_3	a_3	f_4	a_4	f_5
-------	-----	-----	-----	-----	-----	-----	-----	-----

	index	a_1	f_2	a_2	f_3	a_3	f_4	a_4	f_
0	0	369774.0	-0.012993	24998.122651	0.012993	24998.122651	0.000448	15434.851794	-0.00044
1	1	369774.0	-0.012993	24998.122651	0.012993	24998.122651	0.000448	15434.851794	-0.00044
2	2	369774.0	-0.012993	24998.122651	0.012993	24998.122651	0.000448	15434.851794	-0.00044
3	3	369774.0	-0.012993	24998.122651	0.012993	24998.122651	0.000448	15434.851794	-0.00044

In [95]:

```
final_train = pd.concat([df_train, final_fourier_train], axis=1)
final_test = pd.concat([df_test, final_fourier_test], axis=1)

final_train.drop(columns = ['index'], inplace = True)
final_test.drop(columns = ['index'], inplace = True)

print("final train shape",final_train.shape)
print("final test shape",final_test.shape)

final_train.head()
```

final train shape (366760, 18)
final test shape (157200, 18)

Out[95]:

	ft_5	ft_4	ft_3	ft_2	ft_1	lat	lon	weekday	exp_avg	a_1	f_2
0	0	63	217	189	137	40.776228	-73.982119	4	150	369774.0	-0.012993
1	63	217	189	137	135	40.776228	-73.982119	4	139	369774.0	-0.012993
2	217	189	137	135	129	40.776228	-73.982119	4	132	369774.0	-0.012993
3	189	137	135	129	150	40.776228	-73.982119	4	144	369774.0	-0.012993
4	137	135	129	150	164	40.776228	-73.982119	4	158	369774.0	-0.012993

Linear Regression: Grid Search

In [71]:

```
# importing libraries

from sklearn.linear_model import LinearRegression

from sklearn.model_selection import GridSearchCV
```

In [101...]

```
# https://scikit-learn.org/stable/modules/generated
/sklearn.linear_model.LinearRegression.html
lr = LinearRegression()

params = {
    'fit_intercept' : [True, False],
    'normalize': [True, False],
    'copy_X': [True, False],
}

# https://scikit-learn.org/stable/modules/generated
/sklearn.model_selection.GridSearchCV.html
model = GridSearchCV(lr, param_grid= params, cv = 5, n_jobs=-1,
                      return_train_score= True, verbose= 5, scoring=
'neg_mean_squared_error')
model.fit(final_train, tsne_train_output)
```

Fitting 5 folds for each of 8 candidates, totalling 40 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

[Parallel(n_jobs=-1)]: Done 2 tasks | elapsed: 3.2s

[Parallel(n_jobs=-1)]: Done 34 out of 40 | elapsed: 44.5s remaining: 7.8s

[Parallel(n_jobs=-1)]: Done 40 out of 40 | elapsed: 53.0s finished

Out[101...]: GridSearchCV(cv=5, estimator=LinearRegression(), n_jobs=-1,
param_grid={'copy_X': [True, False],
 'fit_intercept': [True, False],
 'normalize': [True, False]},
return_train_score=True, scoring='neg_mean_absolute_error',
verbose=5)

In [102...]

```
print("Estimators: ", model.best_estimator_)
print("Params: ", model.best_params_)
```

Estimators: LinearRegression()

Params: {'copy_X': True, 'fit_intercept': True, 'normalize': False}

```
In [103...]: model = LinearRegression(copy_X= True, fit_intercept = True,
normalize= False)
model.fit(final_train, tsne_train_output)

y_pred = model.predict(final_train)
y_pred_train = [round(value) for value in y_pred]

y_pred = model.predict(final_test)
y_pred_test = [round(value) for value in y_pred]
```

```
In [104...]: train_error = (mean_absolute_error(tsne_train_output,
y_pred_train)) / (sum(tsne_train_output) / len(tsne_train_output))

test_error = (mean_absolute_error(tsne_test_output,
y_pred_test)) / (sum(tsne_test_output) / len(tsne_test_output))

print ("Linear Regression - Train: ",train_error," Test:
",test_error)
```

Linear Regression - Train: 0.13324598902941232 Test: 0.12889723956282564

Random Forest: Random Search

```
In [105...]: # importing libraries
from sklearn.model_selection import RandomizedSearchCV
from sklearn.ensemble import RandomForestRegressor
```

In [109...]

```
# https://analyticsindiamag.com/guide-to-hyperparameters-tuning-
# using-gridsearchcv-and-randomizedsearchcv/
depth = [1, 10, 50, 100, 500, 1000]
n_estimators = [100, 200, 300, 400, 500, 1000]
min_samples_split = [3, 5, 7, 9]

tuned_parameters = {
    'max_depth': depth,
    'n_estimators': n_estimators,
    'min_samples_split': min_samples_split,
}

# Applying RandomizedSearchCV with k folds = 5 and taking
# 'mean_absolute_error' as score metric
clf = RandomizedSearchCV(RandomForestRegressor(), 
param_distributions=tuned_parameters, scoring =
'neg_mean_squared_error',
cv=5, return_train_score=True, verbose =
5, n_jobs=-1)
clf.fit(final_train, tsne_train_output)
```

Fitting 5 folds for each of 10 candidates, totalling 50 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

[Parallel(n_jobs=-1)]: Done 2 tasks | elapsed: 1.3min

[Parallel(n_jobs=-1)]: Done 46 out of 50 | elapsed: 83.3min remaining: 7.2m
in

[Parallel(n_jobs=-1)]: Done 50 out of 50 | elapsed: 114.6min finished

Out[109...]: RandomizedSearchCV(cv=5, estimator=RandomForestRegressor(), n_jobs=-1,
param_distributions={'max_depth': [1, 10, 50, 100, 500,
1000],
'min_samples_split': [3, 5, 7, 9],
'n_estimators': [100, 200, 300, 400,
500, 1000]},
return_train_score=True, scoring='neg_mean_absolute_error',
verbose=5)

In [110...]

```
print('Best hyper parameter: ', clf.best_params_)
print('Model estimator: ', clf.best_estimator_)
```

Best hyper parameter: {'n_estimators': 400, 'min_samples_split': 9, 'max_depth': 10}

Model estimator: RandomForestRegressor(max_depth=10, min_samples_split=9, n_estimators=400)

In [121...]

```
model = RandomForestRegressor(max_features='sqrt', min_samples_leaf=4,
                             min_samples_split=9, n_estimators=1000, n_jobs=-1)

model.fit(final_train, tsne_train_output)

y_pred = model.predict(final_train)
y_pred_train = [round(value) for value in y_pred]

y_pred = model.predict(final_test)
y_pred_test = [round(value) for value in y_pred]
```

In [122...]

```
train_error = (mean_absolute_error(tsne_train_output,
y_pred_train))/(sum(tsne_train_output)/len(tsne_train_output))

test_error = (mean_absolute_error(tsne_test_output,
y_pred_test))/(sum(tsne_test_output)/len(tsne_test_output))

print ("Random Forest Regressor - Train: ",train_error," Test: ",
test_error)
```

Random Forest Regressor - Train: 0.0934906831978537 Test: 0.125479091246191
53

In [123...]

```
#feature importances based on analysis using random forest
print (final_train.columns)
print (model.feature_importances_)

Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'lat', 'lon', 'weekday',
       'exp_avg', 'a_1', 'f_2', 'a_2', 'f_3', 'a_3', 'f_4', 'a_4', 'f_5',
       'a_5'],
      dtype='object')
[0.07467537 0.09424847 0.13697664 0.16205918 0.21850126 0.00113133
 0.00121462 0.00127885 0.26496946 0.01956313 0.00061349 0.00480316
 0.00067343 0.00388331 0.00061237 0.00672539 0.00061124 0.00745931]
```

Xgboost: Random Search

In [124...]

```
# importing libraries
from sklearn.model_selection import RandomizedSearchCV
from xgboost import XGBRegressor
```

In [152...]

```
depth = [1, 2, 3, 5, 7]
n_estimators = [100, 200, 300, 400, 500, 1000]
learning_rate = [0.01, 0.05, 0.1, 0.2, 0.3]
tuned_parameters = {
    'max_depth':depth,
    'n_estimators': n_estimators,
    'learning_rate': learning_rate,
}

# enable gpu, https://xgboost.readthedocs.io/en/latest
#gpu/index.html
model = XGBRegressor(tree_method='gpu_hist', gpu_id = 0)

# Applying RandomizedSearchCV with k folds = 5 and taking
'mean_absolute_error' as score metric
clf = RandomizedSearchCV(model, param_distributions=
tuned_parameters, scoring = 'neg_mean_squared_error',
                           cv=5, return_train_score= True, verbose =
5, n_jobs= -1)
clf.fit(final_train, tsne_train_output)

IPython.display.Audio("notify.mp3", autoplay = True)
```

```
Fitting 5 folds for each of 10 candidates, totalling 50 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done   2 tasks      | elapsed:   6.4s
[Parallel(n_jobs=-1)]: Done   46 out of  50 | elapsed: 20.1min remaining:  1.7m
in
[Parallel(n_jobs=-1)]: Done   50 out of  50 | elapsed: 25.0min finished
```

Out[152...]



0:00 / 0:17



```
In [154...]  
print('Best hyper parameter: ', clf.best_params_)  
print('Model estimator: ', clf.best_estimator_)  
  
Best hyper parameter: {'n_estimators': 500, 'max_depth': 5, 'learning_rate': 0.01}  
Model estimator: XGBRegressor(base_score=0.5, booster='gbtree', colsample_bytree=1,  
                                colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=0,  
                                importance_type='gain', interaction_constraints='',  
                                learning_rate=0.01, max_delta_step=0, max_depth=5,  
                                min_child_weight=1, missing=nan, monotone_constraints='()',  
                                n_estimators=500, n_jobs=-1, num_parallel_tree=1, random_state=0,  
                                reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,  
                                tree_method='gpu_hist', validate_parameters=1, verbosity=None)  
  
In [163...]  
x_model = XGBRegressor(n_estimators= 500, max_depth= 5,  
                       learning_rate= 0.01, tree_method='gpu_hist', gpu_id = 0)  
x_model.fit(final_train, tsne_train_output)  
  
y_pred = x_model.predict(final_train)  
y_pred_train = [round(value) for value in y_pred]  
  
y_pred = x_model.predict(final_test)  
y_pred_test = [round(value) for value in y_pred]  
  
train_error = (mean_absolute_error(tsne_train_output,  
                                   y_pred_train)) / (sum(tsne_train_output) / len(tsne_train_output))  
  
test_error = (mean_absolute_error(tsne_test_output,  
                                   y_pred_test)) / (sum(tsne_test_output) / len(tsne_test_output))  
  
print ("XGB Regressor - Train: ",train_error," Test: ",test_error)
```

XGB Regressor - Train: 0.12893694225353006 Test: 0.12654861251811048

Model comparision

In [164...]

```
# Documentation: https://pypi.org/project/prettytable/
from prettytable import PrettyTable

x = PrettyTable()

x.field_names= ["Model", "Train(MAPE)", "Test(MAPE)"]

x.add_row(["Linear Regression", 0.13324598902941232,
0.12889723956282564])
x.add_row(["Random Forest Regressor", 0.0934906831978537,
0.12547909124619153])
x.add_row(["XGB Regressor", 0.12893694225353006,
0.12654861251811048])
print(x)
```

Model	Train(MAPE)	Test(MAPE)
Linear Regression	0.13324598902941232	0.12889723956282564
Random Forest Regressor	0.0934906831978537	0.12547909124619153
XGB Regressor	0.12893694225353006	0.12654861251811048

Conclusion

The model performance is improved a bit by using fourier transform features(frequency and amplitude) but the mape is not below 12.

Task 3

Holt Winter Forecasting

<https://towardsdatascience.com/holt-winters-exponential-smoothing-d703072c0572>

Exponential smoothing(Triple)

<https://machinelearningmastery.com/exponential-smoothing-for-time-series-forecasting-in-python/>

<https://www.itl.nist.gov/div898/handbook/pmc/section4/pmc435.htm>

In [171...]

```
# This code is copied from : https://github.com/pranaysawant/New-York-Taxi-Demand-Prediction/blob/master/NYC_Final.ipynb

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

def initial_seasonal_components(series, slen):
    seasonals = {}
    season_averages = []
    n_seasons = int(len(series)/slen)
    # compute season averages
    for j in range(n_seasons):

        season_averages.append(sum(series[slen*j:slen*j+slen])/float(slen))
        # compute initial values
        for i in range(slen):
            sum_of_vals_over_avg = 0.0
            for j in range(n_seasons):
                sum_of_vals_over_avg += series[slen*j+i]-
season_averages[j]
            seasonals[i] = sum_of_vals_over_avg/n_seasons
    return seasonals

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

```
trend = beta * (smooth-last_smooth) + (1-beta)*trend  
seasonals[i%slen] = gamma*(val-smooth) + (1-  
gamma)*seasonals[i%slen]  
result.append(smooth+trend+seasonals[i%slen])  
return result
```

Hyper-parameter tuning of alpha, beta, gamma and season length

This code is copied from : [https://github.com/pranaysawant/New-York-Taxi-Demand-Prediction
/blob/master/NYC_Final.ipynb](https://github.com/pranaysawant/New-York-Taxi-Demand-Prediction/blob/master/NYC_Final.ipynb) and manipulated acc. to the requirement.

In [178...]

```
# finding the best alpha, the other parameters are taking as
constant

alpha = [0.1,0.2,0.3,0.4]
beta = 0.25
gamma = 0.25
season_len = 24

for a in alpha:
    y_pred = []
    y_pred_ls = []
    for region in range(0,40):
        y_pred =
triple_exponential_smoothing(regions_cum[region][0:13104],
season_len, a, beta, gamma, 0)
        y_pred_ls.append(y_pred[5:])

tsne_train_flat_triple_avg = [i[:9169] for i in y_pred_ls]
tsne_test_flat_triple_avg = [i[9169:] for i in y_pred_ls]

tsne_train_triple_avg = sum(tsne_train_flat_triple_avg,[])
tsne_test_triple_avg = sum(tsne_test_flat_triple_avg,[])

final_train['triple_exp'] = tsne_train_triple_avg
final_test['triple_exp'] = tsne_test_triple_avg

# checking model performance based on specific alpha
x_model = XGBRegressor(tree_method='gpu_hist', gpu_id = 0)
x_model.fit(final_train, tsne_train_output)

y_pred = x_model.predict(final_train)
y_pred_train = [round(value) for value in y_pred]

y_pred = x_model.predict(final_test)
y_pred_test = [round(value) for value in y_pred]

train_error = (mean_absolute_error(tsne_train_output,
y_pred_train)) / (sum(tsne_train_output)/len(tsne_train_output))
```

```
[ ]
```

```
alpha: 0.1 Train: 0.1127630903009719 Test: 0.11491977873687349
alpha: 0.2 Train: 0.10719369325390786 Test: 0.1227675857929806
alpha: 0.3 Train: 0.12330100786768546 Test: 0.13081121530339576
alpha: 0.4 Train: 0.12157281138842445 Test: 0.12983898405513994
```

In [179...]

```
# finding the best beta, the other parameters are taking as
constant

alpha = 0.1
beta = [0.1, 0.15, 0.20, 0.25]
gamma = 0.25
season_len = 24

for b in beta:
    y_pred = []
    y_pred_ls = []
    for region in range(0, 40):
        y_pred =
    triple_exponential_smoothing(regions_cum[region][0:13104],
    season_len, alpha, b, gamma, 0)
    y_pred_ls.append(y_pred[5:])

tsne_train_flat_triple_avg = [i[:9169] for i in y_pred_ls]
tsne_test_flat_triple_avg = [i[9169:] for i in y_pred_ls]

tsne_train_triple_avg = sum(tsne_train_flat_triple_avg, [])
tsne_test_triple_avg = sum(tsne_test_flat_triple_avg, [])

final_train['triple_exp'] = tsne_train_triple_avg
final_test['triple_exp'] = tsne_test_triple_avg

# checking model performance based on specific beta
x_model = XGBRegressor(tree_method='gpu_hist', gpu_id = 0)
x_model.fit(final_train, tsne_train_output)

y_pred = x_model.predict(final_train)
y_pred_train = [round(value) for value in y_pred]

y_pred = x_model.predict(final_test)
y_pred_test = [round(value) for value in y_pred]

train_error = (mean_absolute_error(tsne_train_output,
y_pred_train)) / (sum(tsne_train_output) / len(tsne_train_output))
```

```
beta: 0.1 Train: 0.11421835061655494 Test: 0.11613351478846587  
beta: 0.15 Train: 0.11217454536034471 Test: 0.11380811119733628  
beta: 0.2 Train: 0.1121997851908555 Test: 0.11384576568586599  
beta: 0.25 Train: 0.1127630903009719 Test: 0.11491977873687349
```

In [180...]

```
# finding the best gamma, the other parameters are taking as
constant
alpha = 0.1
beta = 0.15
gamma = [0.1,0.3,0.4,0.5,0.65,0.75,0.85,0.95]
season_len = 24

for g in gamma:
    y_pred = []
    y_pred_ls = []
    for region in range(0,40):
        y_pred =
triple_exponential_smoothing(regions_cum[region][0:13104],
season_len, alpha, beta, g, 0)
        y_pred_ls.append(y_pred[5:])

tsne_train_flat_triple_avg = [i[:9169] for i in y_pred_ls]
tsne_test_flat_triple_avg = [i[9169:] for i in y_pred_ls]

tsne_train_triple_avg = sum(tsne_train_flat_triple_avg,[])
tsne_test_triple_avg = sum(tsne_test_flat_triple_avg,[])

final_train['triple_exp'] = tsne_train_triple_avg
final_test['triple_exp'] = tsne_test_triple_avg

# checking model performance based on specific gamma
x_model = XGBRegressor(tree_method='gpu_hist', gpu_id = 0)
x_model.fit(final_train, tsne_train_output)

y_pred = x_model.predict(final_train)
y_pred_train = [round(value) for value in y_pred]

y_pred = x_model.predict(final_test)
y_pred_test = [round(value) for value in y_pred]

train_error = (mean_absolute_error(tsne_train_output,
y_pred_train)) / (sum(tsne_train_output)/len(tsne_train_output))
```

```
gamma: 0.1 Train: 0.1191050745677302 Test: 0.12116507192580726
gamma: 0.3 Train: 0.1092357685393098 Test: 0.1105503289511409
gamma: 0.4 Train: 0.10199859091659753 Test: 0.10316441058301382
gamma: 0.5 Train: 0.09314100730344882 Test: 0.09394498621883947
gamma: 0.65 Train: 0.07695162997184185 Test: 0.07699253408973619
gamma: 0.75 Train: 0.06290843668171549 Test: 0.06313596902010406
gamma: 0.85 Train: 0.04662391095565229 Test: 0.046997866882781764
gamma: 0.95 Train: 0.031330170562257814 Test: 0.031386067456964935
```

In [181...]

```
# finding the best season_len, the other parameters are taking as
constant

alpha = 0.1
beta = 0.15
gamma = 0.95
season_len = [24, 72, 144]

for s in season_len:
    y_pred = []
    y_pred_ls = []
    for region in range(0, 40):
        y_pred =
    triple_exponential_smoothing(regions_cum[region][0:13104], s,
alpha, beta, gamma, 0)
    y_pred_ls.append(y_pred[5:])

tsne_train_flat_triple_avg = [i[:9169] for i in y_pred_ls]
tsne_test_flat_triple_avg = [i[9169:] for i in y_pred_ls]

tsne_train_triple_avg = sum(tsne_train_flat_triple_avg, [])
tsne_test_triple_avg = sum(tsne_test_flat_triple_avg, [])

final_train['triple_exp'] = tsne_train_triple_avg
final_test['triple_exp'] = tsne_test_triple_avg

# checking model performance based on specific season_len
x_model = XGBRegressor(tree_method='gpu_hist', gpu_id = 0)
x_model.fit(final_train, tsne_train_output)

y_pred = x_model.predict(final_train)
y_pred_train = [round(value) for value in y_pred]

y_pred = x_model.predict(final_test)
y_pred_test = [round(value) for value in y_pred]

train_error = (mean_absolute_error(tsne_train_output,
y_pred_train)) / (sum(tsne_train_output) / len(tsne_train_output))
```

```
season_len: 24 Train: 0.031330170562257814 Test: 0.031386067456964935
season_len: 72 Train: 0.1229256042022854 Test: 0.14655241619487056
season_len: 144 Train: 0.12197975375240648 Test: 0.1340424291541311
```

In [182...]

```
# preparing data with the best finalized parameters

alpha = 0.1
beta = 0.15
gamma = 0.95
season_len = 24

y_pred = []
y_pred_ls = []

for region in range(0, 40):
    y_pred = triple_exponential_smoothing(regions_cum[region]
[0:13104], season_len, alpha, beta, gamma, 0)
    y_pred_ls.append(y_pred[5:])

tsne_train_flat_triple_avg = [i[:9169] for i in y_pred_ls]
tsne_test_flat_triple_avg = [i[9169:] for i in y_pred_ls]

tsne_train_triple_avg = sum(tsne_train_flat_triple_avg, [])
tsne_test_triple_avg = sum(tsne_test_flat_triple_avg, [])

final_train['triple_exp'] = tsne_train_triple_avg
final_test['triple_exp'] = tsne_test_triple_avg
```

In [185...]

```
print("final train shape", final_train.shape)
print("final test shape", final_test.shape)

final_train.head()
```

```
final train shape (366760, 19)
final test shape (157200, 19)
```

Out[185...]

	ft_5	ft_4	ft_3	ft_2	ft_1	lat	lon	weekday	exp_avg	a_1	f_2
0	0	63	217	189	137	40.776228	-73.982119	4	150	369774.0	-0.012993
1	63	217	189	137	135	40.776228	-73.982119	4	139	369774.0	-0.012993
2	217	189	137	135	129	40.776228	-73.982119	4	132	369774.0	-0.012993
3	189	137	135	129	150	40.776228	-73.982119	4	144	369774.0	-0.012993
4	137	135	129	150	164	40.776228	-73.982119	4	158	369774.0	-0.012993

Linear Regression: Grid Search

In [186...]

```
# https://scikit-learn.org/stable/modules/generated
/sklearn.linear_model.LinearRegression.html
lr = LinearRegression()

params = {
    'fit_intercept' : [True, False],
    'normalize': [True, False],
    'copy_X': [True, False],
}

# https://scikit-learn.org/stable/modules/generated
/sklearn.model_selection.GridSearchCV.html
model = GridSearchCV(lr, param_grid= params, cv = 5, n_jobs= -1,
                      return_train_score= True, verbose= 5, scoring=
'neg_mean_squared_error')
model.fit(final_train, tsne_train_output)

print("Estimators: ", model.best_estimator_)
print("Params: ", model.best_params_)
```

Fitting 5 folds for each of 8 candidates, totalling 40 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.

[Parallel(n_jobs=-1)]: Done 2 tasks | elapsed: 5.2s

[Parallel(n_jobs=-1)]: Done 34 out of 40 | elapsed: 48.4s remaining: 8.5s

[Parallel(n_jobs=-1)]: Done 40 out of 40 | elapsed: 57.8s finished

Estimators: LinearRegression(fit_intercept=False, normalize=True)

Params: {'copy_X': True, 'fit_intercept': False, 'normalize': True}

In [187...]

```
model = LinearRegression(copy_X= True, fit_intercept = False,
normalize= True)

model.fit(final_train, tsne_train_output)

y_pred = model.predict(final_train)
y_pred_train = [round(value) for value in y_pred]

y_pred = model.predict(final_test)
y_pred_test = [round(value) for value in y_pred]

train_error = (mean_absolute_error(tsne_train_output,
y_pred_train)) / (sum(tsne_train_output) / len(tsne_train_output))

test_error = (mean_absolute_error(tsne_test_output,
y_pred_test)) / (sum(tsne_test_output) / len(tsne_test_output))

print ("Linear Regression - Train: ",train_error," Test:
",test_error)
```

Linear Regression - Train: 0.03610488999315242 Test: 0.03305213521975312

Random Forest: Random Search

In [195...]

```
model =  
  
RandomForestRegressor(max_features='sqrt', min_samples_leaf=4,  
  
min_samples_split=9, n_estimators=1000, n_jobs=-1)  
  
model.fit(final_train, tsne_train_output)  
  
y_pred = model.predict(final_train)  
y_pred_train = [round(value) for value in y_pred]  
  
y_pred = model.predict(final_test)  
y_pred_test = [round(value) for value in y_pred]  
  
train_error = (mean_absolute_error(tsne_train_output,  
y_pred_train)) / (sum(tsne_train_output) / len(tsne_train_output))  
  
test_error = (mean_absolute_error(tsne_test_output,  
y_pred_test)) / (sum(tsne_test_output) / len(tsne_test_output))  
  
print ("Random Forest Regressor - Train: ", train_error, " Test:  
", test_error)
```

```
Random Forest Regressor - Train: 0.024735566199262562 Test: 0.0357938407196  
0212
```

In [196...]

```
#feature importances based on analysis using random forest  
print (final_train.columns)  
print (model.feature_importances_)  
  
Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'lat', 'lon', 'weekday',  
       'exp_avg', 'a_1', 'f_2', 'a_2', 'f_3', 'a_3', 'f_4', 'a_4', 'f_5',  
       'a_5', 'triple_exp'],  
      dtype='object')  
[4.9931110e-02 7.10905010e-02 9.88135474e-02 1.25356442e-01  
 1.58997864e-01 7.06199062e-04 5.52054462e-04 2.42641313e-04  
 2.04759788e-01 1.55318155e-02 2.25249047e-04 2.78376737e-03  
 2.37042400e-04 3.16562063e-03 1.70491392e-04 5.37763491e-03  
 1.79124743e-04 5.04321807e-03 2.56835887e-01]
```

XGB-Regressor

In [188...]

```
depth = [1, 2, 3, 5, 7]
n_estimators = [100, 200, 300, 400, 500, 1000]
learning_rate = [0.01, 0.05, 0.1, 0.2, 0.3]
tuned_parameters = {
    'max_depth':depth,
    'n_estimators': n_estimators,
    'learning_rate': learning_rate,
}

# enable gpu, https://xgboost.readthedocs.io/en/latest
#gpu/index.html
model = XGBRegressor(tree_method='gpu_hist', gpu_id = 0)

# Applying RandomizedSearchCV with k folds = 5 and taking
'mean_absolute_error' as score metric
clf = RandomizedSearchCV(model, param_distributions=
tuned_parameters, scoring = 'neg_mean_squared_error',
                           cv=5, return_train_score= True, verbose =
5, n_jobs= -1)
clf.fit(final_train, tsne_train_output)

print('Best hyper parameter: ', clf.best_params_)
print('Model estimator: ', clf.best_estimator_)
```

```
Fitting 5 folds for each of 10 candidates, totalling 50 fits
[Parallel(n_jobs=-1)]: Using backend LokyBackend with 8 concurrent workers.
[Parallel(n_jobs=-1)]: Done   2 tasks      | elapsed:   8.3s
[Parallel(n_jobs=-1)]: Done  46 out of  50 | elapsed:  7.7min remaining:   39.9s
[Parallel(n_jobs=-1)]: Done  50 out of  50 | elapsed:  8.7min finished
Best hyper parameter:  {'n_estimators': 400, 'max_depth': 5, 'learning_rate':
0.05}
Model estimator:  XGBRegressor(base_score=0.5, booster='gbtree', colsample_by
evel=1,
                           colsample_bynode=1, colsample_bytree=1, gamma=0, gpu_id=0,
                           importance_type='gain', interaction_constraints='',
                           learning_rate=0.05, max_delta_step=0, max_depth=5,
                           min_child_weight=1, missing=nan, monotone_constraints='()',
                           n_estimators=400, n_jobs=8, num_parallel_tree=1, random_state=0,
                           reg_alpha=0, reg_lambda=1, scale_pos_weight=1, subsample=1,
                           tree_method='gpu_hist', validate_parameters=1, verbosity=None)
```

In [189]:

```

x_model = XGBRegressor(n_estimators= 400, max_depth= 5,
learning_rate= 0.05, tree_method='gpu_hist', gpu_id = 0)
x_model.fit(final_train, tsne_train_output)

y_pred = x_model.predict(final_train)
y_pred_train = [round(value) for value in y_pred]

y_pred = x_model.predict(final_test)
y_pred_test = [round(value) for value in y_pred]

train_error = (mean_absolute_error(tsne_train_output,
y_pred_train))/(sum(tsne_train_output)/len(tsne_train_output))

test_error = (mean_absolute_error(tsne_test_output,
y_pred_test))/(sum(tsne_test_output)/len(tsne_test_output))

print ("XGB Regressor - Train: ",train_error," Test: ",test_error)

```

XGB Regressor - Train: 0.032564305121960055 Test: 0.0316424811440848

Model comparision

In [197]:

```

x = PrettyTable()

x.field_names= ["Model", "Train(MAPE)", "Test(MAPE)"]
x.add_row(["Linear Regression", 0.03610488999315242,
0.03305213521975312])
x.add_row(["Random Forest Regressor", 0.024735566199262562,
0.03579384071960212])
x.add_row(["XGB Regressor", 0.032564305121960055 ,
0.0316424811440848])
print(x)

```

Model	Train(MAPE)	Test(MAPE)
Linear Regression	0.03610488999315242	0.03305213521975312
Random Forest Regressor	0.024735566199262562	0.03579384071960212
XGB Regressor	0.032564305121960055	0.0316424811440848

In []: