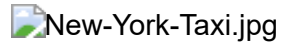


Taxi demand prediction in New York City



```
In [1]: #Importing Libraries
# pip3 install graphviz
#pip3 install dask
#pip3 install toolz
#pip3 install cloudpickle
# https://www.youtube.com/watch?v=ieW3G7ZzRZ0
# https://github.com/dask/dask-tutorial
# please do go through this python notebook: https://github.com/dask/dask-tutorial/blob/master/07_dataframe.ipynb
import dask.dataframe as dd#similar to pandas

import pandas as pd#pandas to create small dataframes

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

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

import time #Convert to unix time

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

# matplotlib: used to plot graphs
import matplotlib
# matplotlib.use('nbagg') : matplotlib uses this protocall which makes plots more user intractive like zoom in a
nd zoom out
matplotlib.use('nbagg')
#!pip install matplotlib.pyplot
%matplotlib inline
import matplotlib.pyplot as plt
import seaborn as sns#Plots
from matplotlib import rcParams#Size of plots

# this lib is used while we calculate the stight line distance between two (lat,lon) pairs in miles
#!pip3 install gpxpy
import gpxpy.geo #Get the haversine distance
```

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

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

#!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")
```

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 (FHVs)

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

Green Taxi: Street Hail Livery (SHL)

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

Credits: Quora

Footnote:

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

Data Collection

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

file name	file name size	number of records	number of features
yellow_tripdata_2016-01	1. 59G	10906858	19
yellow_tripdata_2016-02	1. 66G	11382049	19
yellow_tripdata_2016-03	1. 78G	12210952	19
yellow_tripdata_2016-04	1. 74G	11934338	19
yellow_tripdata_2016-05	1. 73G	11836853	19
yellow_tripdata_2016-06	1. 62G	11135470	19
yellow_tripdata_2016-07	884Mb	10294080	17
yellow_tripdata_2016-08	854Mb	9942263	17
yellow_tripdata_2016-09	870Mb	10116018	17
yellow_tripdata_2016-10	933Mb	10854626	17
yellow_tripdata_2016-11	868Mb	10102128	17
yellow_tripdata_2016-12	897Mb	10449408	17
yellow_tripdata_2015-01	1.84Gb	12748986	19
yellow_tripdata_2015-02	1.81Gb	12450521	19
yellow_tripdata_2015-03	1.94Gb	13351609	19
yellow_tripdata_2015-04	1.90Gb	13071789	19
yellow_tripdata_2015-05	1.91Gb	13158262	19
yellow_tripdata_2015-06	1.79Gb	12324935	19

yellow_tripdata_2015-07	1.68Gb	11562783	19
yellow_tripdata_2015-08	1.62Gb	11130304	19
yellow_tripdata_2015-09	1.63Gb	11225063	19
yellow_tripdata_2015-10	1.79Gb	12315488	19
yellow_tripdata_2015-11	1.65Gb	11312676	19
yellow_tripdata_2015-12	1.67Gb	11460573	19

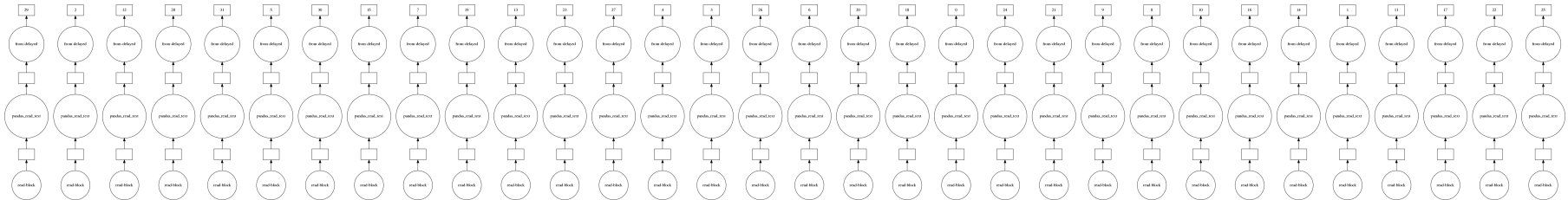
```
In [6]: #Looking at the features
# dask dataframe : # https://github.com/dask/dask-tutorial/blob/master/07_dataframe.ipynb
month = dd.read_csv('yellow_tripdata_2015-01.csv')
print(month.columns)
```

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

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

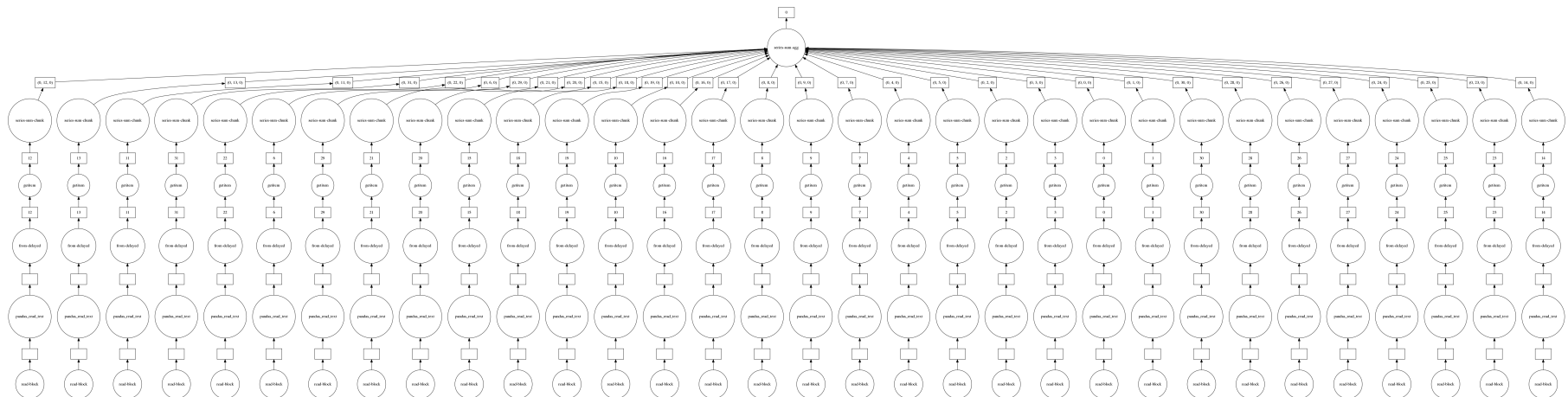
*# to see the visulaization you need to install graphviz
pip3 install graphviz if this doesnt work please check the install_graphviz.jpg in the drive
#!pip3 install graphviz
month.visualize()*

Out[3]:



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

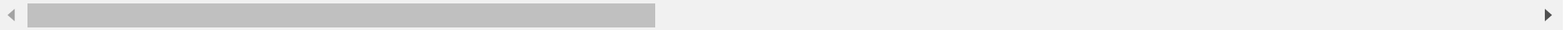
Out[4]:



In [5]: `month.head()`

Out[5]:

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



Features in the dataset:

Field Name	Description
VendorID	A code indicating the TPEP provider that provided the record. 1. Creative Mobile Technologies 2. VeriFone Inc.
tpep_pickup_datetime	The date and time when the meter was engaged.
tpep_dropoff_datetime	The date and time when the meter was disengaged.
Passenger_count	The number of passengers in the vehicle. This is a driver-entered value.
Trip_distance	The elapsed trip distance in miles reported by the taximeter.
Pickup_longitude	Longitude where the meter was engaged.
Pickup_latitude	Latitude where the meter was engaged.
RateCodeID	The final rate code in effect at the end of the trip. 1. Standard rate 2. JFK 3. Newark 4. Nassau or Westchester 5. Negotiated fare 6. Group ride
Store_and_fwd_flag	This flag indicates whether the trip record was held in vehicle memory before sending to the vendor, aka "store and forward," because the vehicle did not have a connection to the server. Y= store and forward trip N= not a store and forward trip
Dropoff_longitude	Longitude where the meter was disengaged.
Dropoff_latitude	Latitude where the meter was disengaged.

Payment_type	A numeric code signifying how the passenger paid for the trip. 1. Credit card 2. Cash 3. No charge 4. Dispute 5. Unknown 6. Voided trip
Fare_amount	The time-and-distance fare calculated by the meter.
Extra	Miscellaneous extras and surcharges. Currently, this only includes. the \$0.50 and \$1 rush hour and overnight charges.
MTA_tax	0.50 MTA tax that is automatically triggered based on the metered rate in use.
Improvement_surcharge	0.30 improvement surcharge assessed trips at the flag drop. the improvement surcharge began being levied in 2015.
Tip_amount	Tip amount – This field is automatically populated for credit card tips.Cash tips are not included.
Tolls_amount	Total amount of all tolls paid in trip.
Total_amount	The total amount charged to passengers. Does not include cash tips.

ML Problem Formulation

Time-series forecasting and Regression

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

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

Performance metrics

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

Data Cleaning

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

In [6]: *#table below shows few datapoints along with all our features*
`month.head(5)`

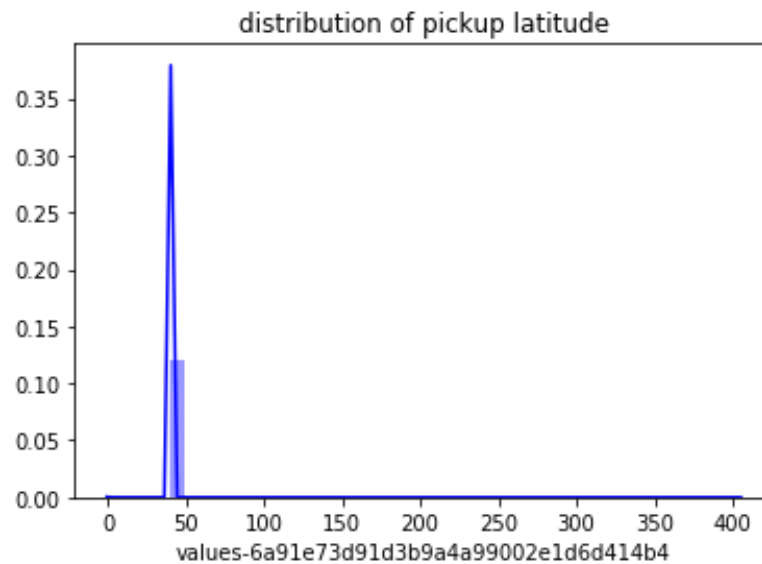
Out[6]:

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

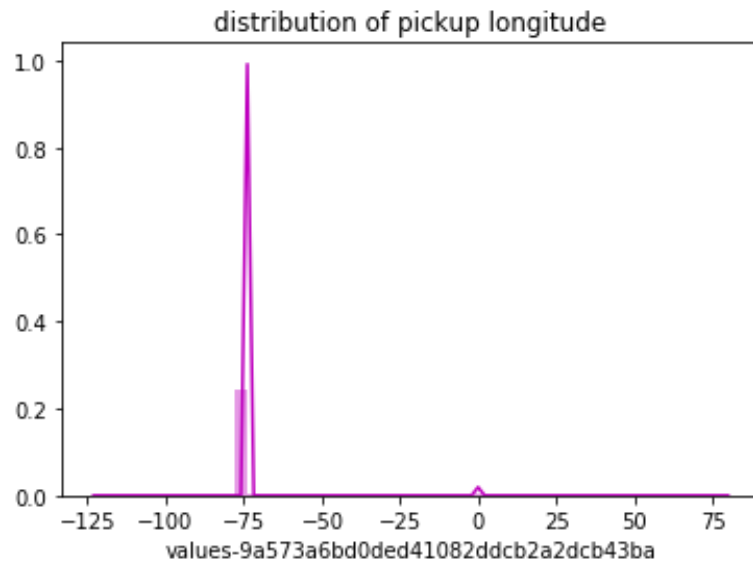
1. Pickup Latitude and Pickup Longitude

It is inferred from the source <https://www.flickr.com/places/info/2459115> (<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 [7]: sns.distplot(month['pickup_latitude'].values, label = 'pickup_latitude',color="b")  
plt.title("distribution of pickup latitude")  
plt.show()
```



```
In [8]: sns.distplot(month['pickup_longitude'].values, label = 'pickup_longitude',color="m")  
plt.title("distribution of pickup longitude")  
#plt.tight_layout()  
plt.show()
```



```
In [6]: outlier_locations = month[((month.pickup_longitude <= -74.15) | (month.pickup_latitude <= 40.5774) | \  
                                   (month.pickup_longitude >= -73.7004) | (month.pickup_latitude >= 40.9176))]
```

```
In [7]: print(outlier_locations.head())
```

	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	\
31	2	2015-01-15 19:05:43	2015-01-15 19:05:44	2	
61	1	2015-01-04 13:44:52	2015-01-04 13:56:49	1	
66	2	2015-01-04 13:44:52	2015-01-04 13:49:03	1	
157	1	2015-01-15 09:47:00	2015-01-15 10:00:07	1	
159	1	2015-01-15 09:47:02	2015-01-15 10:17:47	3	

	trip_distance	pickup_longitude	pickup_latitude	RateCodeID	\
31	0.01	0.0	0.0	5	
61	2.50	0.0	0.0	1	
66	0.85	0.0	0.0	1	
157	1.00	0.0	0.0	1	
159	8.30	0.0	0.0	1	

	store_and_fwd_flag	dropoff_longitude	dropoff_latitude	payment_type	\
31	N	0.0	0.0	1	
61	N	0.0	0.0	1	
66	N	0.0	0.0	2	
157	N	0.0	0.0	2	
159	N	0.0	0.0	1	

	fare_amount	extra	mta_tax	tip_amount	tolls_amount	\
31	60.0	0.0	0.0	0.00	0.00	
61	11.0	0.0	0.5	2.35	0.00	
66	5.5	0.0	0.5	0.00	0.00	
157	10.0	0.0	0.5	0.00	0.00	
159	27.5	0.0	0.5	10.00	5.33	

	improvement_surcharge	total_amount
31	0.3	60.30
61	0.0	14.15
66	0.3	6.30
157	0.3	10.80
159	0.3	43.63

```
In [8]: # Plotting pickup coordinates which are outside the bounding box of New-York
# we will collect all the points outside the bounding box of newyork city to outlier_locations
#outlier_locations = month[((month.pickup_longitude <= -74.15) | (month.pickup_latitude <= 40.5774)| \
#                           (month.pickup_longitude >= -73.7004) | (month.pickup_latitude >= 40.9176))]

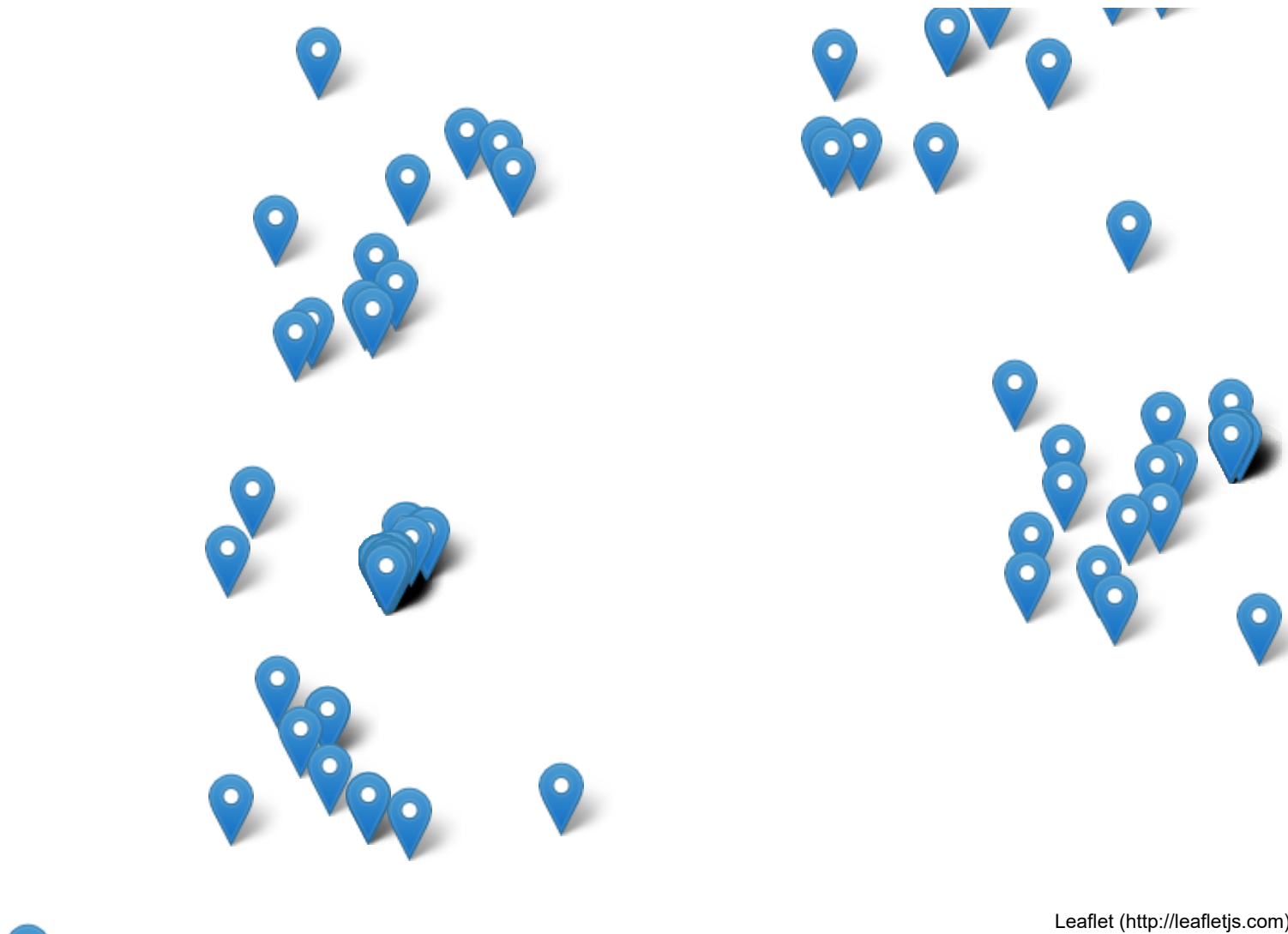
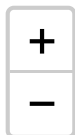
# creating a map with the a base location
# read more about the folium here: http://folium.readthedocs.io/en/latest/quickstart.html

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

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

# we will spot only first 10000 outliers on the map, plotting all the outliers will take more time
sample_locations = outlier_locations.head(1000000)
for i,j in sample_locations.iterrows():
    if int(j['pickup_latitude']) != 0:
        folium.Marker(list((j['pickup_latitude'],j['pickup_longitude']))).add_to(map_osm)
map_osm
```

Out[8]:



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

2. Dropoff Latitude & Dropoff Longitude

It is inferred from the source <https://www.flickr.com/places/info/2459115> (<https://www.flickr.com/places/info/2459115>) that New York is bounded by the location 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 [9]: # 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 indepth knowledge on these maps and plots

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

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

Out[9]:

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

3. Trip Durations:

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

In [7]: *#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 = month[['passenger_count', 'trip_distance', 'pickup_longitude', 'pickup_latitude', 'dropoff_longitude', 'dropoff_latitude', 'total_amount']].compute()
```

```
new_frame['trip_times'] = durations
```

```
new_frame['pickup_times'] = duration_pickup
```

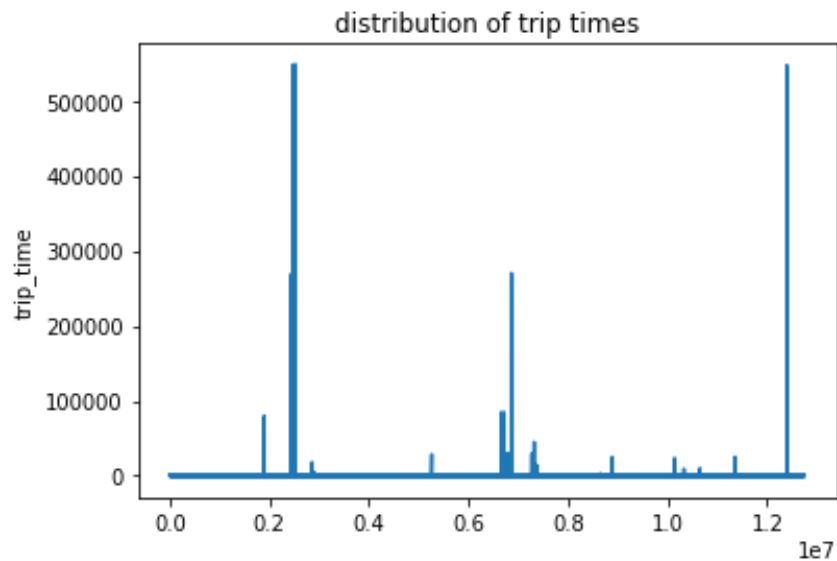
```
new_frame['Speed'] = 60*(new_frame['trip_distance']/new_frame['trip_times'])
```

```
return new_frame
```

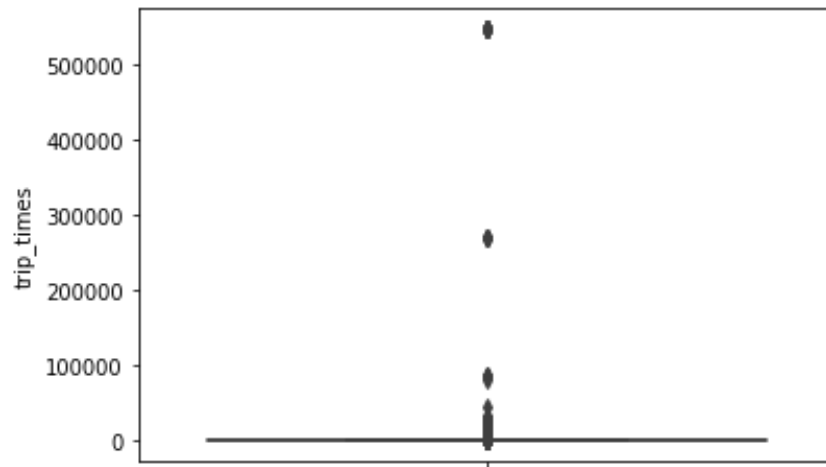
```
# print(frame_with_durations.head())
# passenger_count    trip_distance    pickup_longitude    pickup_latitude    dropoff_longitude    dropoff_
latitude    total_amount    trip_times    pickup_times    Speed
#    1            1.59            -73.993896            40.750111            -73.974785            40.75061
8            17.05            18.050000            1.421329e+09            5.285319
#    1            3.30            -74.001648            40.724243            -73.994415            40.75910
9            17.80            19.833333            1.420902e+09            9.983193
#    1            1.80            -73.963341            40.802788            -73.951820            40.82441
3            10.80            10.050000            1.420902e+09            10.746269
#    1            3.00            -73.971176            40.762428            -74.004181            40.74265
3            16.30            19.316667            1.420902e+09            9.318378
start = time.time()
frame_with_durations = return_with_trip_times(month)
print(time.time() - start)
432.08354020118713
```

```
In [14]: start = time.time()
#sns.set(style="white", palette="muted", color_codes=True)
#f, axes = plt.subplots(1, 1, figsize=(10, 7), sharex=True)
#sns.despine(left=True)
plt.plot(frame_with_durations['trip_times'].values)
plt.title("distribution of trip times")
plt.ylabel("trip_time")
end = time.time()
print("Time taken by above cell is {}".format((end-start)))
plt.show()
```

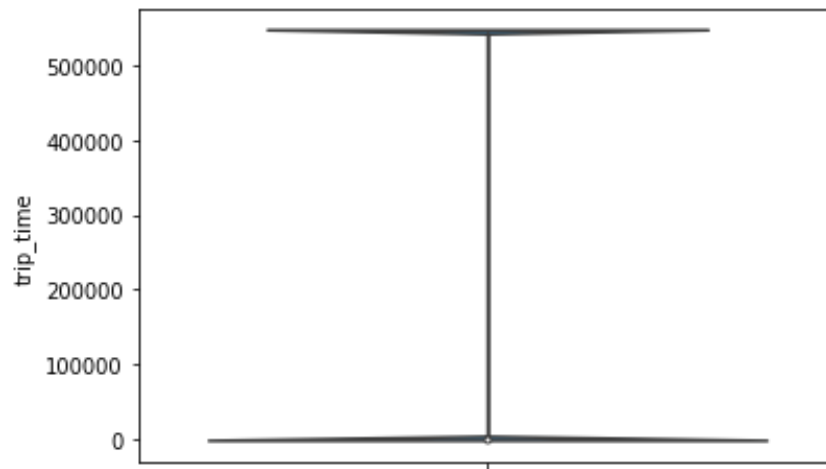
Time taken by above cell is 0.8431565761566162.



```
In [15]: # the skewed box plot shows us the presence of outliers
#plt.figure(figsize=(10,10))
#trip_time = frame_with_durations["trip_times"]
#import matplotlib.pyplot as plt
#plt.subplots(1,2,1)
sns.boxplot(y=frame_with_durations["trip_times"])#, data =frame_with_durations)
plt.ylabel("trip_times")
#trip_time.plot.box()
plt.show()
#plt.show()
```




```
In [16]: sns.violinplot(y=frame_with_durations["trip_times"])#, data =frame_with_durations)
plt.ylabel("trip_time")
plt.show()
```



```
In [17]: #calculating 0-100th percentile to find a the correct percentile value for removal of outliers
for i in range(0,100,10):
    var =frame_with_durations["trip_times"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print ("100 percentile value is ",var[-1])
```

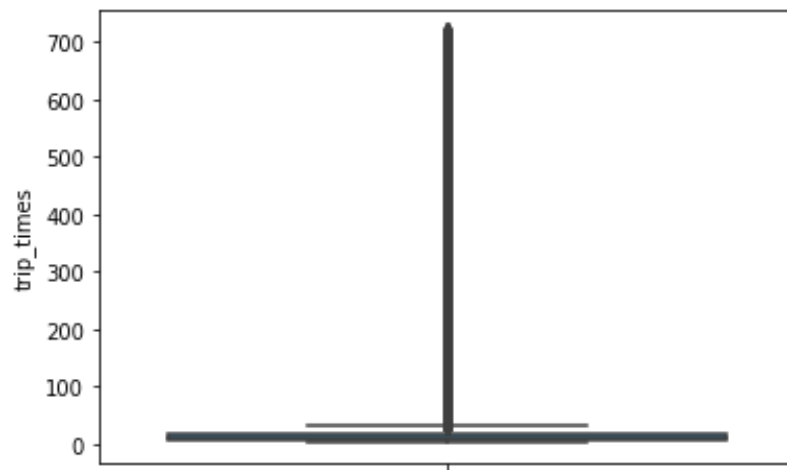
```
0 percentile value is -1211.0166666666667
10 percentile value is 3.8333333333333335
20 percentile value is 5.3833333333333334
30 percentile value is 6.8166666666666666
40 percentile value is 8.3
50 percentile value is 9.95
60 percentile value is 11.866666666666667
70 percentile value is 14.283333333333333
80 percentile value is 17.633333333333333
90 percentile value is 23.45
100 percentile value is 548555.6333333333
```

```
In [18]: #Looking further from the 99th percecntile
for i in range(90,100):
    var =frame_with_durations["trip_times"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print ("100 percentile value is ",var[-1])
```

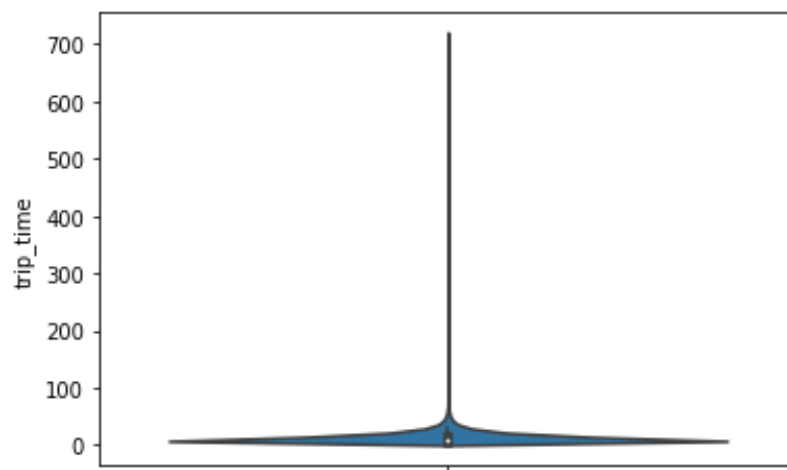
```
90 percentile value is 23.45
91 percentile value is 24.35
92 percentile value is 25.383333333333333
93 percentile value is 26.55
94 percentile value is 27.933333333333334
95 percentile value is 29.583333333333332
96 percentile value is 31.683333333333334
97 percentile value is 34.466666666666667
98 percentile value is 38.716666666666667
99 percentile value is 46.75
100 percentile value is  548555.6333333333
```

```
In [8]: #removing data based on our analysis and TLC regulations
frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_times>1) & (frame_with_durations.t
rip_times<720)]
#we know that new york city cabs are restricted trip to only 12hr, there for we are taking maximum time as 720mi
ns
```

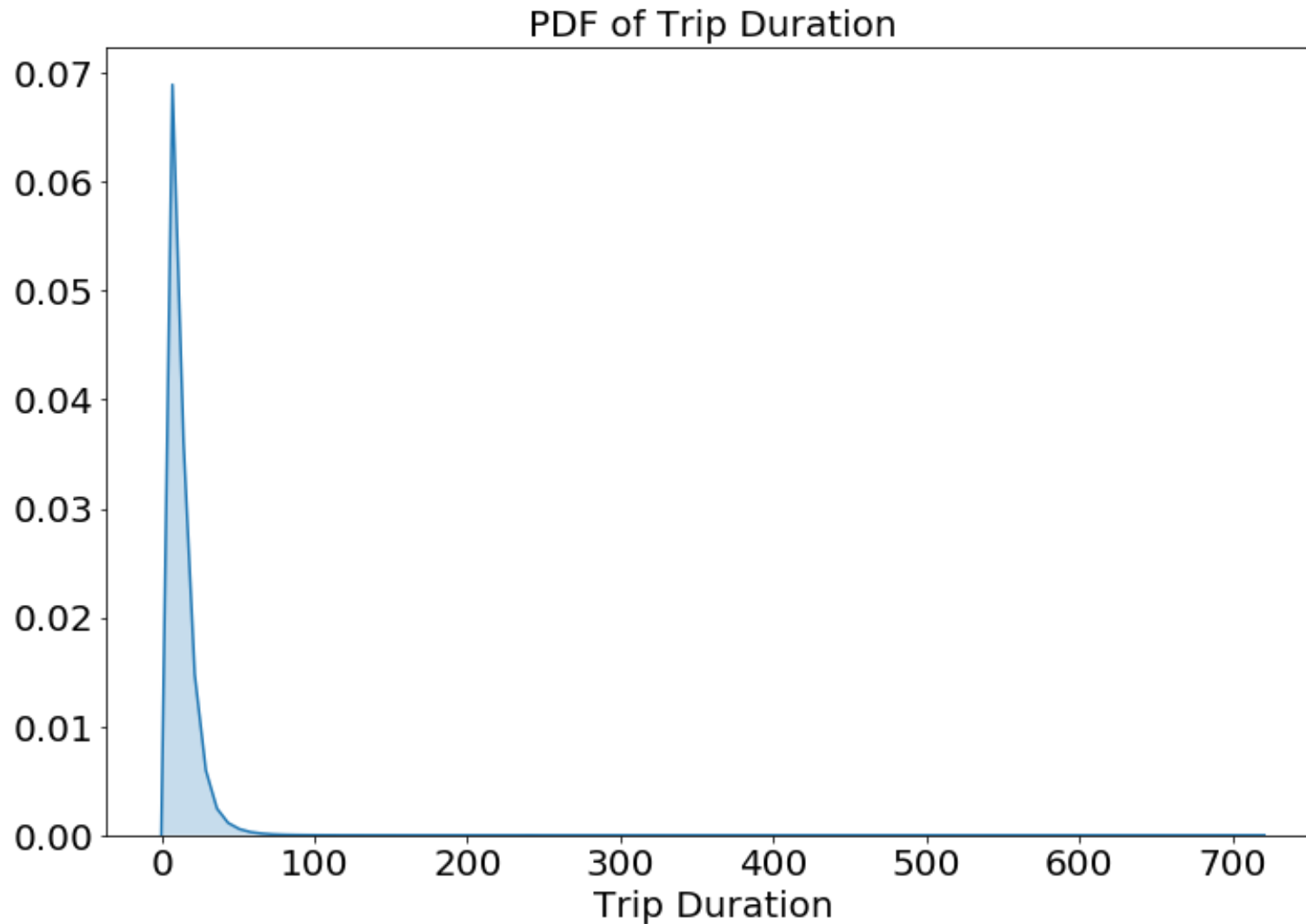
```
In [20]: #box-plot after removal of outliers  
sns.boxplot(y="trip_times", data =frame_with_durations_modified)  
plt.show()
```



```
In [21]: sns.violinplot(y=frame_with_durations_modified["trip_times"], data =frame_with_durations)  
plt.ylabel("trip_time")  
plt.show()
```

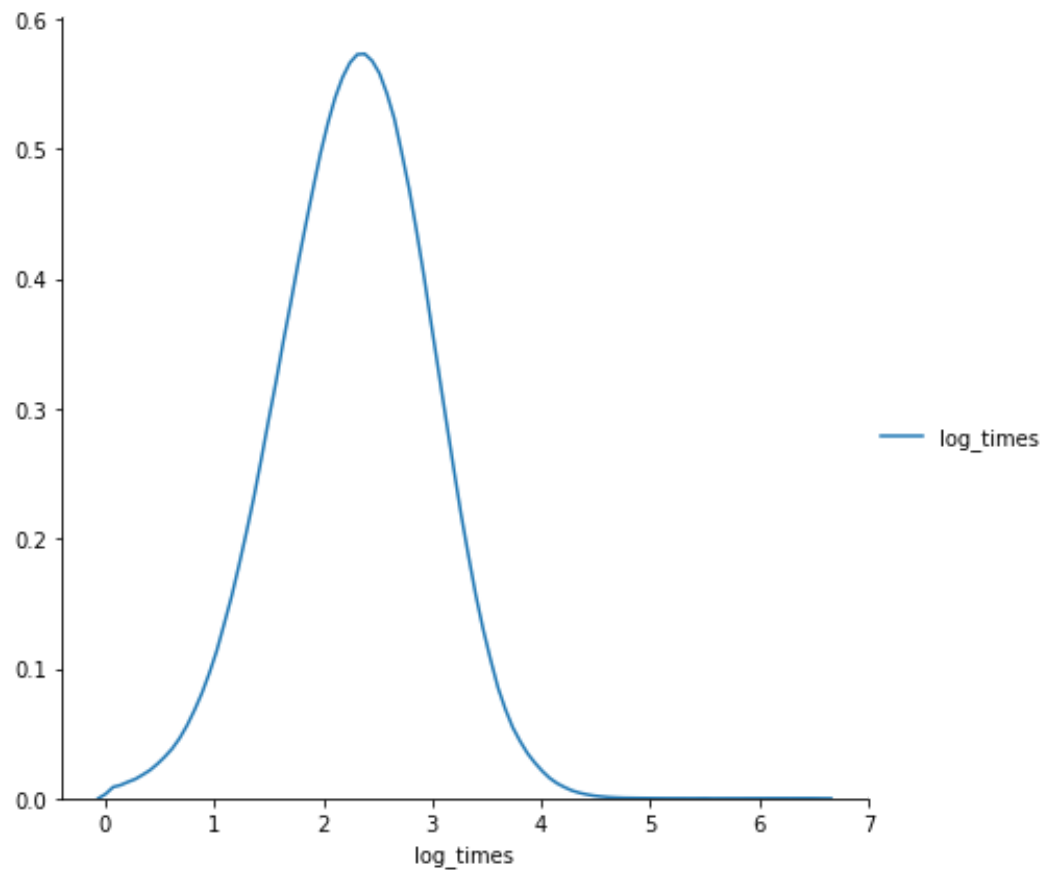


```
In [22]: plt.figure(figsize = (12,8))  
sns.kdeplot(frame_with_durations_modified["trip_times"].values, shade = True, cumulative = False)  
plt.tick_params(labelsize = 20)  
plt.xlabel("Trip Duration", fontsize = 20)  
plt.title("PDF of Trip Duration", fontsize = 20)  
plt.show()
```

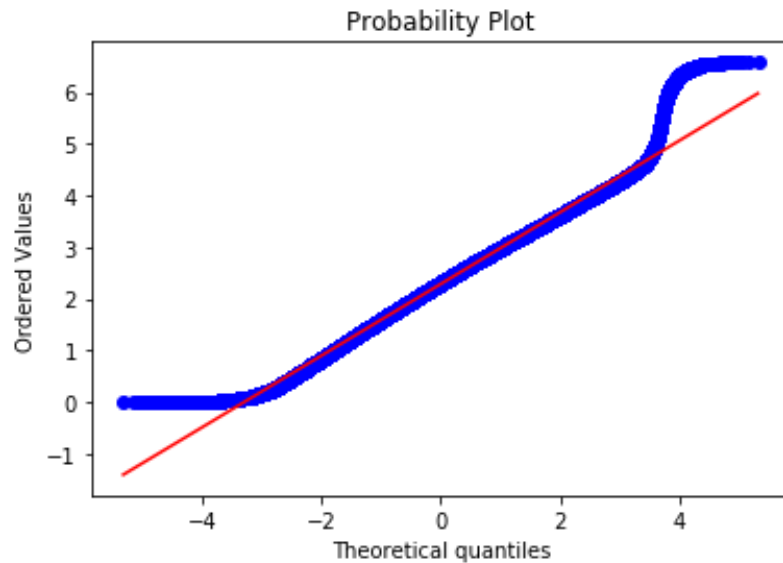


```
In [9]: #converting the values to log-values to chec for log-normal  
import math  
frame_with_durations_modified['log_times']=[math.log(i) for i in frame_with_durations_modified['trip_times'].values]
```

```
In [24]: #pdf of log-values  
sns.FacetGrid(frame_with_durations_modified,size=6) \  
    .map(sns.kdeplot,"log_times") \  
    .add_legend();  
plt.show();
```



```
In [25]: #Q-Q plot for checking if trip-times is log-normal
import scipy
scipy.stats.probplot(frame_with_durations_modified['log_times'].values, plot=plt)
plt.show()
```



4. Speed

```
In [26]: frame_with_durations_modified.head(2)
```

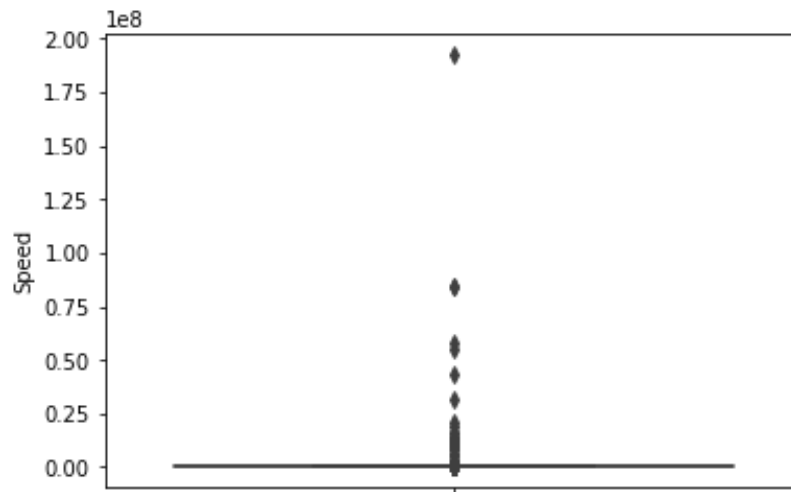
Out[26]:

	passenger_count	trip_distance	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	total_amount	trip_time
0	1	1.59	-73.993896	40.750111	-73.974785	40.750618	17.05	18.0500
1	1	3.30	-74.001648	40.724243	-73.994415	40.759109	17.80	19.8333

```
In [27]: frame_with_durations_modified['Speed'].describe()
```

```
Out[27]: count    1.263525e+07  
mean      6.952009e+01  
std       7.049486e+04  
min       0.000000e+00  
25%      8.387097e+00  
50%     1.106866e+01  
75%     1.475589e+01  
max      1.928571e+08  
Name: Speed, dtype: float64
```

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



```
In [29]: #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 [30]: #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 [31]: #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 [11]: #removing further outliers based on the 99.9th percentile value
frame_with_durations_modified=frame_with_durations[(frame_with_durations.Speed>0) & (frame_with_durations.Speed<
45.31)]
```

```
In [33]: #avg.speed of cabs in New-York
sum(frame_with_durations_modified["Speed"]) / float(len(frame_with_durations_modified["Speed"]))
```

```
Out[33]: 12.450173996027528
```

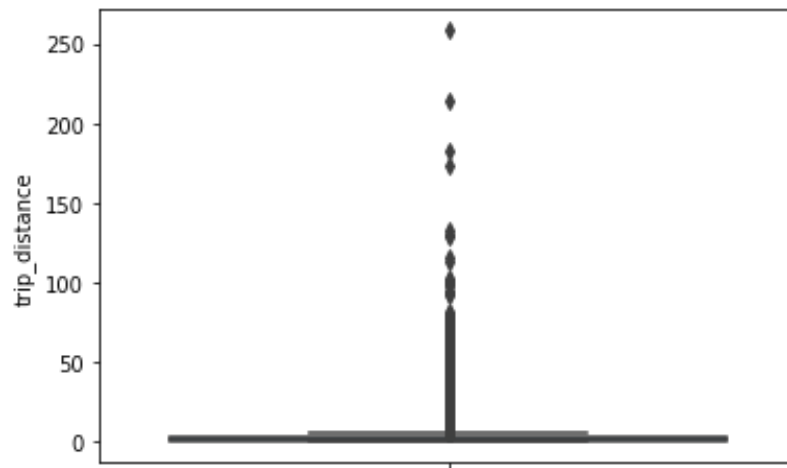
The avg speed in Newyork speed is 12.45miles/hr, so a cab driver can travel 2 miles per 10min on avg.

4. Trip Distance

```
In [34]: frame_with_durations_modified['trip_distance'].describe()
```

```
Out[34]: count      1.264716e+07  
mean        2.783871e+00  
std         3.336452e+00  
min         1.000000e-02  
25%         1.000000e+00  
50%         1.690000e+00  
75%         3.000000e+00  
max         2.589000e+02  
Name: trip_distance, dtype: float64
```

```
In [35]: # 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 [36]: #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 [37]: #calculating trip distance 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["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 [38]: #calculating trip distance values at each percentile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100
for i in np.arange(0.0, 1.0, 0.1):
    var =frame_with_durations_modified["trip_distance"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))
print("100 percentile value is ",var[-1])
```

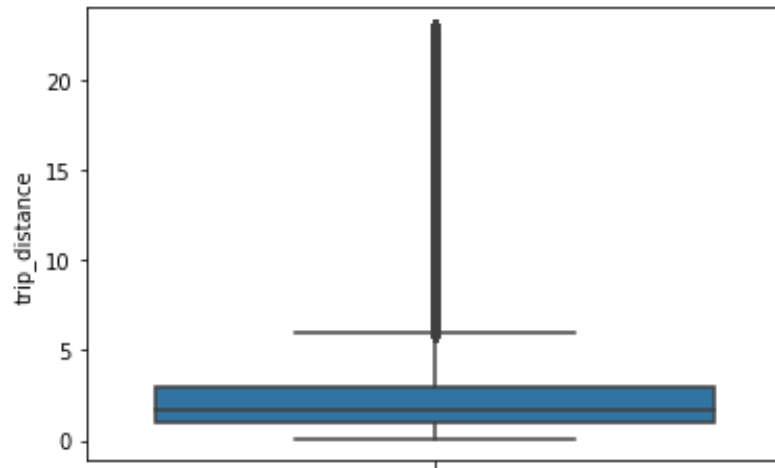
```
99.0 percentile value is 18.17
99.1 percentile value is 18.37
99.2 percentile value is 18.6
99.3 percentile value is 18.83
99.4 percentile value is 19.13
99.5 percentile value is 19.5
99.6 percentile value is 19.96
99.7 percentile value is 20.5
99.8 percentile value is 21.22
99.9 percentile value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

```
In [12]: #removing further outliers based on the 99.9th percentile value
frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_distance>0) & (frame_with_durations.trip_distance<23)]
```

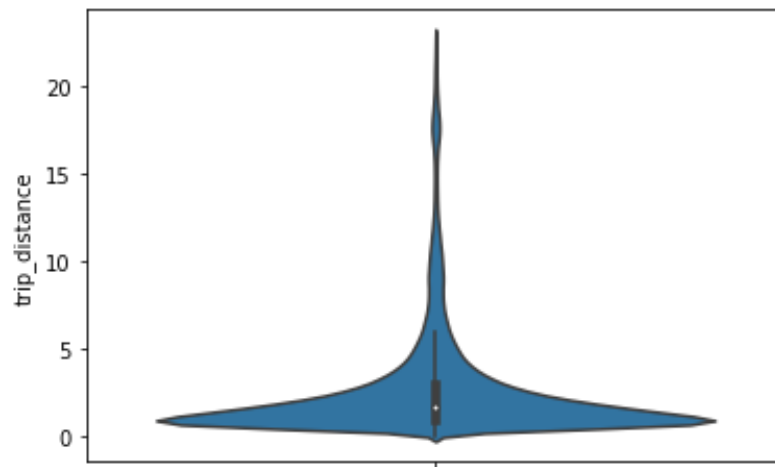
```
In [40]: frame_with_durations_modified['trip_distance'].describe()
```

```
Out[40]: count      1.265639e+07
mean        2.773040e+00
std         3.268942e+00
min         1.000000e-02
25%         1.000000e+00
50%         1.690000e+00
75%         3.000000e+00
max         2.299000e+01
Name: trip_distance, dtype: float64
```

```
In [41]: #box-plot after removal of outliers  
sns.boxplot(y="trip_distance", data = frame_with_durations_modified)  
plt.show()
```



```
In [42]: sns.violinplot(y="trip_distance", data = frame_with_durations_modified)  
plt.show()
```

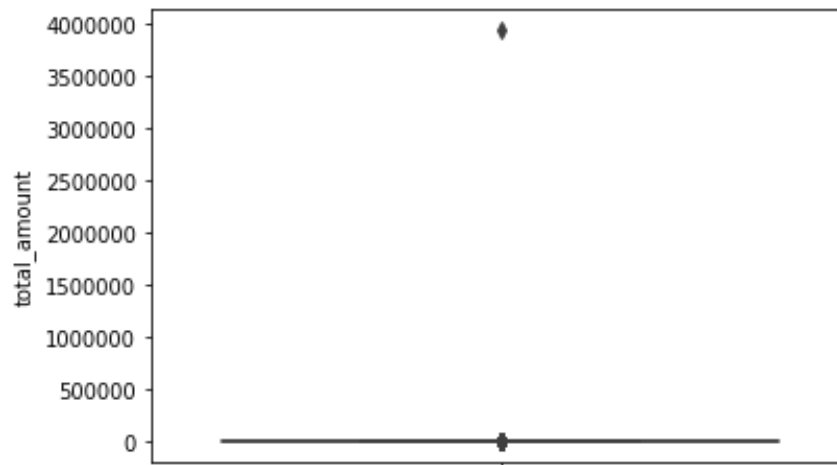


5. Total Fare

```
In [43]: frame_with_durations_modified['total_amount'].describe()
```

```
Out[43]: count      1.265639e+07  
mean        1.497958e+01  
std         1.110534e+03  
min        -2.425500e+02  
25%         8.300000e+00  
50%        1.116000e+01  
75%        1.630000e+01  
max         3.950612e+06  
Name: total_amount, dtype: float64
```

```
In [44]: # 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 [45]: *#calculating total fare amount values at each percentile 0,10,20,30,40,50,60,70,80,90,100*

```
for i in range(0,100,10):  
    var = frame_with_durations_modified["total_amount"].values  
    var = np.sort(var,axis = None)  
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))  
print("100 percentile value is ",var[-1])
```

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

In [46]: *#calculating total fare amount values at each percentile 90,91,92,93,94,95,96,97,98,99,100*

```
for i in range(90,100):  
    var = frame_with_durations_modified["total_amount"].values  
    var = np.sort(var,axis = None)  
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))  
print("100 percentile value is ",var[-1])
```

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

```
In [47]: #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

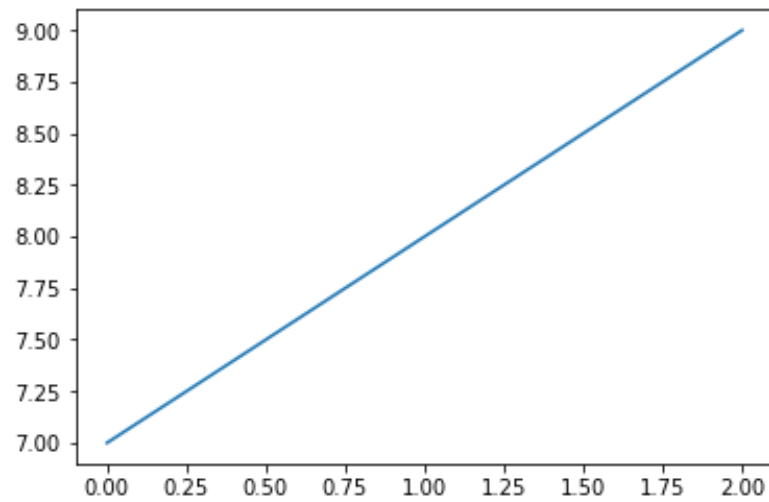
```
In [48]: lst = [i for i in range(10)]
print(lst)

[0, 1, 2, 3, 4, 5, 6, 7, 8, 9]
```

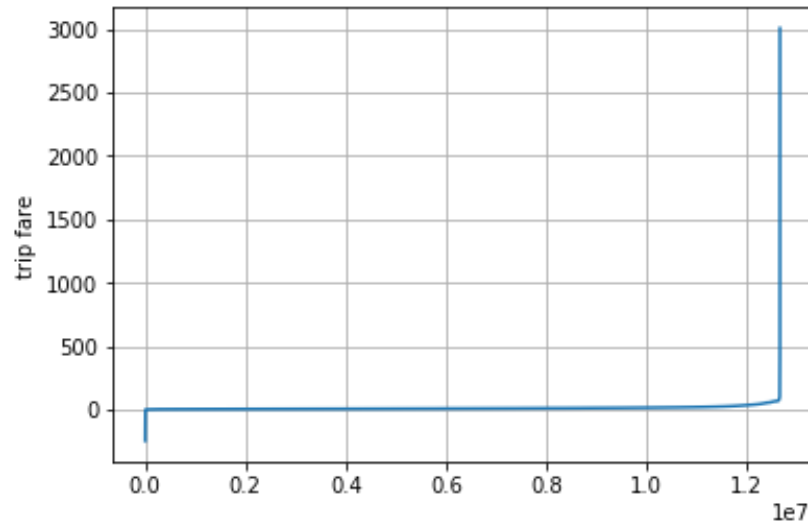


```
In [49]: plt.plot(lst[-3:])  
plt.show
```

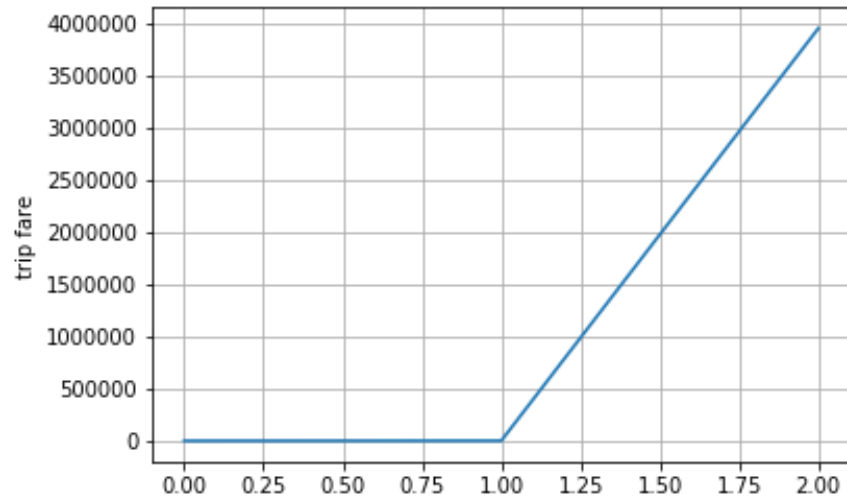
```
Out[49]: <function matplotlib.pyplot.show(*args, **kw)>
```



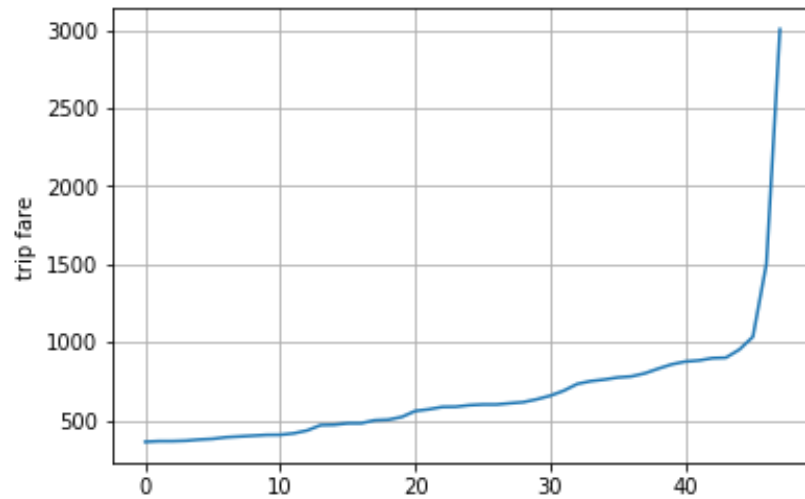
```
In [50]: #below plot shows us the fare values(sorted) to find a sharp increase to remove those values as outliers  
# plot the fare amount excluding last two values in sorted data  
plt.plot(var[:-2])  
plt.grid()  
plt.ylabel("trip fare")  
plt.show()
```



```
In [51]: # 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.grid()  
plt.ylabel("trip fare")  
plt.show()
```



```
In [52]: #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.grid()  
plt.ylabel("trip fare")  
plt.show()
```



Remove all outliers/erronous points.

In [13]: *#removing all outliers based on our univariate analysis above*

```
def remove_outliers(new_frame):

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

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

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

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

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

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

```

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

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

```

```

In [14]: print ("Removing outliers in the month of Jan-2015")
print ("----")
frame_with_durations_outliers_removed = remove_outliers(frame_with_durations)
print("fraction of data points that remain after removing outliers", float(len(frame_with_durations_outliers_removed)/len(frame_with_durations)))

```

Removing outliers in the month of Jan-2015

Number of pickup records = 12748986
 Number of outlier coordinates lying outside NY boundaries: 293919
 Number of outliers from trip times analysis: 23889
 Number of outliers from trip distance analysis: 92597
 Number of outliers from speed analysis: 24473
 Number of outliers from fare analysis: 5275
 Total outliers removed 377910

fraction of data points that remain after removing outliers 0.9703576425607495

Data-preperation

Clustering/Segmentation

```

In [15]: #trying different cluster sizes to choose the right K in K-means
coords = frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']].values
neighbours=[]

def find_min_distance(cluster_centers, cluster_len):
    nice_points = 0
    wrong_points = 0
    less2 = []
    more2 = []
    min_dist=1000
    for i in range(0, cluster_len):
        nice_points = 0
        wrong_points = 0
        for j in range(0, cluster_len):
            if j!=i:
                distance = gpxpy.geo.haversine_distance(cluster_centers[i][0], cluster_centers[i][1],cluster_centers[j][0], cluster_centers[j][1])# Latitude_1, Longitude_1, Latitude_2, Longitude_2
                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)
        #print("less2:/n",less2)
        #print("more2:/n",more2)
        neighbours.append(less2)
    print ("On choosing a cluster size of ",cluster_len,"\nAvg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2):", np.ceil(sum(less2)/len(less2)), "\nAvg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):", np.ceil(sum(more2)/len(more2)), "\nMin inter-cluster distance = ",min_dist,"\n---")

def find_clusters(increment):
    kmeans = MiniBatchKMeans(n_clusters=increment, batch_size=10000,random_state=42).fit(coords)
    frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_longitude']])
    #print("line1:",frame_with_durations_outliers_removed['pickup_cluster'])
    cluster_centers = kmeans.cluster_centers_
    cluster_len = len(cluster_centers)
    #print("cluster_centers:/n",cluster_centers)

```

```
#print("cluster_len:/n",cluster_len)
return cluster_centers, cluster_len

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



```
On choosing a cluster size of 10
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 2.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 8.0
Min inter-cluster distance = 1.0945442325142543
---
On choosing a cluster size of 20
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 4.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 16.0
Min inter-cluster distance = 0.7131298007387813
---
On choosing a cluster size of 30
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 22.0
Min inter-cluster distance = 0.5185088176172206
---
On choosing a cluster size of 40
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 32.0
Min inter-cluster distance = 0.5069768450363973
---
On choosing a cluster size of 50
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 12.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 38.0
Min inter-cluster distance = 0.365363025983595
---
On choosing a cluster size of 60
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 14.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 46.0
Min inter-cluster distance = 0.34704283494187155
---
On choosing a cluster size of 70
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 16.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 54.0
Min inter-cluster distance = 0.30502203163244707
---
On choosing a cluster size of 80
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 18.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 62.0
Min inter-cluster distance = 0.29220324531738534
---
```

```
On choosing a cluster size of 90
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 21.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 69.0
Min inter-cluster distance = 0.18257992857034985
---
```

Inference:

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

```
In [16]: # 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:

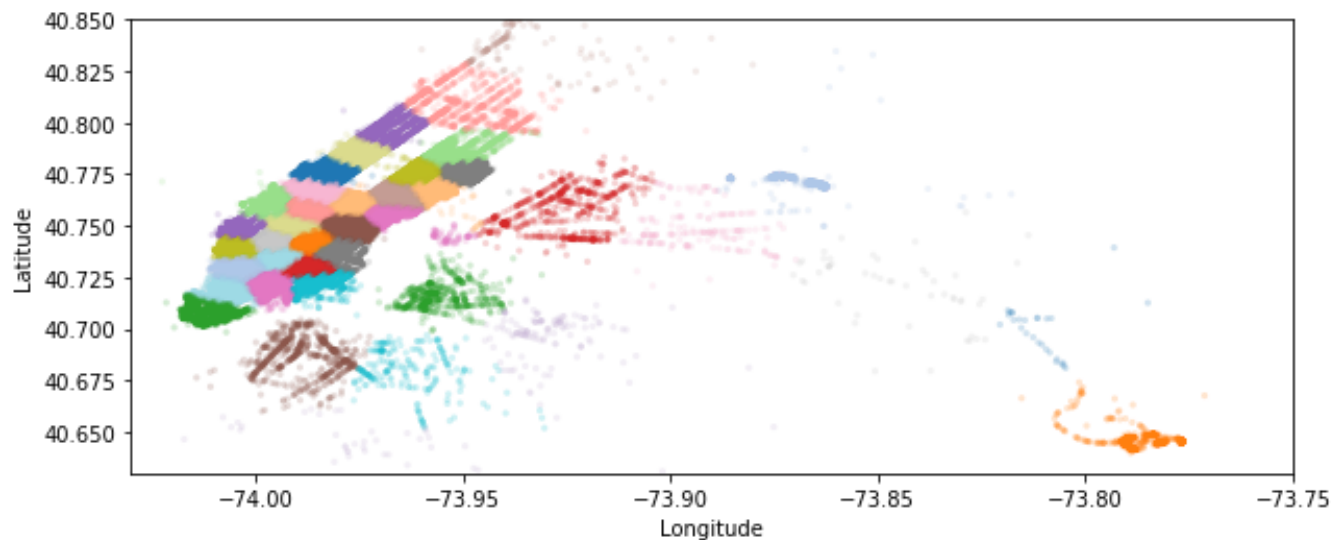
```
In [16]: # Plotting the cluster centers on OSM
cluster_centers = kmeans.cluster_centers_
cluster_len = len(cluster_centers)
map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')
for i in range(cluster_len):
    folium.Marker(list((cluster_centers[i][0],cluster_centers[i][1])), popup=(str(cluster_centers[i][0])+str(cluster_centers[i][1]))).add_to(map_osm)
map_osm
```

Out[16]:

Plotting the clusters:

```
In [17]: #Visualising the clusters on a map
def plot_clusters(frame):
    city_long_border = (-74.03, -73.75)
    city_lat_border = (40.63, 40.85)
    #plt.figure(figsize=(100,20))
    fig, ax = plt.subplots(ncols=1, nrows=1,figsize=(10,4))
    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 [17]: #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)
    #print(frame.head())
    return frame
```

```
In [18]: # 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']])
#print(frame_with_durations_outliers_removed.head())
jan_2015_frame = add_pickup_bins(frame_with_durations_outliers_removed,1,2015)
jan_2015_groupby = jan_2015_frame[['pickup_cluster','pickup_bins','trip_distance']].groupby(['pickup_cluster','pickup_bins']).count()
```

```
In [16]: # 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[16]:

	passenger_count	trip_distance	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude	total_amount	trip_time
0	1	1.59	-73.993896	40.750111	-73.974785	40.750618	17.05	18.0500
1	1	3.30	-74.001648	40.724243	-73.994415	40.759109	17.80	19.8333
2	1	1.80	-73.963341	40.802788	-73.951820	40.824413	10.80	10.0500
3	1	0.50	-74.009087	40.713818	-74.004326	40.719986	4.80	1.86666
4	1	3.00	-73.971176	40.762428	-74.004181	40.742653	16.30	19.3166

```
In [17]: # hear the trip_distance represents the number of pickups that are happend in that particular 10min intravel  
# this data frame has two indices  
# primary index: pickup_cluster (cluster number)  
# secondary index : pickup_bins (we devid whole months time into 10min intravels  $24*31*60/10 = 4464$ bins)  
jan_2015_groupby.head()
```

Out[17]:

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

```

In [19]: # 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_bin', 'trip_distance']].groupby(['pickup_cluster', 'pickup_bin']).count()

    return final_updated_frame, final_groupby_frame

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

jan_2016_frame, jan_2016_groupby = datapreparation(month_jan_2016, kmeans, 1, 2016)

```



```
feb_2016_frame,feb_2016_groupby = datapreparation(month_feb_2016,kmeans,2,2016)
mar_2016_frame,mar_2016_groupby = datapreparation(month_mar_2016,kmeans,3,2016)
Return with trip times..
Remove outliers..
Number of pickup records = 10906858
Number of outlier coordinates lying outside NY boundaries: 214677
Number of outliers from trip times analysis: 27190
Number of outliers from trip distance analysis: 79742
Number of outliers from speed analysis: 21047
Number of outliers from fare analysis: 4991
Total outliers removed 297784
---
Estimating clusters..
Final groupbying..
Return with trip times..
Remove outliers..
Number of pickup records = 11382049
Number of outlier coordinates lying outside NY boundaries: 223161
Number of outliers from trip times analysis: 27670
Number of outliers from trip distance analysis: 81902
Number of outliers from speed analysis: 22437
Number of outliers from fare analysis: 5476
Total outliers removed 308177
---
Estimating clusters..
Final groupbying..
Return with trip times..
Remove outliers..
Number of pickup records = 12210952
Number of outlier coordinates lying outside NY boundaries: 232444
Number of outliers from trip times analysis: 30868
Number of outliers from trip distance analysis: 87318
Number of outliers from speed analysis: 23889
Number of outliers from fare analysis: 5859
Total outliers removed 324635
---
Estimating clusters..
Final groupbying..
```

Smoothing

```
In [20]: # 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 [21]: # 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 [22]: # for each cluster number of 10min intravels with 0 pickups
         for i in range(40):
             print("for the ",i,"th cluster number of 10min intavels with zero pickups: ",4464 - len(set(jan_2015_unique[
i])))
             print('- '*60)
```

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

```
for the 20 th cluster number of 10min intavels with zero pickups: 53
-----
for the 21 th cluster number of 10min intavels with zero pickups: 29
-----
for the 22 th cluster number of 10min intavels with zero pickups: 29
-----
for the 23 th cluster number of 10min intavels with zero pickups: 163
-----
for the 24 th cluster number of 10min intavels with zero pickups: 35
-----
for the 25 th cluster number of 10min intavels with zero pickups: 41
-----
for the 26 th cluster number of 10min intavels with zero pickups: 31
-----
for the 27 th cluster number of 10min intavels with zero pickups: 214
-----
for the 28 th cluster number of 10min intavels with zero pickups: 36
-----
for the 29 th cluster number of 10min intavels with zero pickups: 41
-----
for the 30 th cluster number of 10min intavels with zero pickups: 1180
-----
for the 31 th cluster number of 10min intavels with zero pickups: 42
-----
for the 32 th cluster number of 10min intavels with zero pickups: 44
-----
for the 33 th cluster number of 10min intavels with zero pickups: 43
-----
for the 34 th cluster number of 10min intavels with zero pickups: 39
-----
for the 35 th cluster number of 10min intavels with zero pickups: 42
-----
for the 36 th cluster number of 10min intavels with zero pickups: 36
-----
for the 37 th cluster number of 10min intavels with zero pickups: 321
-----
for the 38 th cluster number of 10min intavels with zero pickups: 36
-----
for the 39 th cluster number of 10min intavels with zero pickups: 43
-----
```

there are two ways to fill up these values

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

```
In [23]: # Fills a value of zero for every bin where no pickup data is present
# the count_values: number pickps 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 [23]: def smoothing1(numberOfPickups, correspondingTimeBin):
    ind = 0
    repeat = 0
    smoothed_region = []
    for cluster in range(0, 40):
        smoothed_bin = []
        for t1 in range(4464):
            if repeat != 0: #this will ensure that we shall not fill the pickup values again which we already
filled by smoothing
                repeat -= 1
            else:
                if t1 in correspondingTimeBin[cluster]:
                    smoothed_bin.append(numberOfPickups[ind])
                    ind += 1
                else:
                    if t1 == 0:
#<-----CASE-1:Pickups missing in the beginning----->
                        for t2 in range(t1, 4464):
                            if t2 not in correspondingTimeBin[cluster]:
                                continue
                            else:
                                right_hand_limit = t2
                                smoothed_value = (numberOfPickups[ind]*1.0)/((right_hand_limit + 1)*1.0)
                                for i in range(right_hand_limit + 1):
                                    smoothed_bin.append(math.ceil(smoothed_value))
                                ind += 1
                                repeat = right_hand_limit - t1

                    if t1 != 0:
                        right_hand_limit = 0
                        for t2 in range(t1, 4464):
                            if t2 not in correspondingTimeBin[cluster]:
                                continue
                            else:
                                right_hand_limit = t2
                                break
                        if right_hand_limit == 0:
#<-----CASE-2: Pickups MISSING IN THE END----->
                            smoothed_value = (numberOfPickups[ind-1]*1.0)/(((4464 - t1)+1)*1.0)

```



```
        del smoothed_bin[-1]
        for i in range((4464 - t1)+1):
            smoothed_bin.append(math.ceil(smoothed_value))
        repeat = (4464 - t1) - 1
        #<-----CASE-3: Pickups MISSING IN MIDDLE OF TWO VALUES----->
        else:
            smoothed_value = ((numberOfPickups[ind-1] + numberOfPickups[ind])*1.0)/(((right_hand
_limit - t1)+2)*1.0)
            del smoothed_bin[-1]
            for i in range((right_hand_limit - t1)+2):
                smoothed_bin.append(math.ceil(smoothed_value))
            ind += 1
            repeat = right_hand_limit - t1
        smoothed_region.extend(smoothed_bin)
    return smoothed_region
```

```

In [24]: # Fills a value of zero for every bin where no pickup data is present
# the count_values: number pickups that are happened in each region for each 10min intravel
# there wont be any value if there are no pickups.
# values: number of unique bins

# for every 10min intravel(pickup_bin) we will check it is there in our unique bin,
# if it is there we will add the count_values[index] to smoothed data
# if not we add smoothed data (which is calculated based on the methods that are discussed in the above markdown
cell)
# we finally return smoothed data
def smoothing(count_values,values):
    smoothed_regions=[] # stores list of final smoothed values of each 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
                        else:
                            right_hand_limit=j
                            break
                    if right_hand_limit==0:
                        #Case 1: When we have the last/last few values are found to be missing,hence we have no righ
t-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 le
ft-limit here
        right_hand_limit=0
        for j in range(i,4464):
            if j not in values[r]:
                continue
            else:
                right_hand_limit=j
                break
        smoothed_value=count_values[ind]*1.0/((right_hand_limit-i)+1)*1.0
        for j in range(i,right_hand_limit+1):
            smoothed_bins.append(math.ceil(smoothed_value))
        repeat=(right_hand_limit-i)

    ind+=1
    smoothed_regions.extend(smoothed_bins)
    return smoothed_regions

```

```
In [25]: #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)
print("fill_missing completed")
#Smoothing Missing values of Jan-2015
jan_2015_smooth = smoothing(jan_2015_groupby['trip_distance'].values,jan_2015_unique)
print("smoothing completed")
```

```
fill_missing completed
smoothing completed
```

```
In [26]: # 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 [27]: def countZeros(num):
        count = 0
        for i in num:
            if i == 0:
                count += 1
        return count
```

```
In [28]: print("Number of values filled with zero in zero fill data= "+str(countZeros(jan_2015_fill)))
```

```
Number of values filled with zero in zero fill data= 9451
```

```
In [29]: print("Sanity check for number of zeros in smoothed data = "+str(countZeros(jan_2015_smooth)))
```

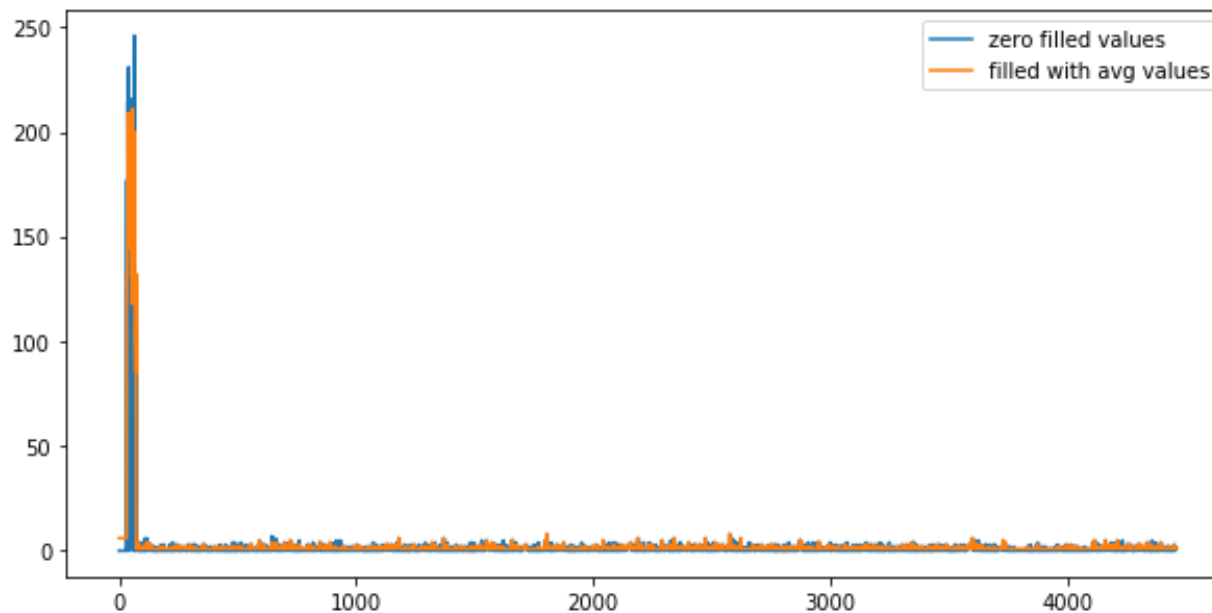
```
Sanity check for number of zeros in smoothed data = 0
```

```
In [30]: print("Total number of pickup values = "+str(len(jan_2015_fill)))  
print("Total number of pickup values = "+str(len(jan_2015_smooth)))
```

Total number of pickup values = 178560

Total number of pickup values = 178560

```
In [31]: # 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 [0]: *# 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  
# when you are using smoothing we are looking at the future number of pickups which might cause a data leakage.  
  
# so we use smoothing for jan 2015th data since it acts as our training data  
# and we use simple fill_misssing method for 2016th data.
```

```

In [33]: # 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 =  $24 \times 31 \times 60 / 10 = 4464$ 
# number of 10min indices for jan 2016 =  $24 \times 31 \times 60 / 10 = 4464$ 
# number of 10min indices for feb 2016 =  $24 \times 29 \times 60 / 10 = 4176$ 
# number of 10min indices for march 2016 =  $24 \times 31 \times 60 / 10 = 4464$ 
# 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
regions_sum = np.array(regions_cum)
#np.save("regions_cum.npy", regions_cum)

```

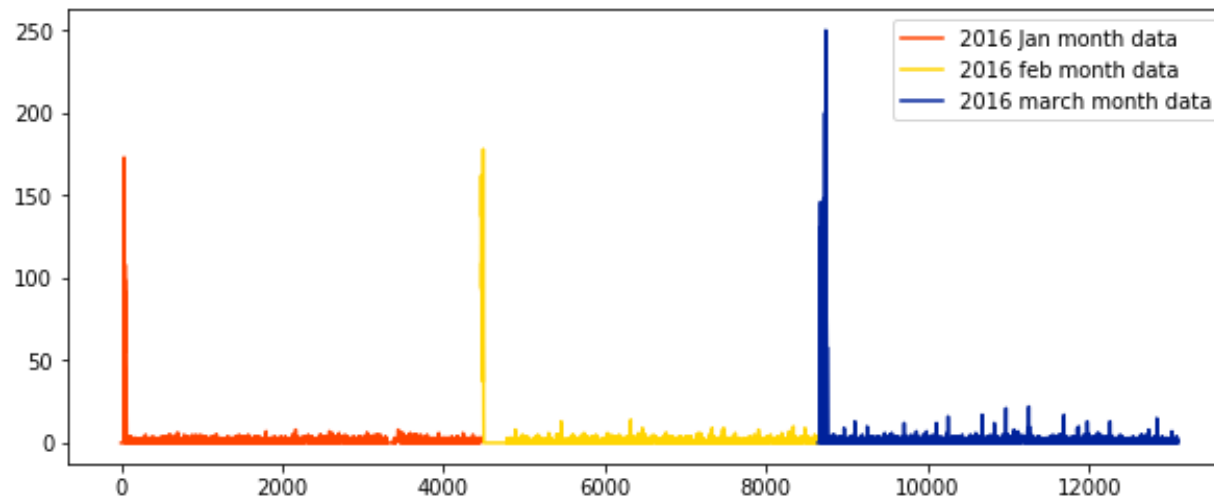
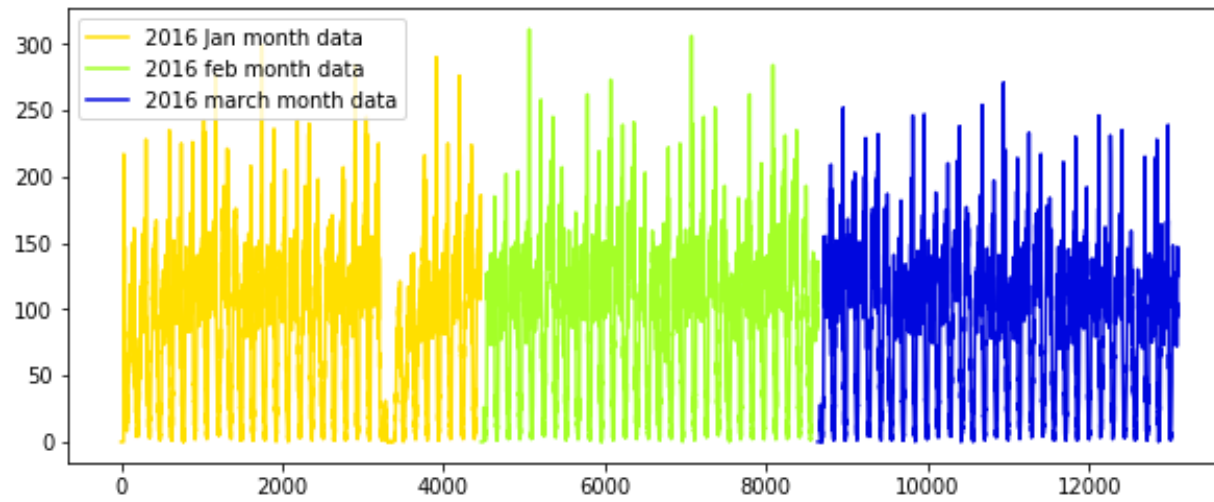
```

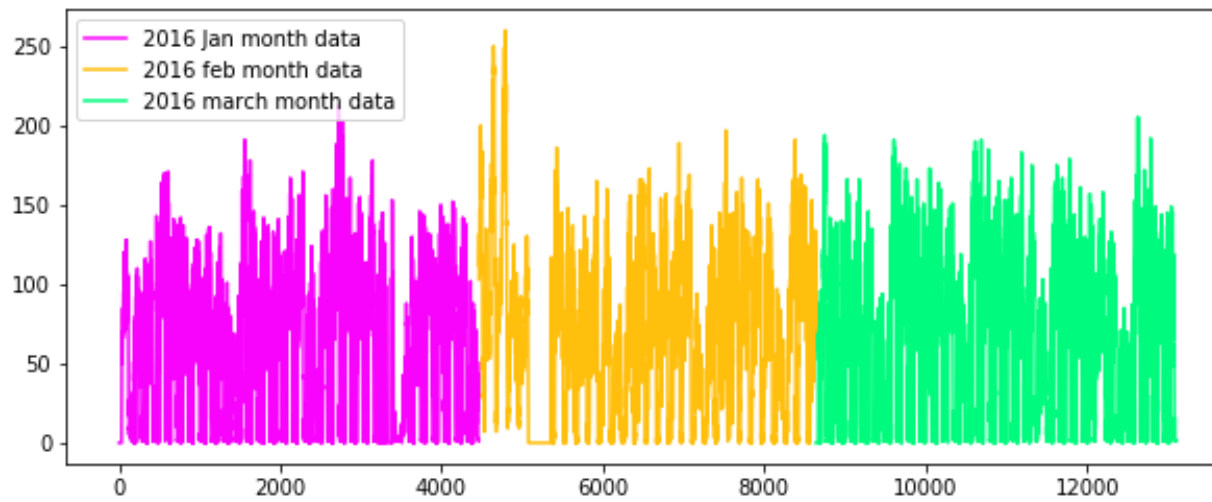
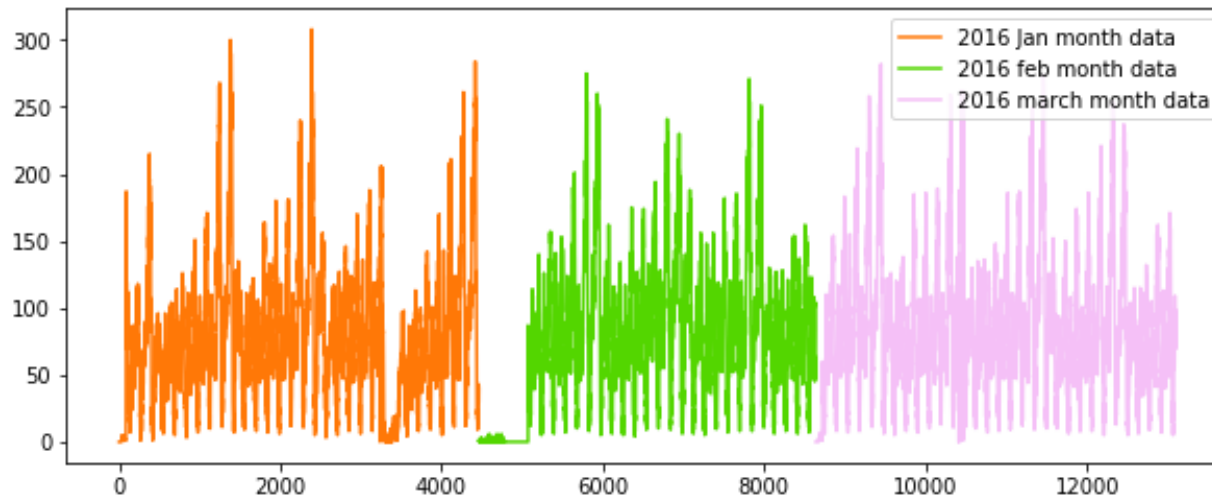
In [53]: jan_2016_fill = fill_missing(jan_2016_groupby['trip_distance'].values, jan_2016_unique)
feb_2016_fill = fill_missing(feb_2016_groupby['trip_distance'].values, feb_2016_unique)
mar_2016_fill = fill_missing(mar_2016_groupby['trip_distance'].values, mar_2016_unique)

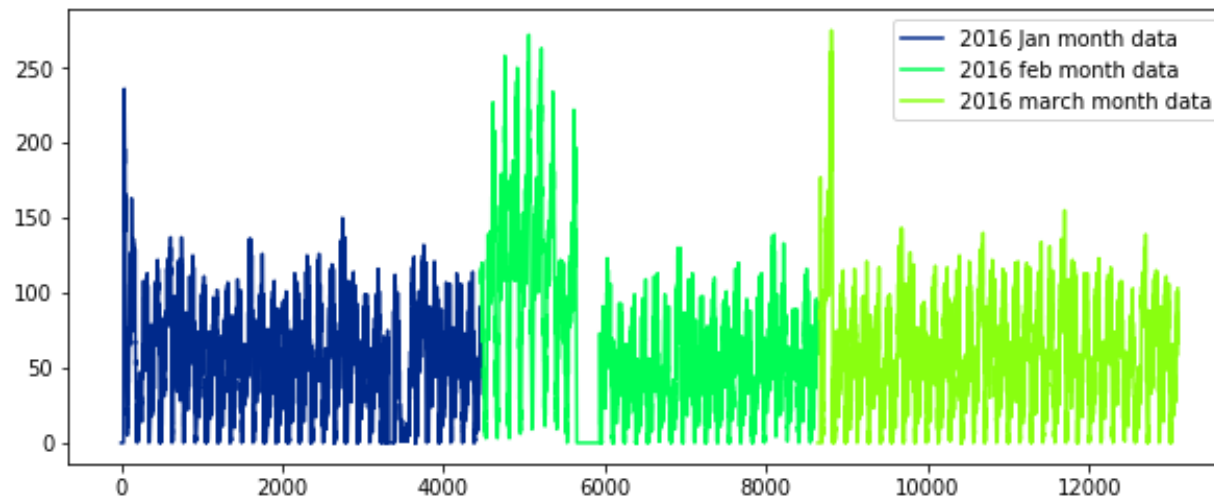
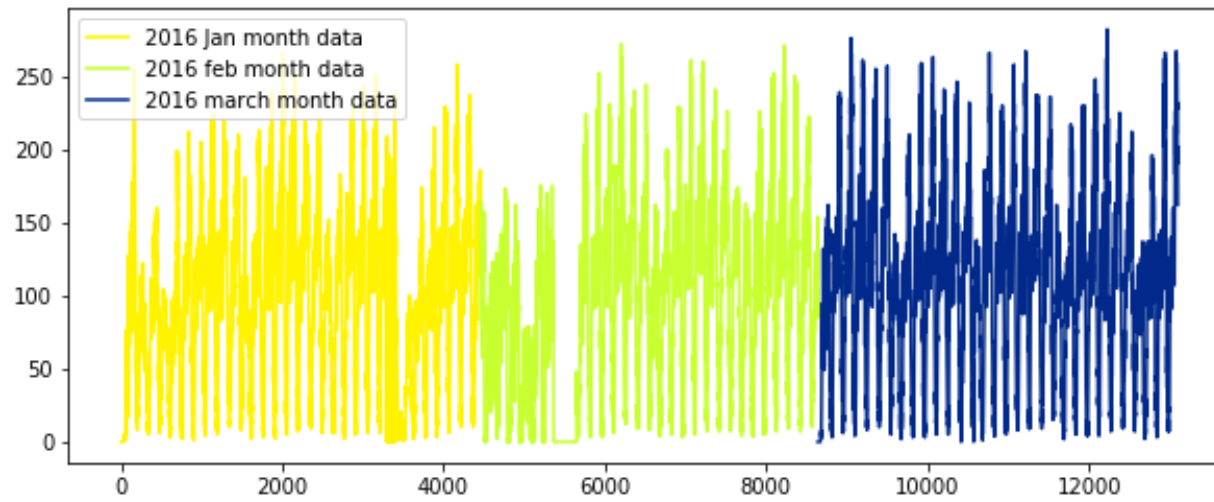
```

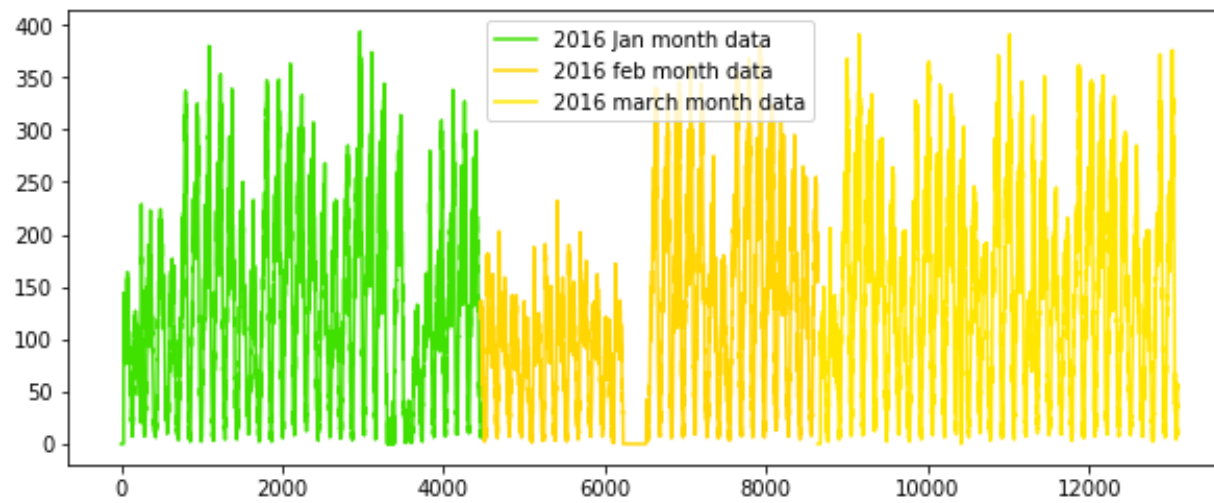
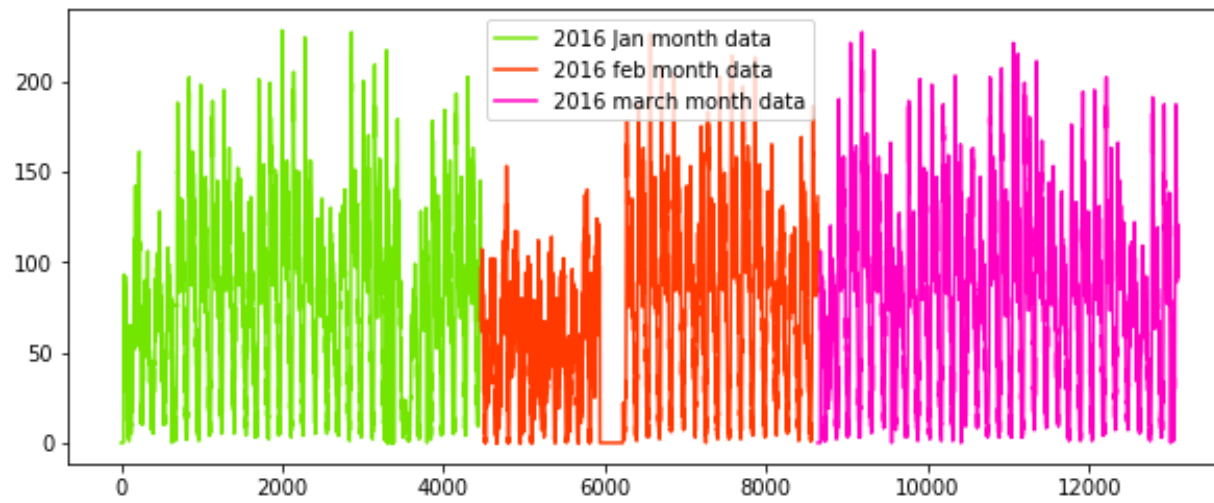
Time series and Fourier Transforms

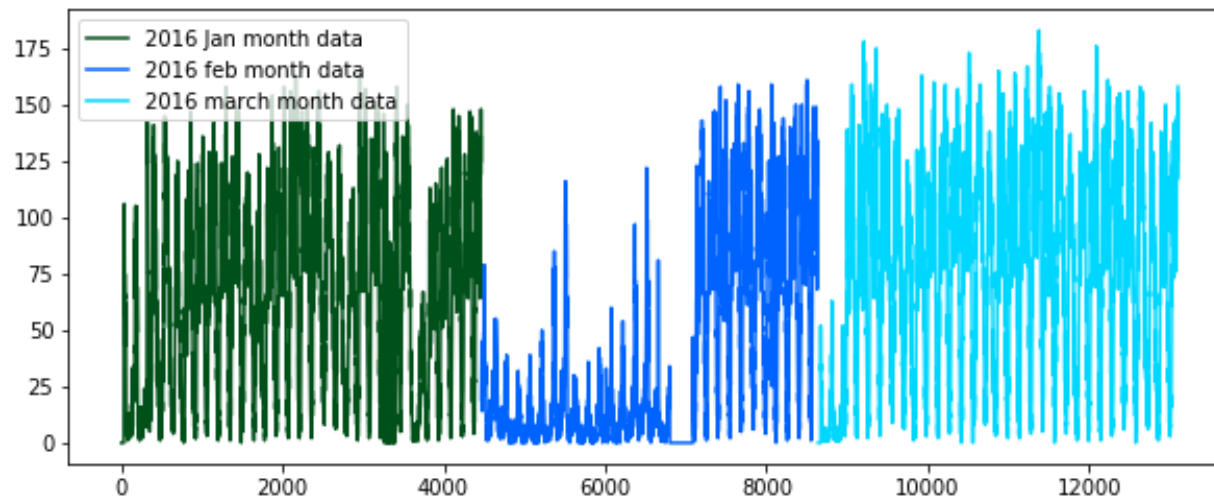
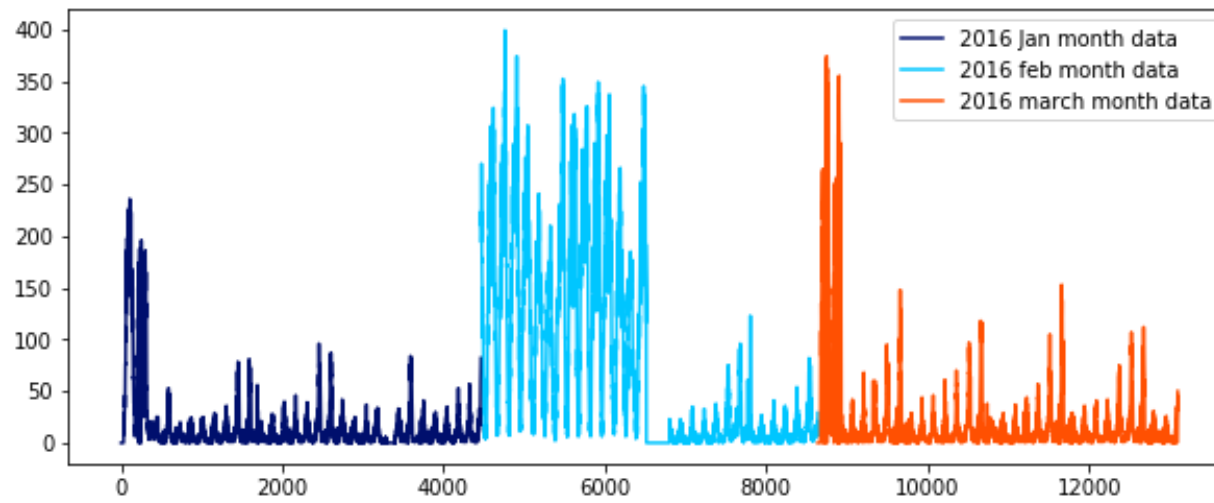

```
In [34]: def uniqueish_color():  
        """There're better ways to generate unique colors, but this isn't awful."""  
        return plt.cm.gist_ncar(np.random.random())  
first_x = list(range(0,4464))  
second_x = list(range(4464,8640))  
third_x = list(range(8640,13104))  
for i in range(40):  
    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()
```

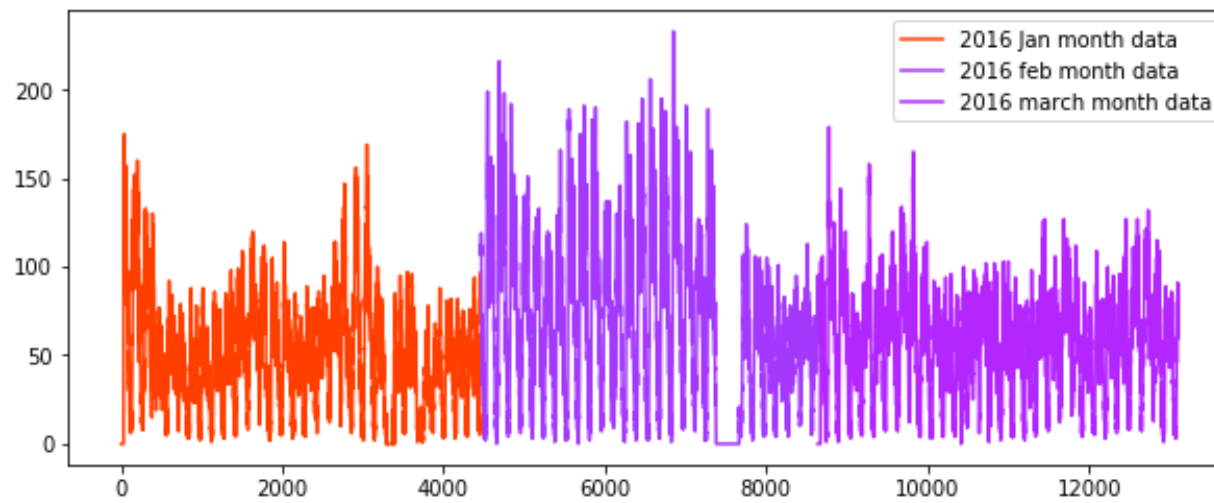
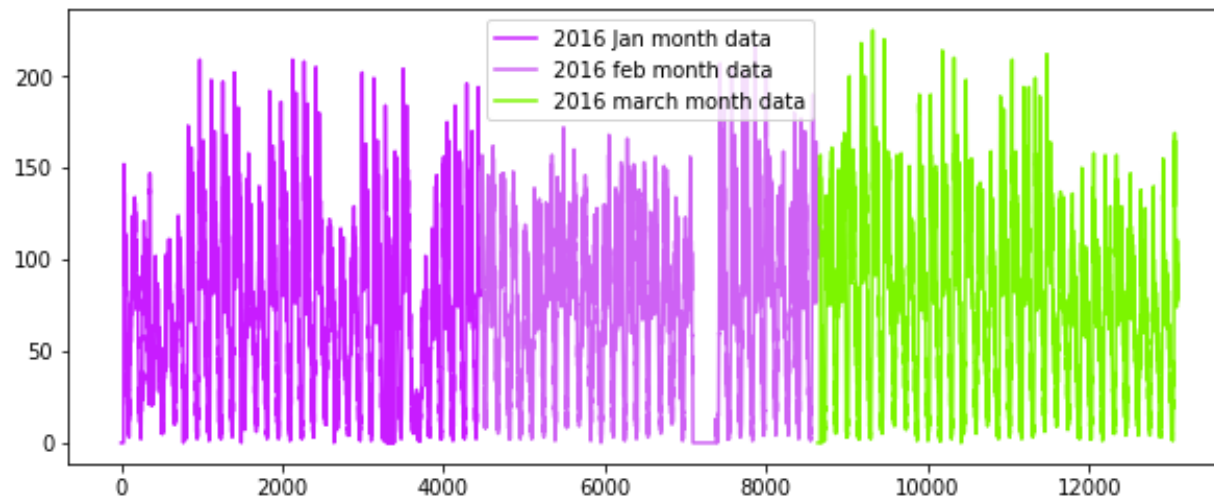


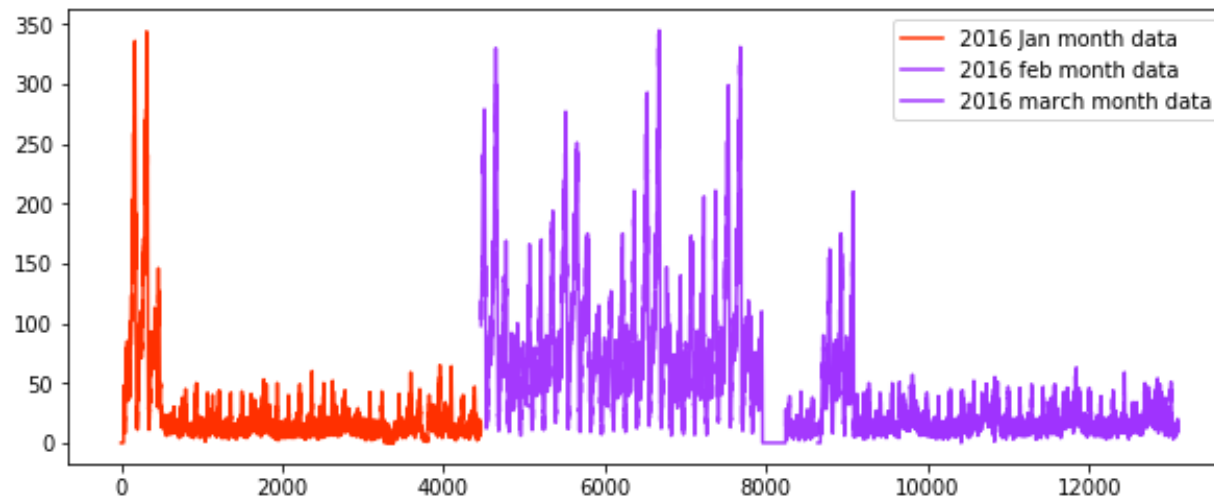
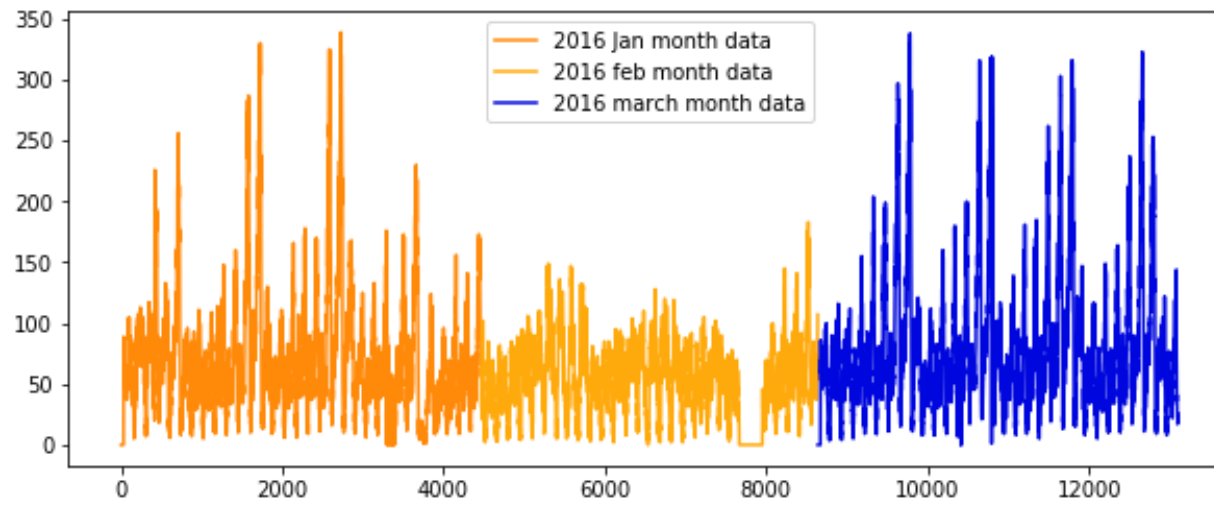


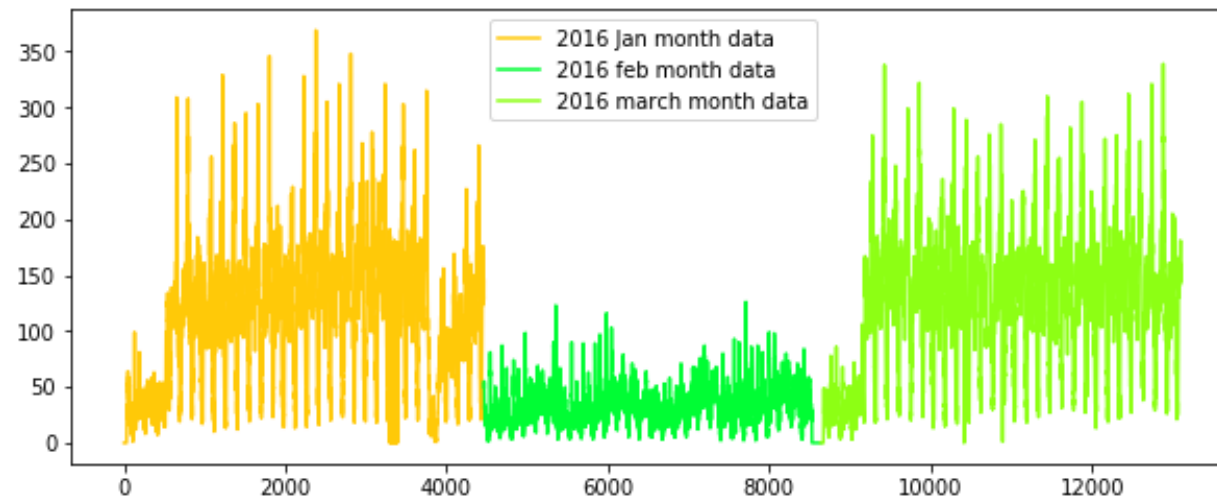
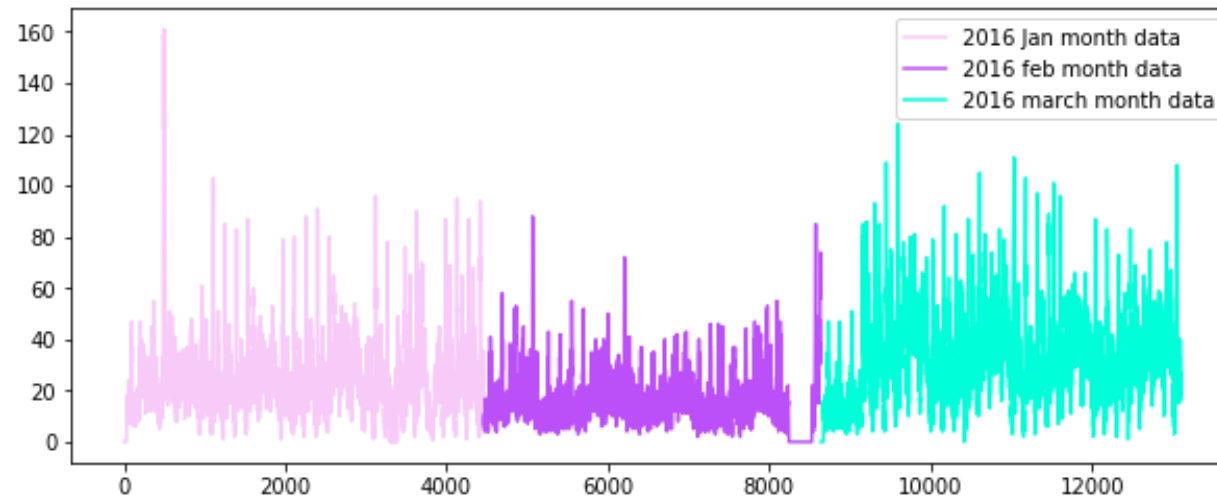


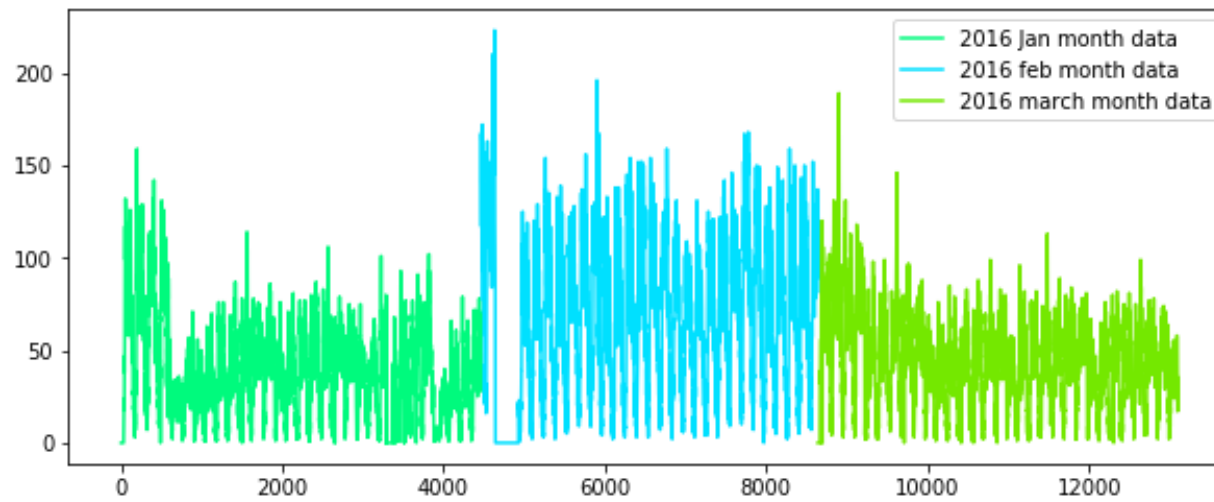
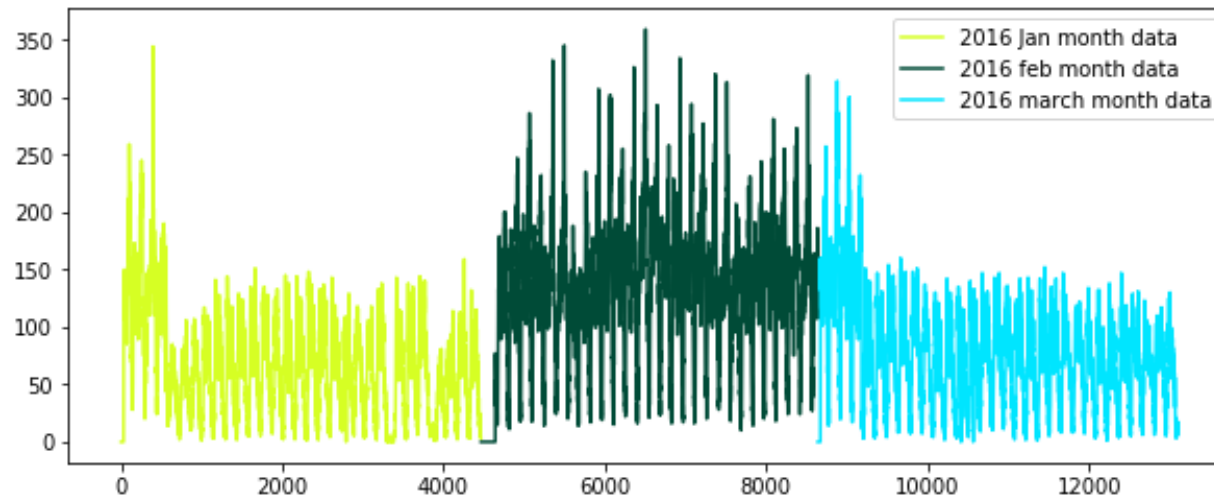


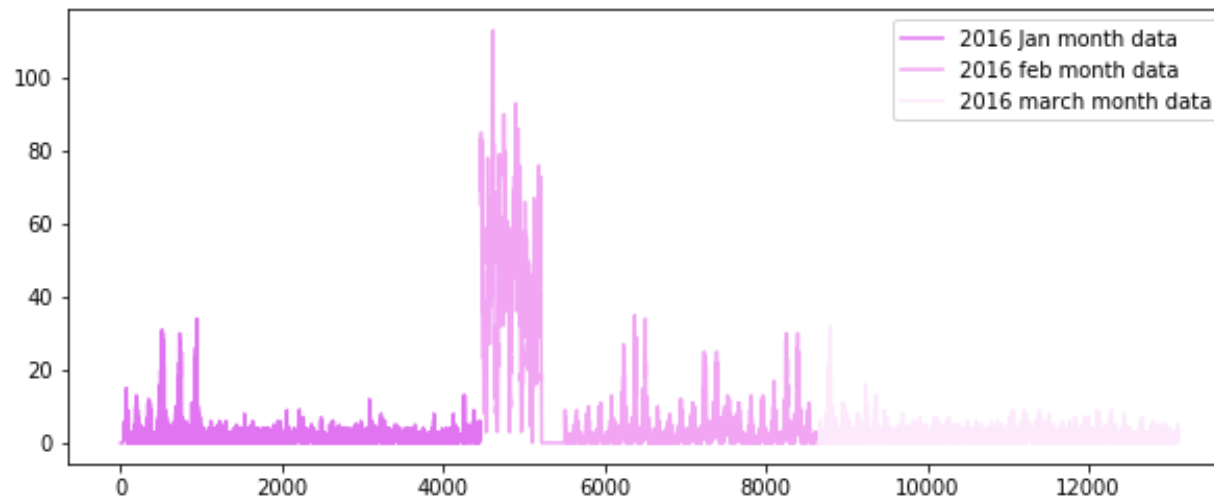
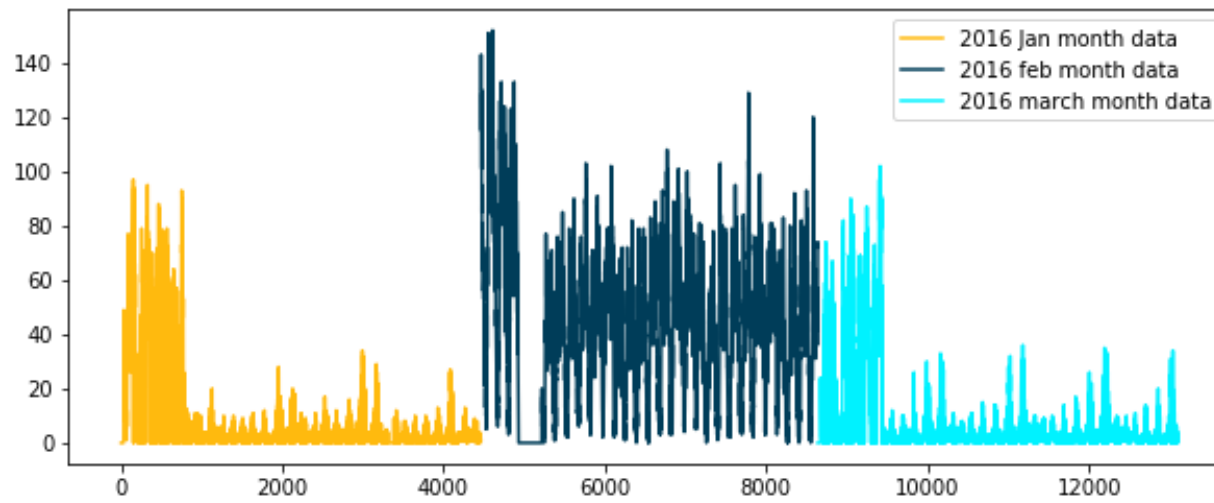


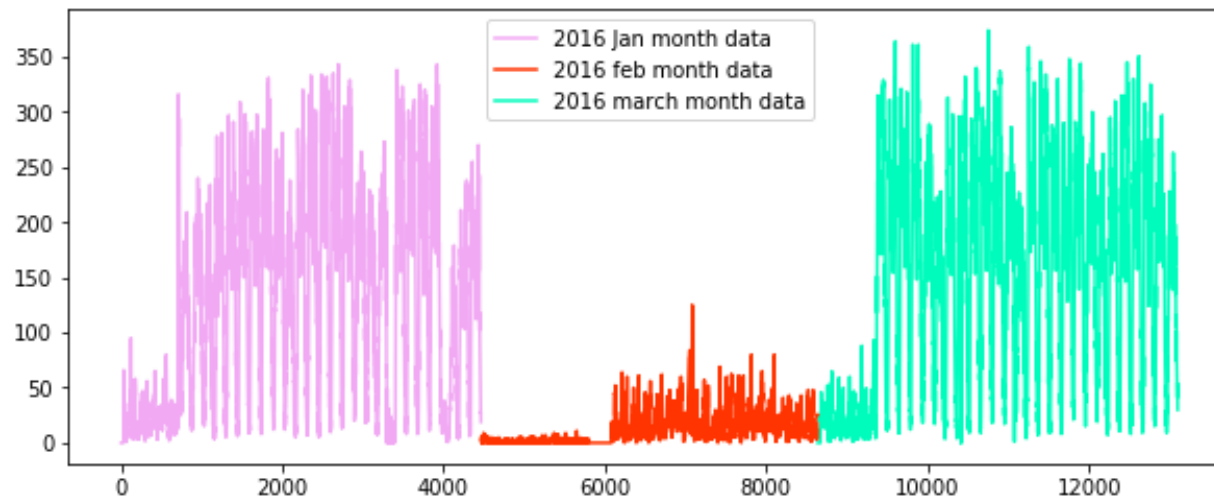
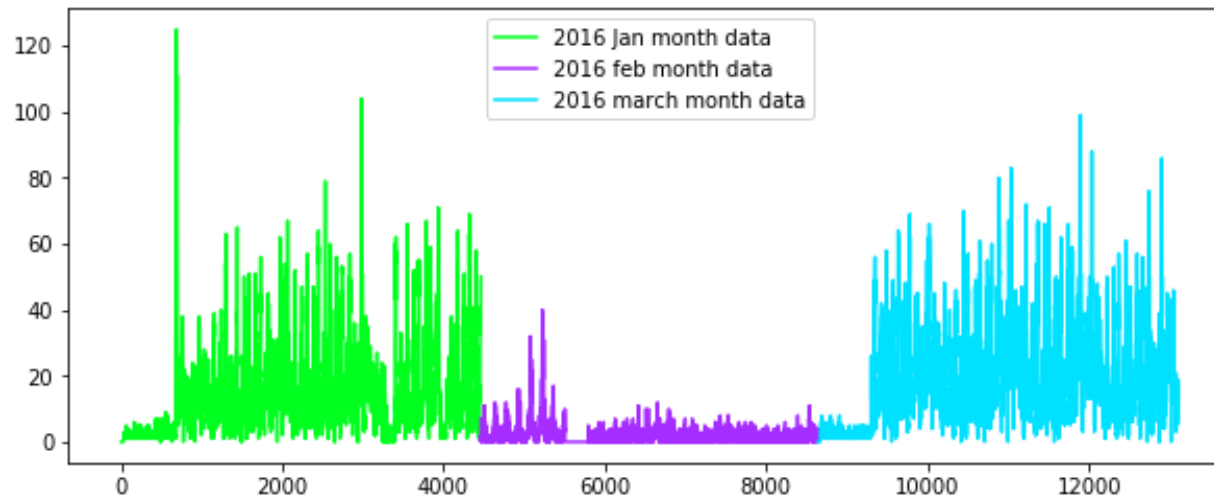


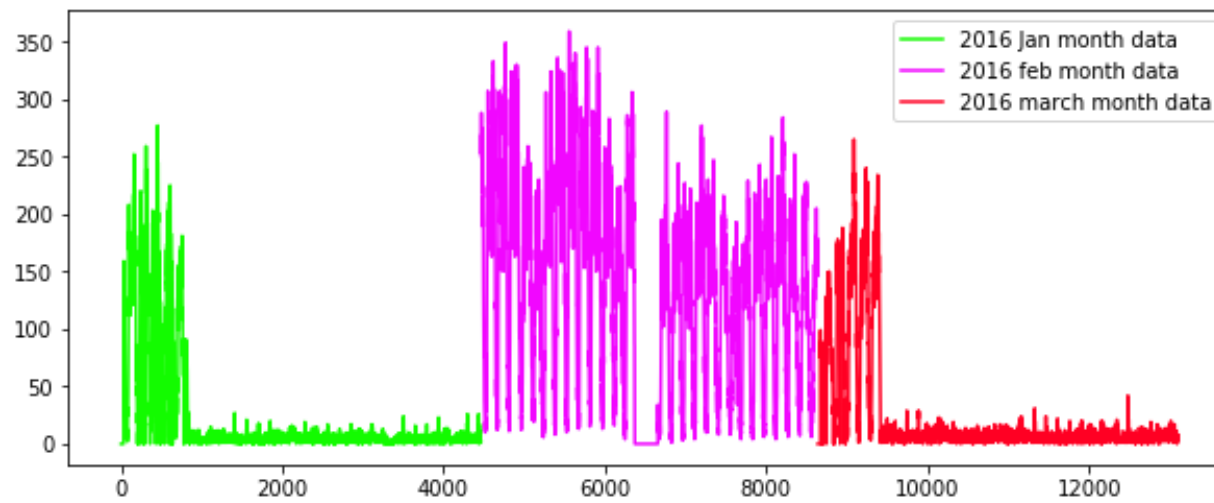
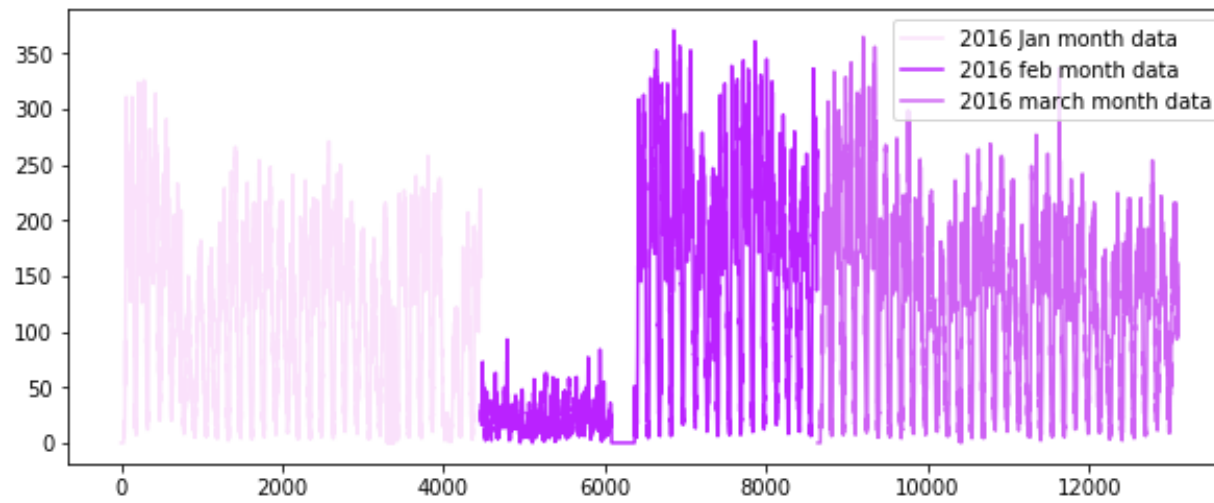


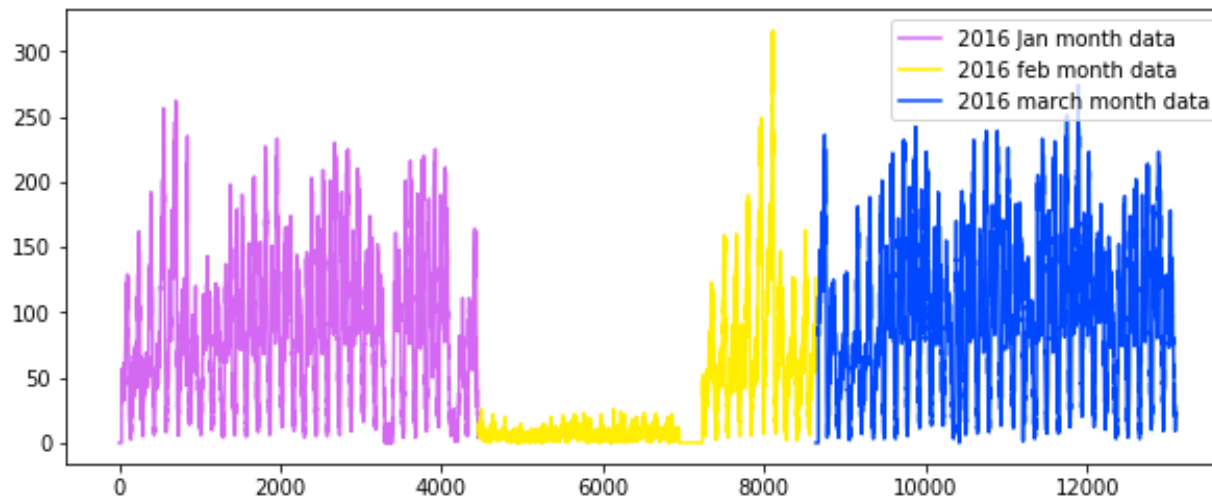
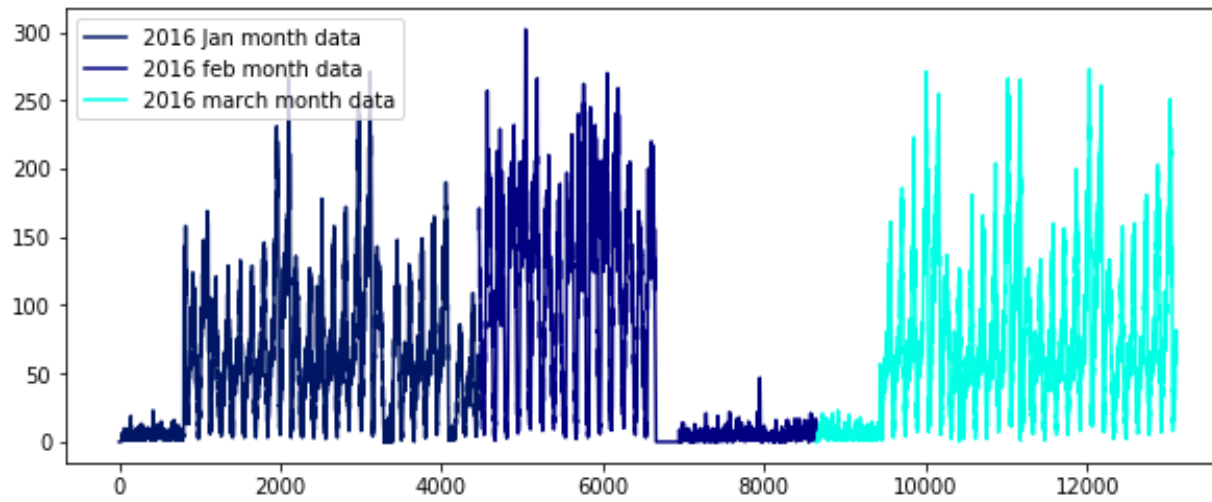


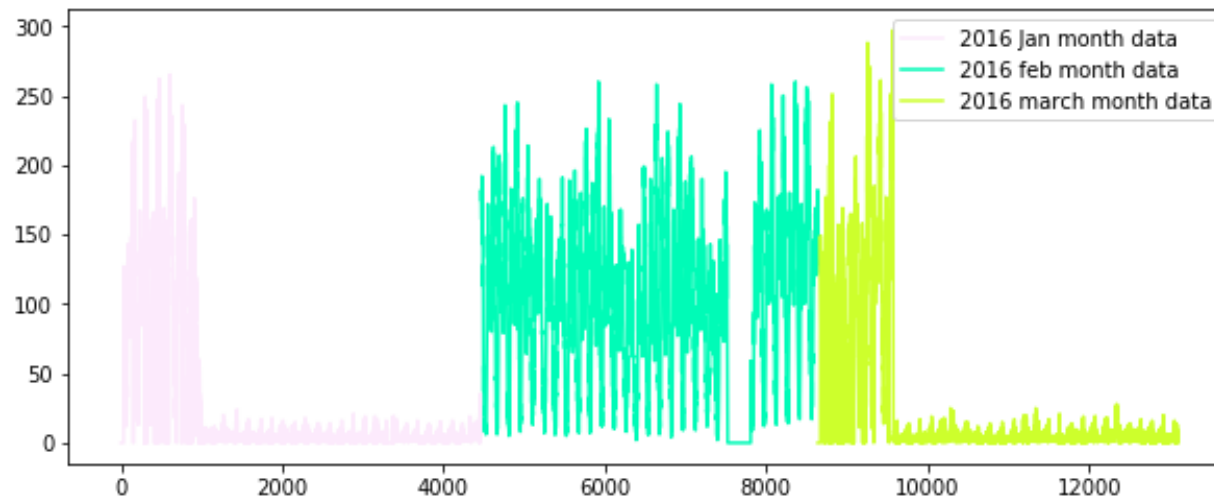
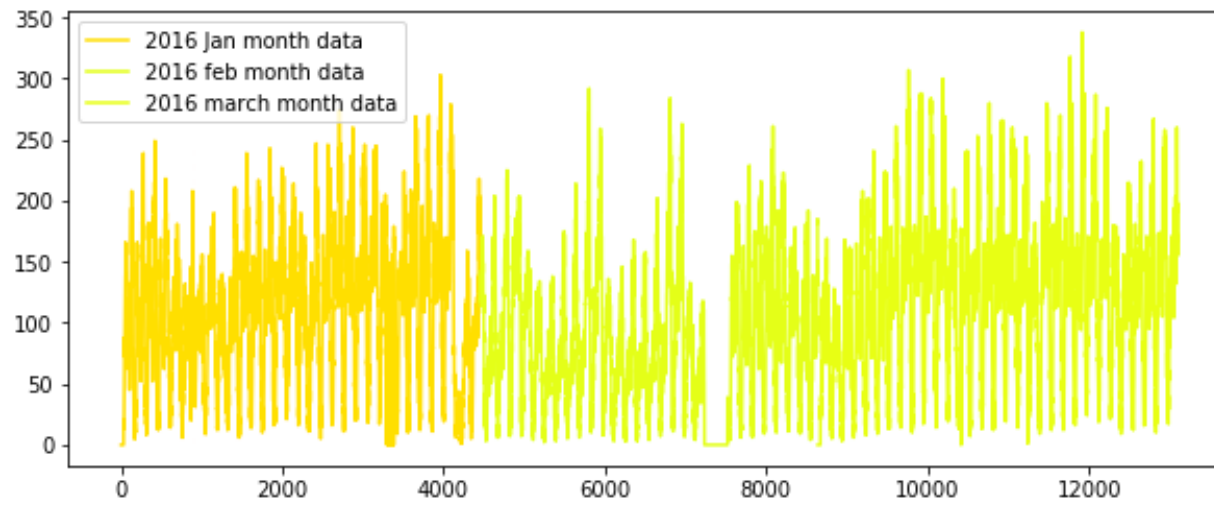


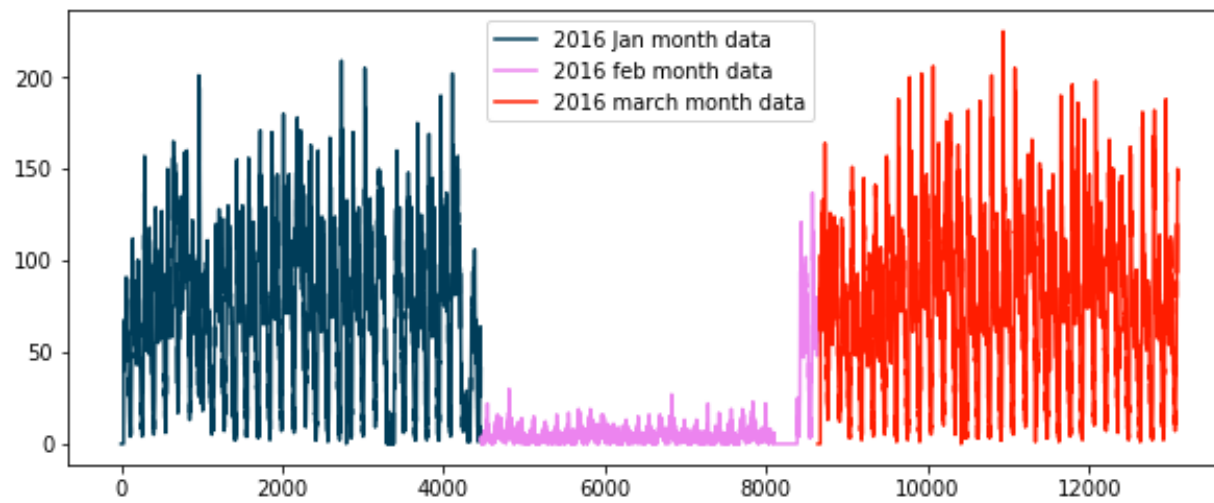
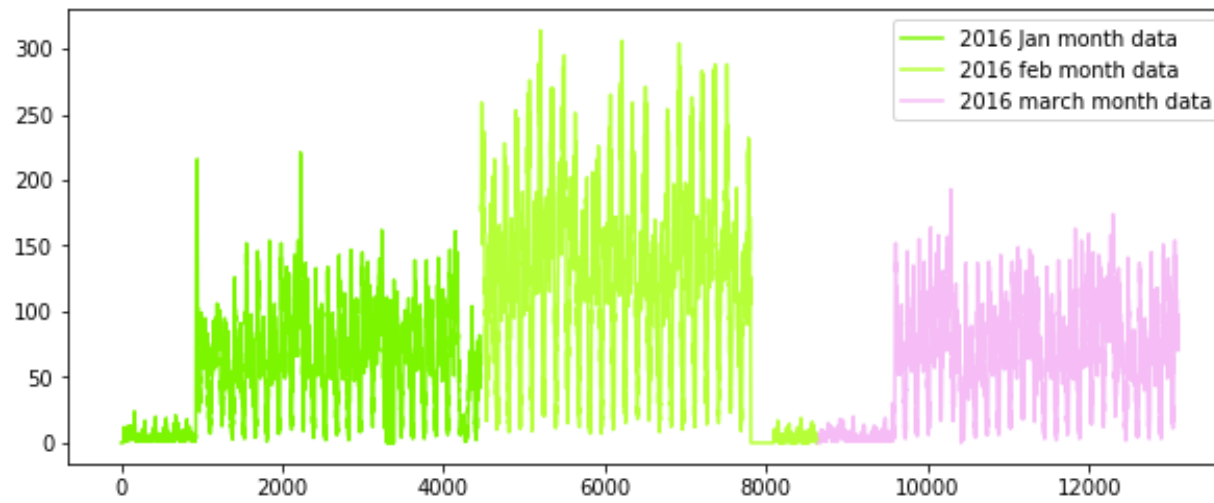


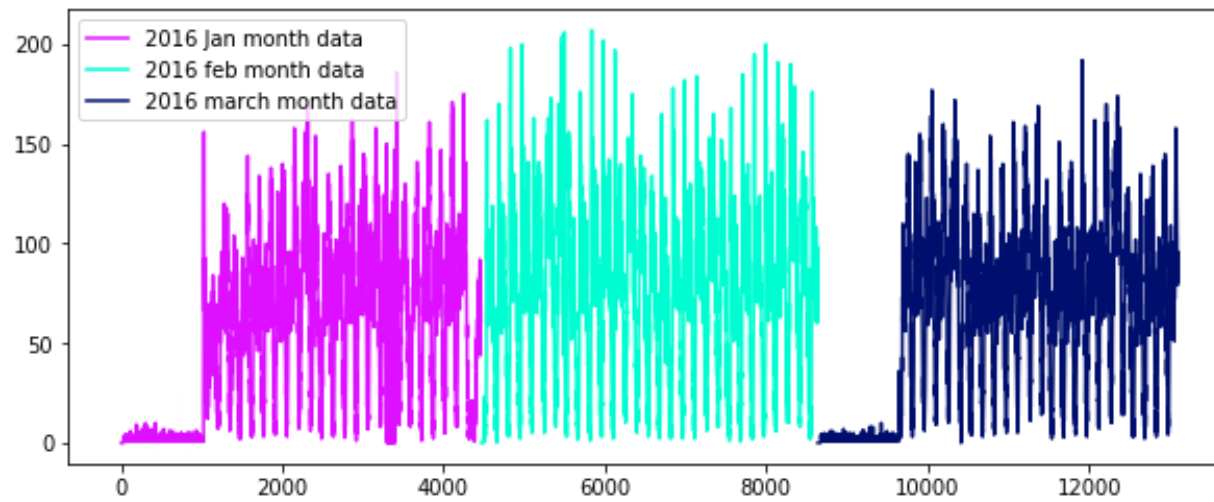
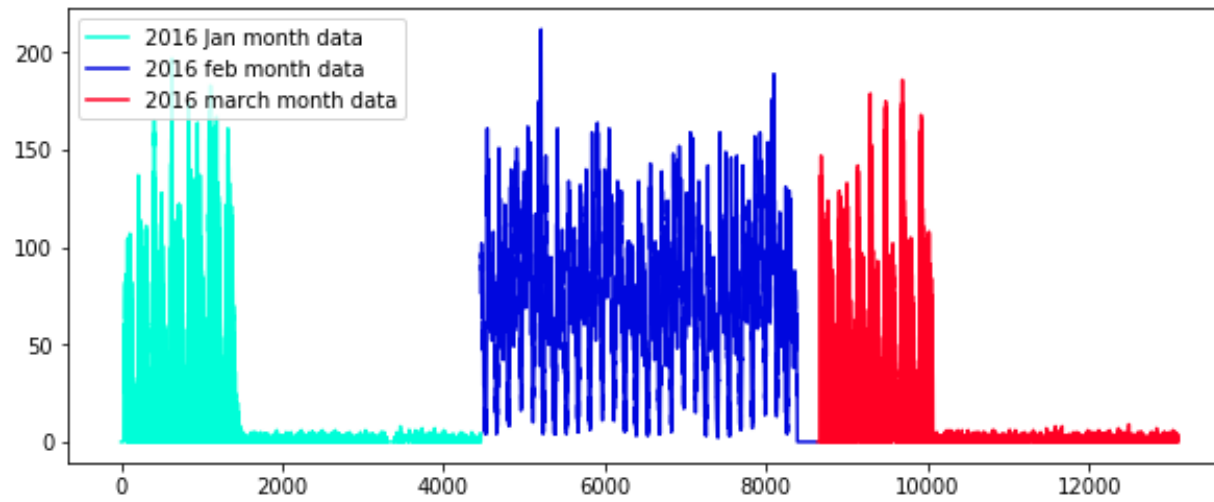


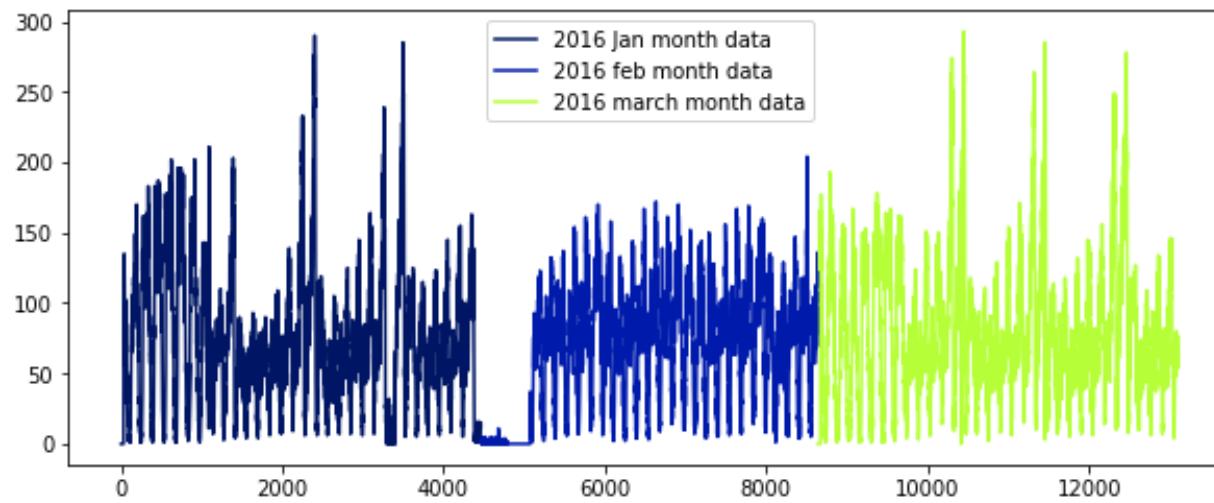
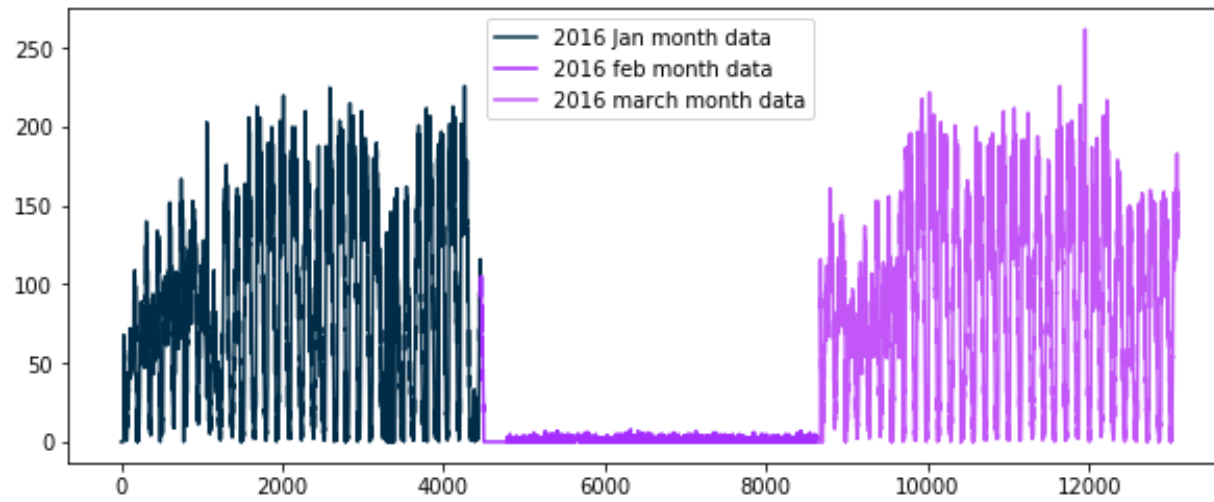


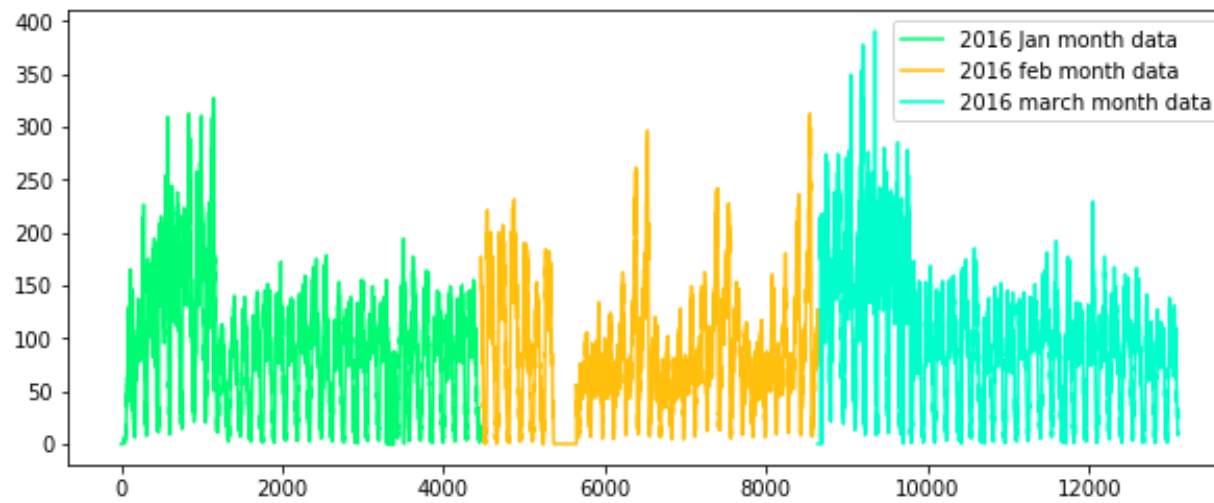
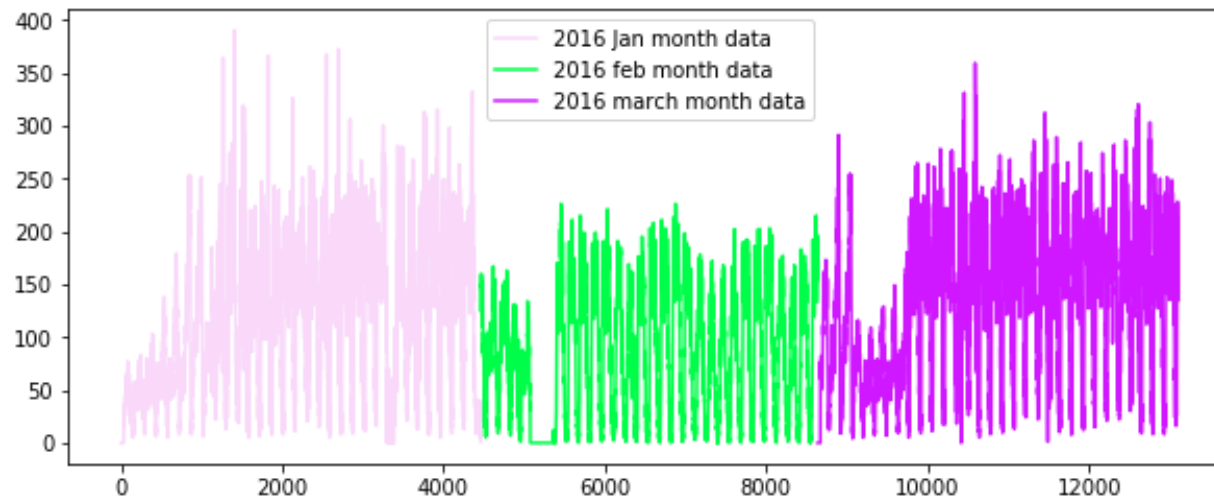


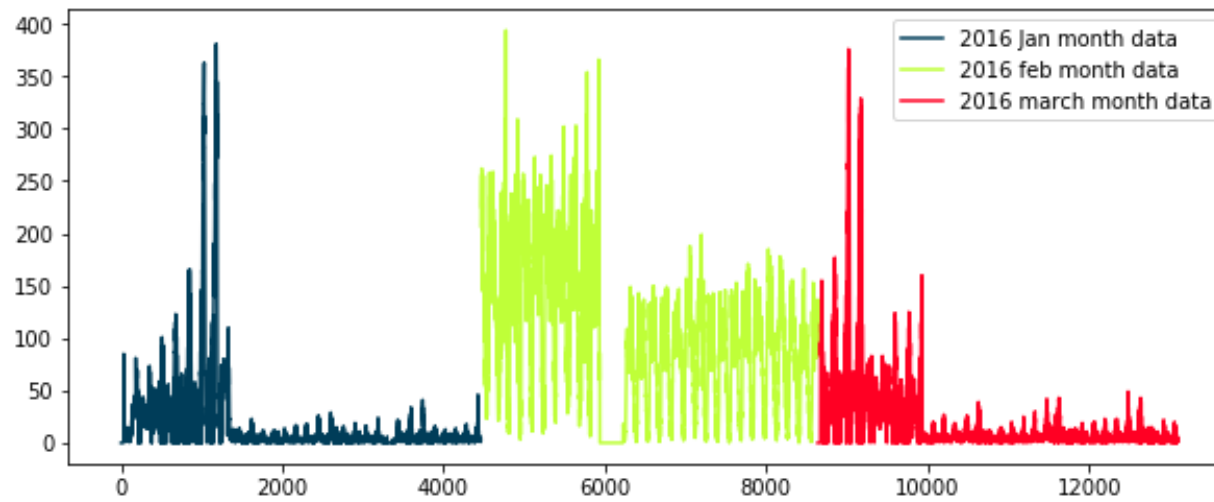
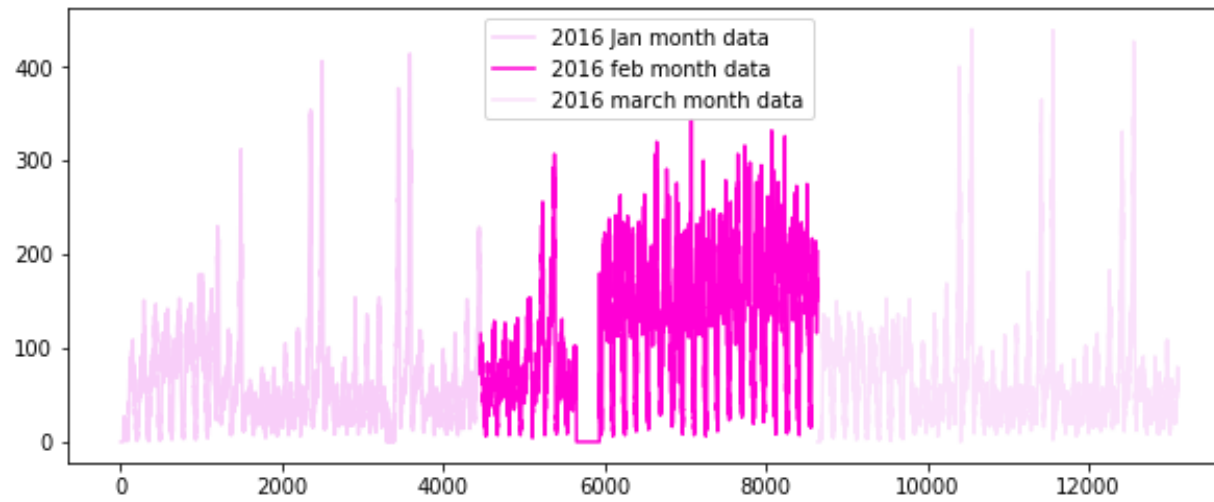


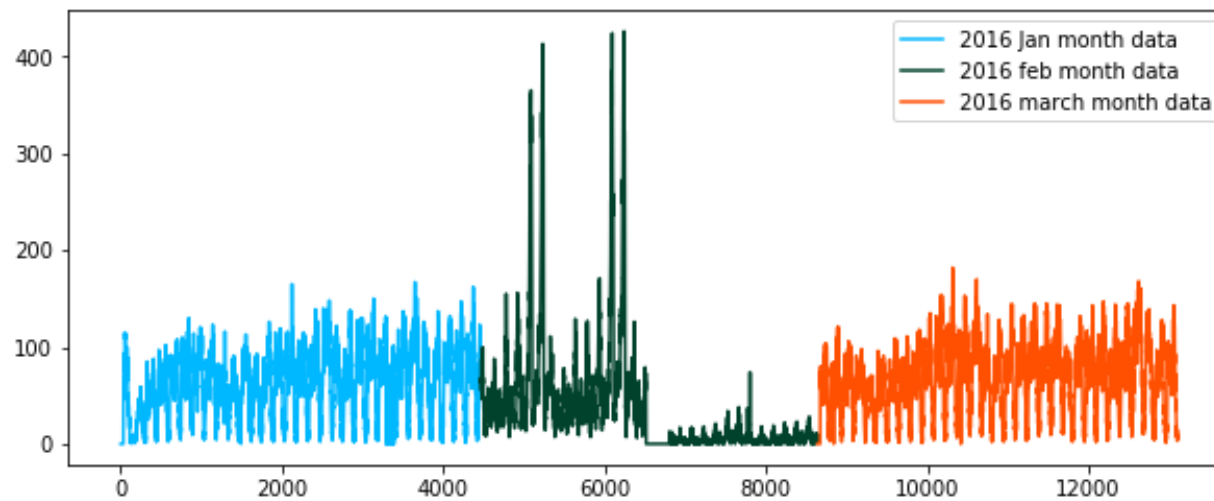
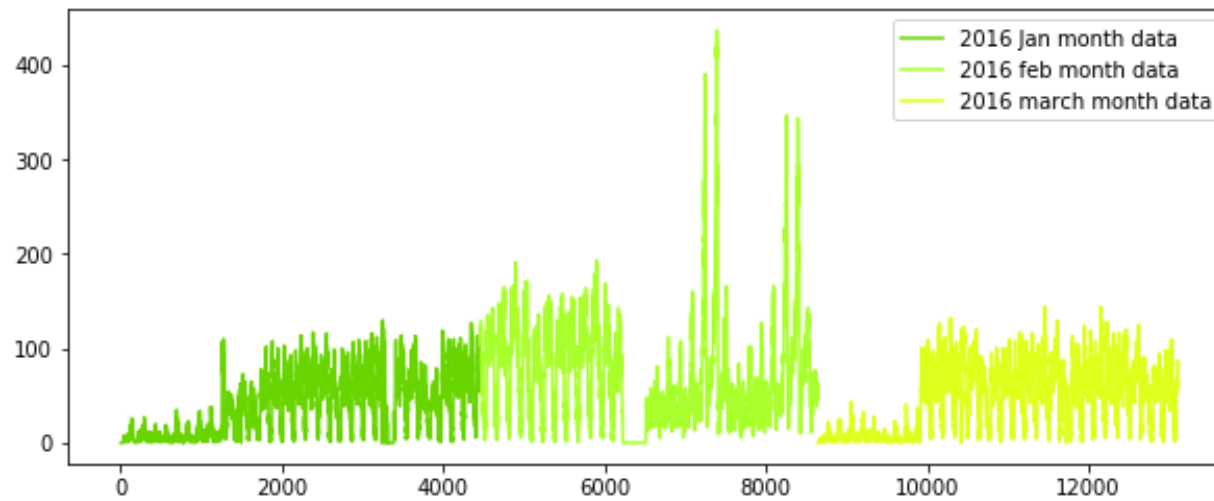




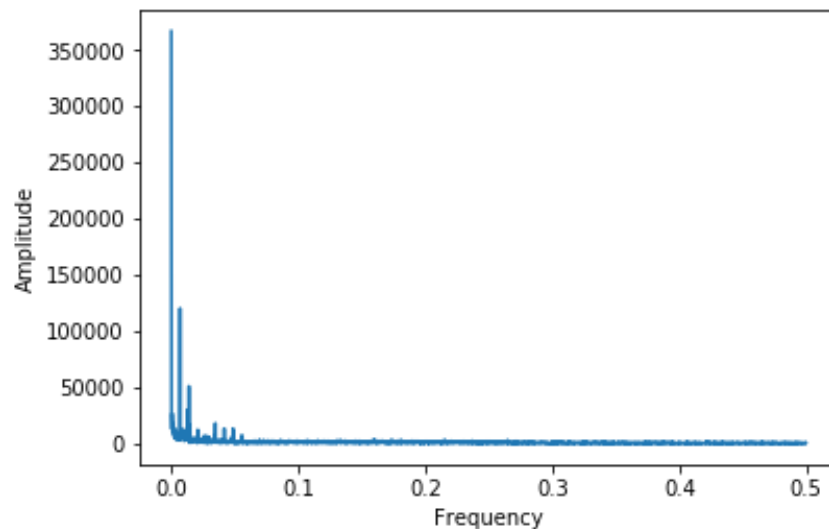








```
In [35]: # getting peaks: https://blog.ytotech.com/2015/11/01/findpeaks-in-python/
# read more about fft function : https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fft.html
Y      = np.fft.fft(np.array(jan_2016_smooth)[0:4460])
# read more about the fftfreq: https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fftfreq.html
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 [36]: #Preparing the Dataframe only with x(i) values as jan-2015 data and y(i) values as jan-2016
ratios_jan = pd.DataFrame()
ratios_jan['Given']=jan_2015_smooth
ratios_jan['Prediction']=jan_2016_smooth
ratios_jan['Ratios']=ratios_jan['Prediction']*1.0/ratios_jan['Given']*1.0
```

In [37]: `ratios_jan.tail()`

Out[37]:

	Given	Prediction	Ratios
178555	100	83	0.830000
178556	97	83	0.855670
178557	105	86	0.819048
178558	103	79	0.766990
178559	111	84	0.756757

In [38]: `print("Total Number of zeros in Prediction = "+str(ratios_jan["Prediction"].value_counts()[0]))`

Total Number of zeros in Prediction = 11250

In [39]: `print("Total Number of zeros in ratios = "+str(ratios_jan["Ratios"].value_counts()[0]))`

Total Number of zeros in ratios = 11250

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
$$R_t = P^{2016}_t / P^{2015}_t$$
2. Using Previous known values of the 2016 data itself to predict the future values

Simple Moving Averages

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

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

```
In [41]: 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
$$R_t = (R_{t-1} + R_{t-2} + R_{t-3})/3$$

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

```
In [42]: def MA_P_Predictions(ratios,month):
    predicted_value=(ratios['Prediction'].values)[0]
    error=[]
    predicted_values=[]
    window_size=1
    predicted_ratio_values=[]
    for i in range(0,4464*40):
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
        if i+1>=window_size:
            predicted_value=int(sum((ratios['Prediction'].values)[(i+1)-window_size:(i+1)])/window_size)
        else:
            predicted_value=int(sum((ratios['Prediction'].values)[0:(i+1)])/(i+1))

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

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

Weighted Moving Averages

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

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

```

In [43]: 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)/len(ratios['Prediction'].values))
    mse_err = sum([e**2 for e in error])/len(error)
    return ratios,mape_err,mse_err

```

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

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

```

In [44]: def WA_P_Predictions(ratios,month):
    predicted_value=(ratios['Prediction'].values)[0]
    error=[]
    predicted_values=[]
    window_size=2
    for i in range(0,4464*40):
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
        if i+1>=window_size:
            sum_values=0
            sum_of_coeff=0
            for j in range(window_size,0,-1):
                sum_values += j*(ratios['Prediction'].values)[i-window_size+j]
                sum_of_coeff+=j
            predicted_value=int(sum_values/sum_of_coeff)

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

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

```

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

Exponential Weighted Moving Averages

https://en.wikipedia.org/wiki/Moving_average#Exponential_moving_average (https://en.wikipedia.org/wiki/Moving_average#Exponential_moving_average)

Through weighted averaged we have satisfied the analogy of giving higher weights to the latest value and decreasing weights to the subsequent ones but we still do not know which is the correct weighting scheme as there are infinitely many possibilities in which we can assign weights in a non-increasing order and tune the the hyperparameter window-size. To simplify this process we use Exponential Moving Averages which is a more logical way towards assigning weights and at the same time also using an optimal window-size.

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

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

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

```

In [45]: 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

```

$$P^{\{t\}}_{\{t\}} = \alpha * P^{\{t-1\}}_{\{t-1\}} + (1-\alpha) * P^{\{t\}}_{\{t-1\}}$$

```
In [46]: 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 [47]: 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 [48]: 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 ("")
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 ("")
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:  0.22785156353133512      MSE:  1196.2953853046595
Moving Averages (2016 Values) -                         MAPE:  0.15583458712025738      MSE:  254.6630936379928
3
-----
Weighted Moving Averages (Ratios) -                     MAPE:  0.22706529144871415      MSE:  1053.083529345878
Weighted Moving Averages (2016 Values) -                 MAPE:  0.1479482182992932      MSE:  224.81054547491038
-----
Exponential Moving Averages (Ratios) -                     MAPE:  0.2275474636148534      MSE:  1019.3071012544802
Exponential Moving Averages (2016 Values) -             MAPE:  0.1475381297798153      MSE:  222.35159610215055
```

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

From the above matrix it is inferred that the best forecasting model for our prediction would be:-
$$P_t = \alpha P_{t-1} + (1-\alpha)P_{t-1}$$
 i.e Exponential Moving Averages using 2016 Values

Regression Models

Train-Test Split

Before we start predictions using the tree based regression models we take 3 months of 2016 pickup data and split it such that for every region we have 80% data in train and 20% in test, ordered date-wise for every region

```
In [49]: # 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 variable
# it is list of lists
# it will contain number of pickups 13099 for each cluster
output = []

# tsne_lat will contain 13104-5=13099 times latitude 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]
# it is list of lists
tsne_lon = []

# we will code each day
# sunday = 0, monday=1, tue = 2, wed=3, thur=4, fri=5,sat=6
# for every cluster we will be adding 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..x13
    104], [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(regi
    ons_cum[i])-number_of_time_stamps)]))
    output.append(regions_cum[i][5:])
tsne_feature = tsne_feature[1:]

```

```

In [50]: 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[50]: True

```

```

In [49]: len(tsne_feature)

```

```

Out[49]: 523960

```

```
In [50]: tsne_feature
```

```
Out[50]: array([[ 0,  0,  0,  0,  0],  
                [ 0,  0,  0,  0,  0],  
                [ 0,  0,  0,  0,  0],  
                ...,  
                [11,  7, 15,  9,  3],  
                [ 7, 15,  9,  3,  7],  
                [15,  9,  3,  7,  6]])
```

```

In [51]: # 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 latitude
# 2. cluster center longitude
# 3. day of the week
# 4. f_t_1: number of pickups that are happened previous t-1th 10min intravel
# 5. f_t_2: number of pickups that are happened previous t-2th 10min intravel
# 6. f_t_3: number of pickups that are happened previous t-3th 10min intravel
# 7. f_t_4: number of pickups that are happened previous t-4th 10min intravel
# 8. f_t_5: number of pickups that are happened previous t-5th 10min intravel

# from the baseline models we said the exponential weighted moving 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]
predict_list = []
tsne_flat_exp_avg = []
for r in range(0,40):
    for i in range(0,13104):
        if i==0:
            predicted_value= regions_cum[r][0]
            predicted_values.append(0)
            continue
        predicted_values.append(predicted_value)
        predicted_value =int((alpha*predicted_value) + (1-alpha)*(regions_cum[r][i]))
    predict_list.append(predicted_values[5:])
    predicted_values=[]

```

```
In [210]: len(predict_list[0])
```

```
Out[210]: 13099
```

```
In [211]: #predicted_value1=[]
#error=[]
predict_list1=[]
window_size=2
#predicted_value=[]
for r in range(0,40):
    lst = []
    for i in range(0,13104):
        predicted_value=0
        #error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
        if i+1>=window_size:
            #predicted_value.append(regions_cum[r][0])
            sum_values=0
            sum_of_coeff=0
            for j in range(window_size,0,-1):
                sum_values += j*(regions_cum[r][i-window_size+j])
                sum_of_coeff+=j
            predicted_value=int(sum_values/sum_of_coeff)
            #print("if:",predicted_value)
            #lst.append(predicted_values)

        else:
            #predicted_values.append(regions_cum[r][0])
            sum_values=0
            sum_of_coeff=0
            for j in range(i+1,0,-1):
                sum_values += j*(regions_cum[r][j-1])
                sum_of_coeff+=j
            predicted_value=int(sum_values/sum_of_coeff)
            #print("else:",predicted_value)
        lst.append(predicted_value)
    #print(lst)
    predict_list1.append(lst[5:])
#predicted_value = []
```

```
In [212]: len(predict_list1[0])
```

```
Out[212]: 13099
```

```
In [140]: regions_cum[10][1000]
```

```
Out[140]: 118
```

```
In [52]: len(predict_list[0])*len(predict_list) == 13099*40
```

```
Out[52]: True
```

```
In [53]: amplitude = []
frequency = []
for i in range(40):
    ampli = np.abs(np.fft.fft(regions_cum[i][:13099]))
    freq = np.abs(np.fft.fftfreq(13099,1))
    ampli_indices = np.argsort(-ampli)[1:]      #it will return an array of indices for which corresponding am
plitude values are sorted in reverse order.
    amplitude_values = []
    frequency_values = []
    for j in range(0, 9, 2): #taking top five amplitudes and frequencies
        amplitude_values.append(ampli[ampli_indices[j]])
        frequency_values.append(freq[ampli_indices[j]])
    for k in range(13099): #those top 5 frequencies and amplitudes are same for all the points in one cluster
        amplitude.append(amplitude_values)
        frequency.append(frequency_values)
```

```
In [54]: # 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*0.8))
print("size of test data :", int(13099*0.2))
```

```
size of train data : 10479
```

```
size of test data : 2619
```



```
In [55]: # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our training data
train_features = [tsne_feature[i*13099:(13099*i+10479)] for i in range(0,40)]
# temp = [0]*(12955 - 9068)
test_features = [tsne_feature[(13099*(i))+10479:13099*(i+1)] for i in range(0,40)]
```

```
In [56]: 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 10479 Each data point contains 5 features
 Number of data clusters 40 Number of data points in test data 2620 Each data point contains 5 features

```
In [57]: train_fourier_freq = [frequency[i*13099:(13099*i+10479)] for i in range(40)]
test_fourier_freq = [frequency[(i*13099)+10479:(13099*(i+1))] for i in range(40)]
```

```
In [58]: train_fourier_amp = [amplitude[i*13099:(13099*i+10479)] for i in range(40)]
test_fourier_amp = [amplitude[(i*13099)+10479:(13099*(i+1))] for i in range(40)]
```

```
In [59]: # extracting first 10479 timestamp values i.e 80% of 13099 (total timestamps) for our training data
tsne_train_flat_lat = [i[:10479] for i in tsne_lat]
tsne_train_flat_lon = [i[:10479] for i in tsne_lon]
tsne_train_flat_weekday = [i[:10479] for i in tsne_weekday]
tsne_train_flat_output = [i[:10479] for i in output]
tsne_train_flat_exp_avg = [i[:10479] for i in predict_list]
```

```
In [213]: tsne_train_flat_weg_avg = [i[:10479] for i in predict_list1]
```

```
In [60]: # extracting the rest of the timestamp values i.e 30% of 12956 (total timestamps) for our test data
tsne_test_flat_lat = [i[10479:] for i in tsne_lat]
tsne_test_flat_lon = [i[10479:] for i in tsne_lon]
tsne_test_flat_weekday = [i[10479:] for i in tsne_weekday]
tsne_test_flat_output = [i[10479:] for i in output]
tsne_test_flat_exp_avg = [i[10479:] for i in predict_list]
```

```
In [214]: tsne_test_flat_weg_avg = [i[10479:] for i in predict_list1]
```

```
In [61]: # convert from lists of lists of list to lists of list
train_pickups = []
test_pickups = []
train_freq = []
test_freq = []
train_amp = []
test_amp = []
for i in range(40):
    train_pickups.extend(train_features[i])
    test_pickups.extend(test_features[i])
    train_freq.extend(train_fourier_freq[i])
    test_freq.extend(test_fourier_freq[i])
    train_amp.extend(train_fourier_amp[i])
    test_amp.extend(test_fourier_amp[i])
```

```
In [62]: print(len(train_pickups))
print(len(test_pickups))
print(len(train_freq))
print(len(test_freq))
print(len(train_amp))
print(len(test_amp))
```

```
419160
104800
419160
104800
419160
104800
```

```
In [63]: train_features_freq_amp = np.hstack((train_pickups, train_freq, train_amp))
test_features_freq_amp = np.hstack((test_pickups, test_freq, test_amp))
```

```
In [64]: # 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_freq_amp[i])
         test_new_features = []
         for i in range(0,40):
             test_new_features.extend(test_features_freq_amp[i])
```

```
In [65]: # converting lists of lists into single 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 [215]: tsne_train_weg_avg = sum(tsne_train_flat_weg_avg,[])
```

```
In [66]: # converting lists of lists into single 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 [216]: tsne_test_weg_avg = sum(tsne_test_flat_weg_avg,[])
```

```
In [217]: # Preparing the data frame for our train data
import math
columns = ['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'freq1', 'freq2', 'freq3', 'freq4', 'freq5', 'Amp1', 'Amp2', 'Amp3',
'Amp4', 'Amp5']
df_train = pd.DataFrame(data=train_features_freq_amp, 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
df_train["weg_avg"] = tsne_train_weg_avg
#df_train['exp_avg_log'] = df_train['exp_avg'].apply(lambda num : Log(num))

print(df_train.shape)
```

(419160, 20)

```
In [218]: from math import log
df_train['exp_avg_log'] = df_train['exp_avg'].apply(lambda x : math.log(x+0.001))
```

```
In [220]: df_train.tail()
```

Out[220]:

	ft_5	ft_4	ft_3	ft_2	ft_1	freq1	freq2	freq3	freq4	freq5	...	Amp2	Amp3	
419155	111.0	119.0	112.0	133.0	127.0	0.000153	0.006871	0.007023	0.000076	0.013894	...	119614.989009	106952.567546	10
419156	119.0	112.0	133.0	127.0	150.0	0.000153	0.006871	0.007023	0.000076	0.013894	...	119614.989009	106952.567546	10
419157	112.0	133.0	127.0	150.0	118.0	0.000153	0.006871	0.007023	0.000076	0.013894	...	119614.989009	106952.567546	10
419158	133.0	127.0	150.0	118.0	114.0	0.000153	0.006871	0.007023	0.000076	0.013894	...	119614.989009	106952.567546	10
419159	127.0	150.0	118.0	114.0	103.0	0.000153	0.006871	0.007023	0.000076	0.013894	...	119614.989009	106952.567546	10

5 rows × 21 columns



```
In [221]: # Preparing the data frame for our train data
df_test = pd.DataFrame(data=test_features_freq_amp, 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
df_test["weg_avg"] = tsne_test_weg_avg
df_test['exp_avg_log'] = df_test['exp_avg'].apply(lambda x : math.log(x+0.001))
print(df_test.shape)
```

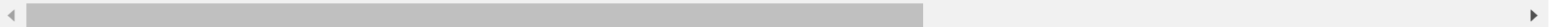
(104800, 21)

```
In [222]: df_test.head()
```

Out[222]:

	ft_5	ft_4	ft_3	ft_2	ft_1	freq1	freq2	freq3	freq4	freq5	...	Amp2	Amp3	Amp4
0	151.0	154.0	121.0	161.0	177.0	0.006947	0.013894	0.012902	0.034735	0.00794	...	180348.100917	82999.303583	63169.450
1	154.0	121.0	161.0	177.0	153.0	0.006947	0.013894	0.012902	0.034735	0.00794	...	180348.100917	82999.303583	63169.450
2	121.0	161.0	177.0	153.0	132.0	0.006947	0.013894	0.012902	0.034735	0.00794	...	180348.100917	82999.303583	63169.450
3	161.0	177.0	153.0	132.0	158.0	0.006947	0.013894	0.012902	0.034735	0.00794	...	180348.100917	82999.303583	63169.450
4	177.0	153.0	132.0	158.0	161.0	0.006947	0.013894	0.012902	0.034735	0.00794	...	180348.100917	82999.303583	63169.450

5 rows × 21 columns



```
In [223]: import pickle
with open("df_train.pkl","wb") as f:
    pickle.dump(df_train,f)
```

```
In [224]: with open("df_test.pkl","wb") as f:
    pickle.dump(df_test,f)
```

Using Linear Regression

```

In [245]: # find more about LinearRegression function here http://scikit-learn.org/stable/modules/generated/sklearn.Linear\_model.LinearRegression.html
# -----
# default paramters
# sklearn.linear_model.LinearRegression(fit_intercept=True, normalize=False, copy_X=True, n_jobs=1)

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

from sklearn.linear_model import SGDRegressor

from sklearn.model_selection import GridSearchCV
#hyper-paramater tuning
sgd = SGDRegressor(loss = "squared_loss", penalty = "l2")
values = [10**-4,10**-3,10**-2,10**-1,10**0,10**1,10**2]
param = {"alpha": values}
clf = GridSearchCV(sgd, param, scoring = "neg_mean_absolute_error", cv = 5,n_jobs=-1)
clf.fit(df_train, tsne_train_output)
alpha = clf.best_params_["alpha"]

#applying linear regression with best hyper-parameter
clf = SGDRegressor(loss = "squared_loss", penalty = "l2", alpha = alpha)
clf.fit(df_train, tsne_train_output)
y_pred = clf.predict(df_test)
lr_test_predictions = [round(value) for value in y_pred]
y_pred = clf.predict(df_train)
lr_train_predictions = [round(value) for value in y_pred]

```

Using Random Forest Regressor


```

In [246]: # Training a hyper-parameter tuned random forest regressor on our train data
# find more about LinearRegression function here http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestRegressor.html
# -----
# default paramters
# sklearn.ensemble.RandomForestRegressor(n_estimators=10, criterion='mse', max_depth=None, min_samples_split=2,
# min_samples_leaf=1, min_weight_fraction_leaf=0.0, max_features='auto', max_leaf_nodes=None, min_impurity_decrease=0.0,
# min_impurity_split=None, bootstrap=True, oob_score=False, n_jobs=1, random_state=None, verbose=0, warm_start=False)

# some of methods of RandomForestRegressor()
# apply(X)      Apply trees in the forest to X, return leaf indices.
# decision_path(X)      Return the decision path in the forest
# fit(X, y[, sample_weight])      Build a forest of trees from the training set (X, y).
# get_params([deep])      Get parameters for this estimator.
# predict(X)      Predict regression target for X.
# score(X, y[, sample_weight])      Returns the coefficient of determination R^2 of the prediction.
# -----
# video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/regression-using-decision-trees-2/
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/lessons/what-are-ensembles/
# -----
from sklearn.model_selection import RandomizedSearchCV
#values = [10, 40, 100, 300, 500, 1000]
rf = RandomForestRegressor(n_jobs=-1)
params = {"n_estimators": [10, 100, 300, 500],
          "min_samples_split": [2, 3, 5, 6],
          "min_samples_leaf": [2, 4, 6, 8],
          "max_depth": [3, 5, 9, 11],
          "max_features": ["auto", "sqrt", "log2"]}
clf = RandomizedSearchCV(rf, params, scoring = "neg_mean_absolute_error", cv = 5, n_jobs=-1)
clf.fit(df_train, tsne_train_output)
print(clf.best_estimator_)
#estimators = best_parameter.best_params_["n_estimators"]

#applying random forest with best hyper-parameter
#clf = RandomForestRegressor(n_estimators = estimators)
#clf.fit(df_train, tsne_train_output)

```

```
RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=11,
                        max_features='auto', max_leaf_nodes=None,
                        min_impurity_decrease=0.0, min_impurity_split=None,
                        min_samples_leaf=4, min_samples_split=6,
                        min_weight_fraction_leaf=0.0, n_estimators=300, n_jobs=-1,
                        oob_score=False, random_state=None, verbose=0, warm_start=False)
```

In [247]: *# 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
clf =RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=11,
                            max_features='auto', max_leaf_nodes=None,
                            min_impurity_decrease=0.0, min_impurity_split=None,
                            min_samples_leaf=4, min_samples_split=6,
                            min_weight_fraction_leaf=0.0, n_estimators=300, n_jobs=-1,
                            oob_score=False, random_state=None, verbose=0, warm_start=False)
clf.fit(df_train, tsne_train_output)
y_pred = clf.predict(df_test)
rndf_test_predictions = [round(value) for value in y_pred]
y_pred = clf.predict(df_train)
rndf_train_predictions = [round(value) for value in y_pred]
```

In [248]: *#feature importances based on analysis using random forest*

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

```
Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'freq1', 'freq2', 'freq3',
      'freq4', 'freq5', 'Amp1', 'Amp2', 'Amp3', 'Amp4', 'Amp5', 'lat', 'lon',
      'weekday', 'exp_avg', 'weg_avg', 'exp_avg_log'],
      dtype='object')
[7.58855194e-07 8.73654689e-07 8.59624795e-07 3.32226253e-06
 5.55295339e-03 1.54096063e-07 1.36401418e-07 1.22231706e-07
 1.68983090e-07 1.41828427e-07 1.58862727e-07 1.18867762e-07
 9.31812847e-08 1.31723760e-07 1.46364389e-07 2.82217671e-07
 2.93448110e-07 2.09494950e-07 8.93954381e-05 9.94264225e-01
 8.54538056e-05]
```

Using XgBoost Regressor

```
In [249]: # Training a hyper-parameter tuned Xg-Boost regressor on our train data

# find more about XGBRegressor function here http://xgboost.readthedocs.io/en/latest/python/python\_api.html?#module-xgboost.sklearn
# -----
# default paramters
# xgboost.XGBRegressor(max_depth=3, learning_rate=0.1, n_estimators=100, silent=True, objective='reg:Linear',
# booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1, max_delta_step=0, 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/
# -----

params = {"learning_rate" : [0.0001,0.001,0.01,0.1,1],
          "n_estimators":[10, 40, 100, 500, 1000],
          "max_depth" : [3,5,7,9,11]}

x_model = xgb.XGBRegressor(n_jobs=-1)
clf = RandomizedSearchCV(x_model, params, scoring="neg_mean_absolute_error",n_jobs=-1,cv=5)
clf.fit(df_train, tsne_train_output)
clf.best_estimator_
```

```
Out[249]: XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                        colsample_bytree=1, gamma=0, importance_type='gain',
                        learning_rate=0.1, max_delta_step=0, max_depth=11,
                        min_child_weight=1, missing=None, n_estimators=100, n_jobs=-1,
                        nthread=None, objective='reg:linear', random_state=0, reg_alpha=0,
                        reg_lambda=1, scale_pos_weight=1, seed=None, silent=True,
                        subsample=1)
```

```
In [250]: #predicting with our trained Xg-Boost regressor
# the models x_model is already hyper parameter tuned
# the parameters that we got above are found using grid search
x_model = xgb.XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                           colsample_bytree=1, gamma=0, importance_type='gain',
                           learning_rate=0.1, max_delta_step=0, max_depth=11,
                           min_child_weight=1, missing=None, n_estimators=100, n_jobs=-1,
                           nthread=None, objective='reg:linear', random_state=0, reg_alpha=0,
                           reg_lambda=1, scale_pos_weight=1, seed=None, silent=True,
                           subsample=1)
x_model.fit(df_train,tsne_train_output)
y_pred = x_model.predict(df_test)
xgb_test_predictions = [round(value) for value in y_pred]
y_pred = x_model.predict(df_train)
xgb_train_predictions = [round(value) for value in y_pred]
```

Calculating the error metric values for various models

```
In [251]: train_mape=[]
          test_mape=[]

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

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

Error Metric Matrix

```
In [252]: 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 regressor - Train: ",train_mape[3]," Test: ",test_mape[3])
```

Error Metric Matrix (Tree Based Regression Methods) - MAPE

```
-----
Baseline Model - Train: 0.14869676078402586 Test: 0.1389957811283114
Exponential Averages Forecasting - Train: 0.1412521852007101 Test: 0.1315394376406458
Linear Regression - Train: 7.435871286648411e+16 Test: 6.9046310066689816e+16
Random Forest Regression - Train: 0.006352796187657793 Test: 0.006357003436047449
Xgboost regressor - Train: 0.0037375479778321645 Test: 0.004907429885864914
```

XGBoost is as far as the best model from the above cell with 0.4% MAPE.

SUMMARY:

Step-by-Step procedure used in sloving this problem:

1.Loading the data

- In this case study we used 2015 jan data which is it self is huge so we use dask to overcome memryo constraint. Dask is similar to pandas but it has a nice ferature which adds key value pairs and do operations.
- pose the problem type in this case it is a regression and time series problem.

2. Data cleaning:

- The most important features are pickup latitude and longitude, dropoff latitude and longitude, pickup and dropoff datetime, trip distance, total amount.

A. Pickup Latitude and Longitude:

- From the distribution plots we can see that the range of latitudes and longitude are very high but in reality between -74.15 to 40.5774 latitude and -73.7004 to 40.9176 longitude.
- so, apply this values to get only pickups happened in the city only.

B. Dropoff Latitude and Longitude:

- Apply the same above conditions and to get the drops happened only int he city.

C. Trip Durations:

- According to the newyork city rules a cab can be drived only 12hrs a day. So, consider the trip time happened at max 12hrs.
- But, before that month date time should be converted into unix time.
- And added Speed feature to a dataframe.

D. Speed:

- By calculatings the percentile it is founded that the avg speed ofthe cab in the new york city is 12.45miles/hr i.e., on avg 2 miles per 10 mins.

E. Trip Distanace:

- By seeing 99.9 percntile 22.57 miles is best best trip distance.

F. Total Fare:

- By calculating percentiles and plotings graphs 1000 can be the max value by keeping trip duration in mind.

3. Data Preperation:

- Preparing the data...

A. Clustering:

- find clusters using k-means clusterings and divide them into regions.

B. Time Binning:

- Using pickup timing calculate pickup bins.
- finally, apply the above all data cleaning methods to 2016-jan,feb,mar this is the final data we work on.
- smooth the data.

4. Modeling

- Using some of the popular baseline models like simple moving averages, weighted moving averages, exponential weightedmoving averages

5. Regression Models:

- split the data into 80:20
- consider pickups happened at the 5 10mins interval as a features.
- And top 5 freq and amplitudes in each region.
- finally, combine all features and apply linear regression. RandomForestRegressor and XGBoostRegressor.