

## ¬ Taxi demand prediction in New York City



```
#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
import dask.dataframe as dd#similar to pandas

import pandas as pd#pandas to create small dataframes

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

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

import time #Convert to unix time

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

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

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

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

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

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

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

## ¬ Data Information

Get the data from : [http://www.nyc.gov/html/tlc/html/about/trip\\_record\\_data.shtml](http://www.nyc.gov/html/tlc/html/about/trip_record_data.shtml) (2016 data) The data used in 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 medallions issued by the TLC. You access this mode of transportation by standing in the street and hailing an av arranged.

### **For Hire Vehicles (FHV)**

FHV transportation is accessed by a pre-arrangement with a dispatcher or limo company. These FHVs are not pe 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 bran 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

## ▼ 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 f
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

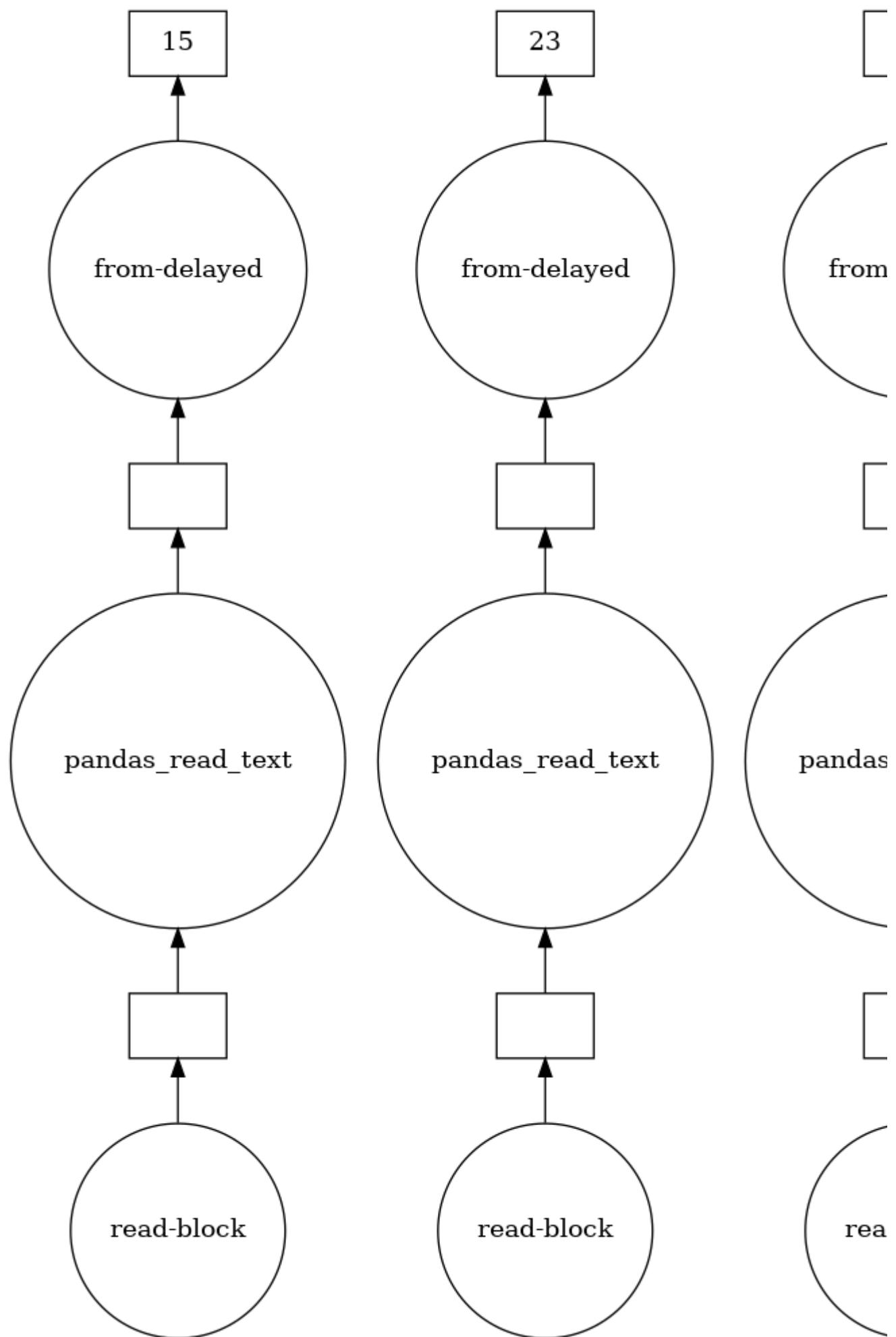
```
#Looking at the features
# dask_dataframe : # https://github.com/dask/dask-tutorial/blob/master/07_dataframe.ipynb
month = dd.read_csv('yellow_tripdata_2015-01.csv')
print(month.columns)
```

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')

```
# 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.
```

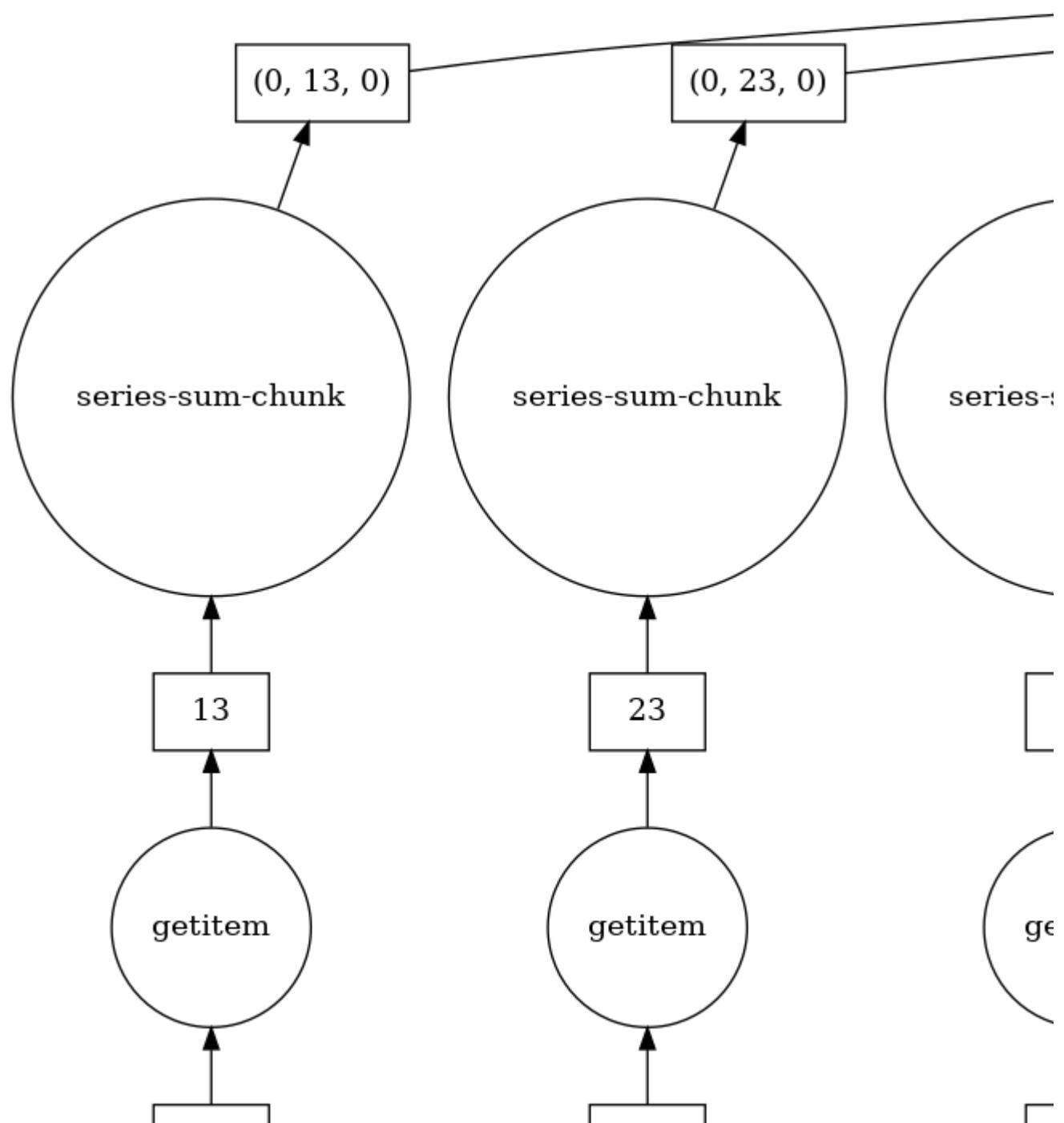
```
# to see the visualization you need to install graphviz
# pip3 install graphviz if this doesnt work please check the install_graphviz.jpg in the drive
month.visualize()
```

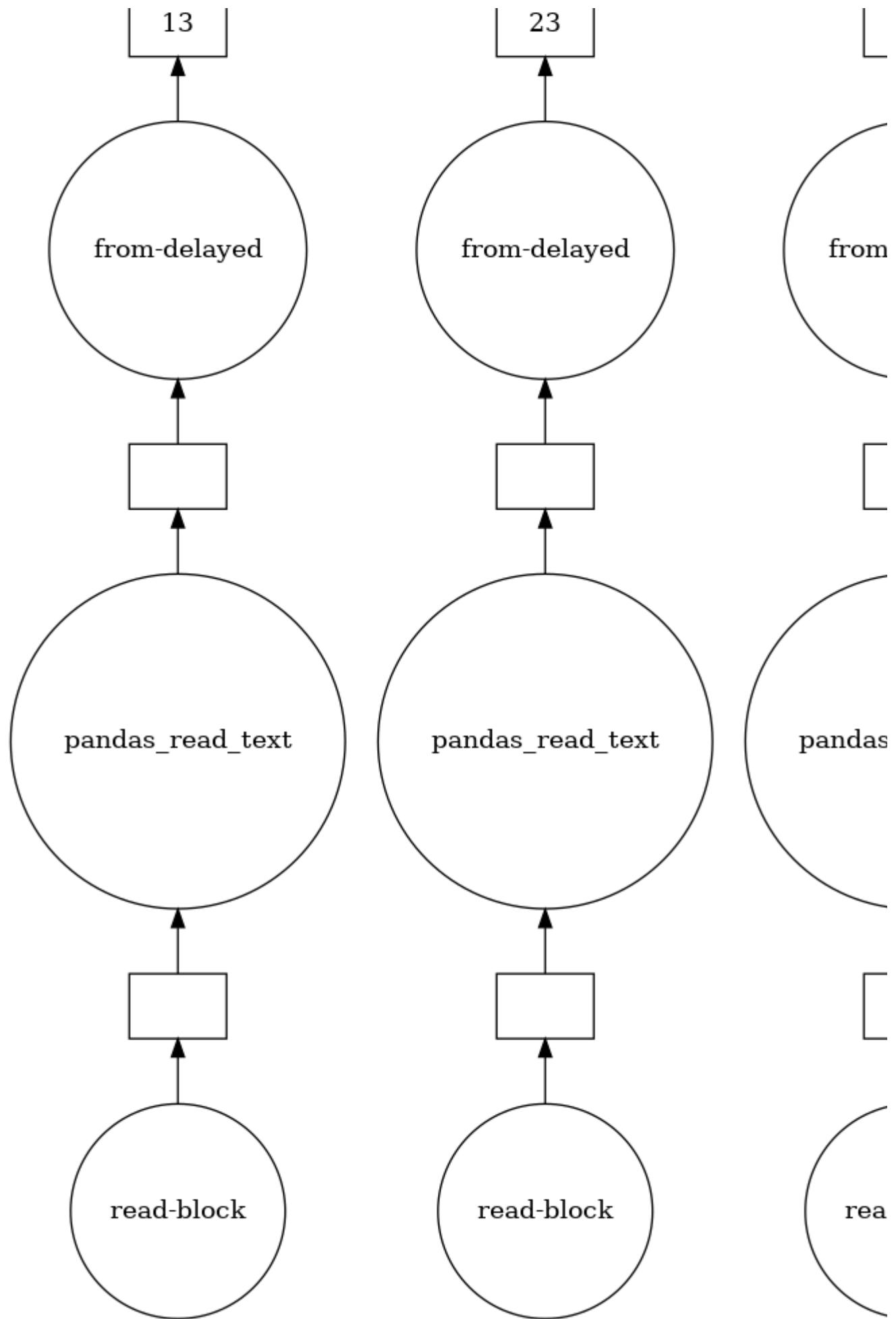




```
month.fare_amount.sum().visualize()
```







# 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 subsequently predict the same.
- 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 causing the problem.

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

	VendorID	tpep_pickup_datetime	tpep_dropoff_datetime	passenger_count	trip_distance
0	2	2015-01-15 19:05:39	2015-01-15 19:23:42	1	
1	1	2015-01-10 20:33:38	2015-01-10 20:53:28	1	
2	1	2015-01-10 20:33:38	2015-01-10 20:43:41	1	
3	1	2015-01-10 20:33:39	2015-01-10 20:35:31	1	
4	1	2015-01-10 20:33:39	2015-01-10 20:52:58	1	

## ▼ 1. Pickup Latitude and Pickup Longitude

It is inferred from the source <https://www.flickr.com/places/info/2459115> that New York is bounded by the latitude (40.9176,-73.7004) so hence any coordinates not within these coordinates are not considered by us as we are only interested in New York.

```
# 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

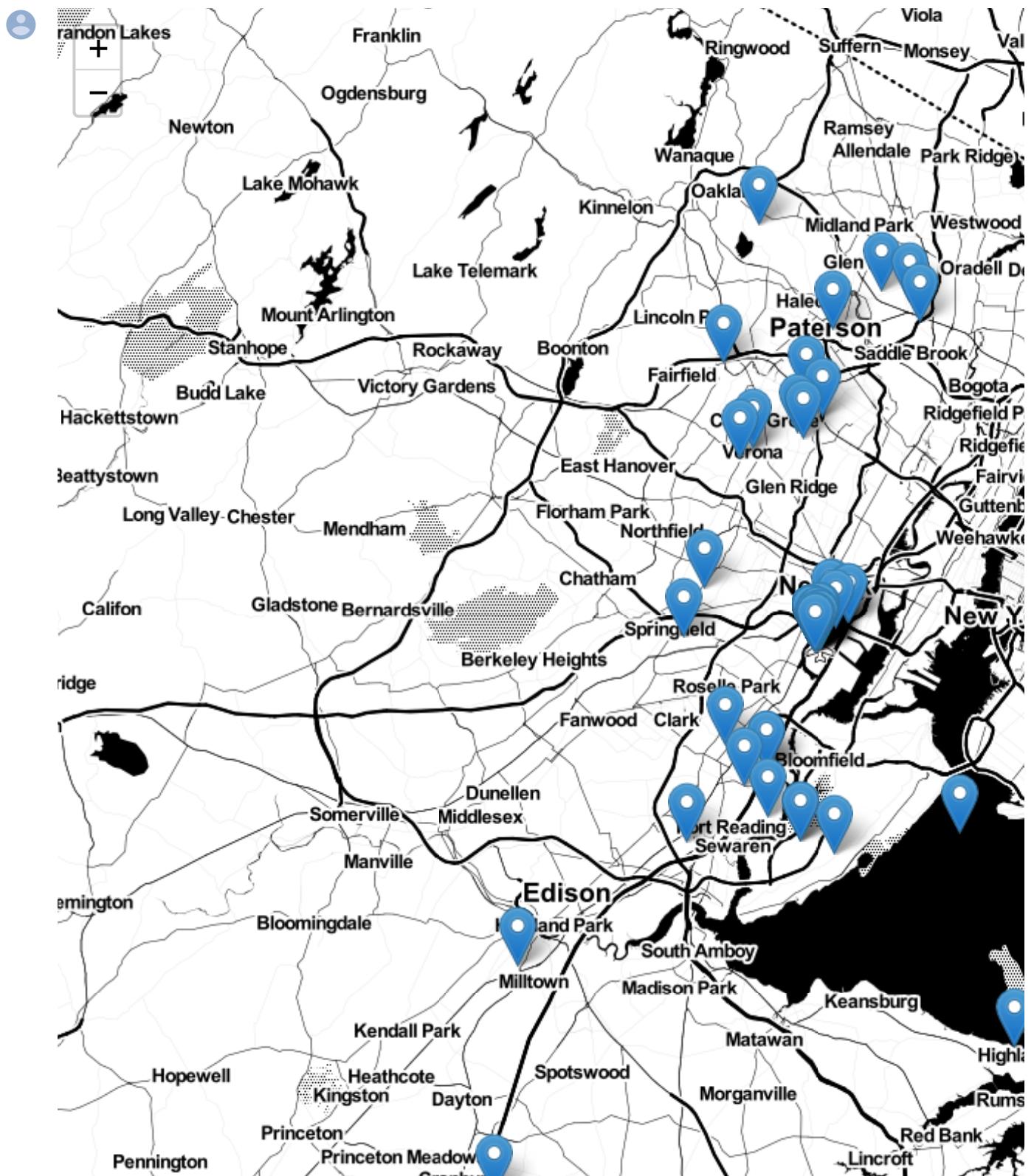
# note: you dont need to remember any of these, you dont need indepth knowledge on these maps as they are just for reference
map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')

# we will spot only first 100 outliers on the map, plotting all the outliers will take more time
sample_locations = outlier_locations.head(10000)
```

```

for i,j in sample_locations.itervalues():
    if int(j['pickup_latitude']) != 0:
        folium.Marker(list((j['pickup_latitude'],j['pickup_longitude']))).add_to(map_osm)
map_osm

```



**Observation:-** As you can see above that there are some points just outside the boundary but there are a few tha

## ▼ 2. Dropoff Latitude & Dropoff Longitude

It is inferred from the source <https://www.flickr.com/places/info/2459115> that New York is bounded by the loca (40.9176,-73.7004) so hence any cordinates not within these cordinates are not considered by us as we are only

York.

```
# 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 <= 40.57)
                             (month.dropoff_longitude >= -73.7004) | (month.dropoff_latitude >= 40.9176))]

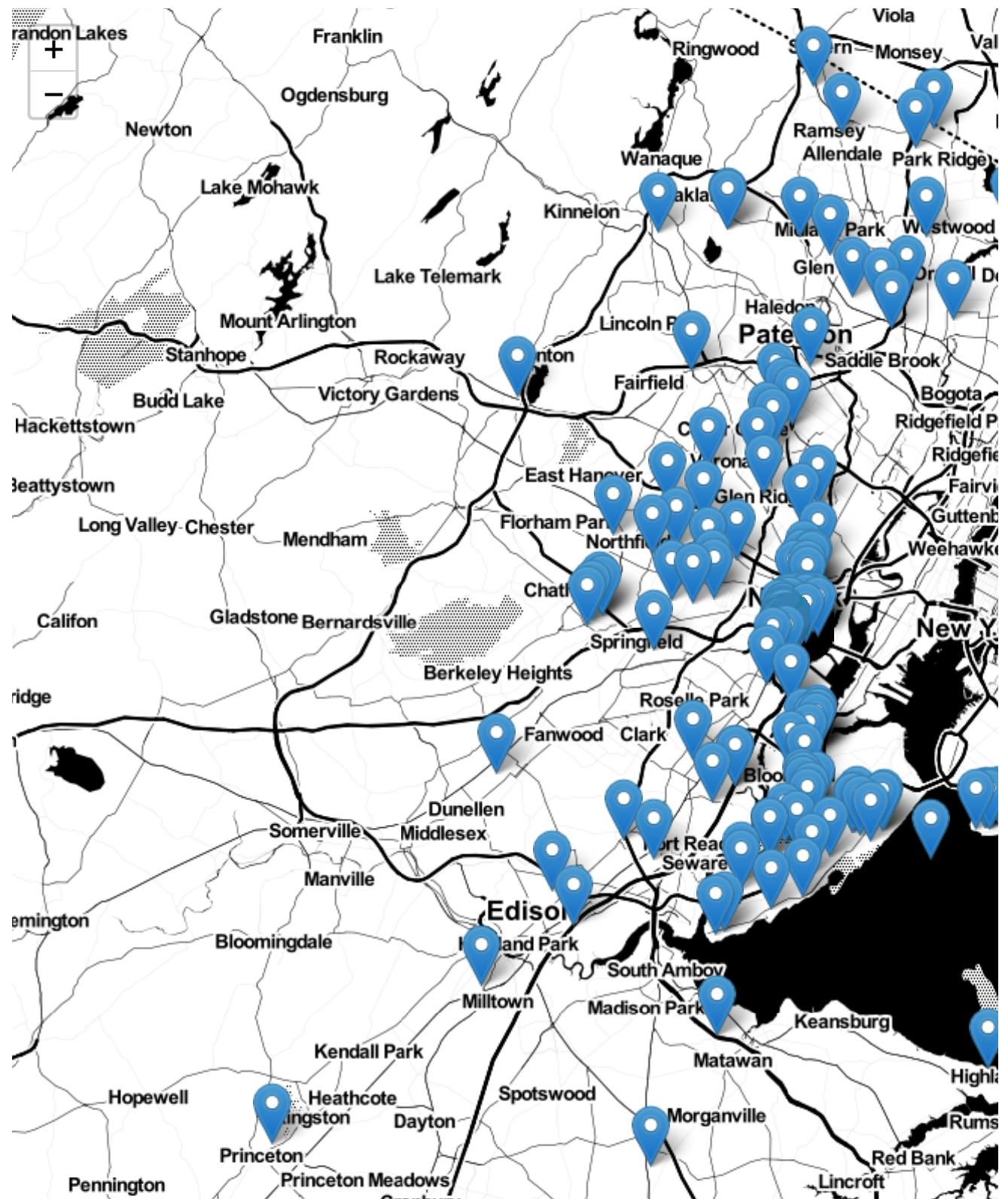
# creating a map with the a base location
# read more about the folium here: http://folium.readthedocs.io/en/latest/quickstart.html

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

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

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





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

### ▼ 3. Trip Durations:

According to NYC Taxi & Limousine Commission Regulations the maximum allowed trip duration in a 24 hour int

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

```
# in out data we have time in the formate "YYYY-MM-DD HH:MM:SS" we convert this sting to python
# 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','dr
    new_frame['trip_times'] = durations
    new_frame['pickup_times'] = duration_pickup
    new_frame['Speed'] = 60*(new_frame['trip_distance']/new_frame['trip_times']))

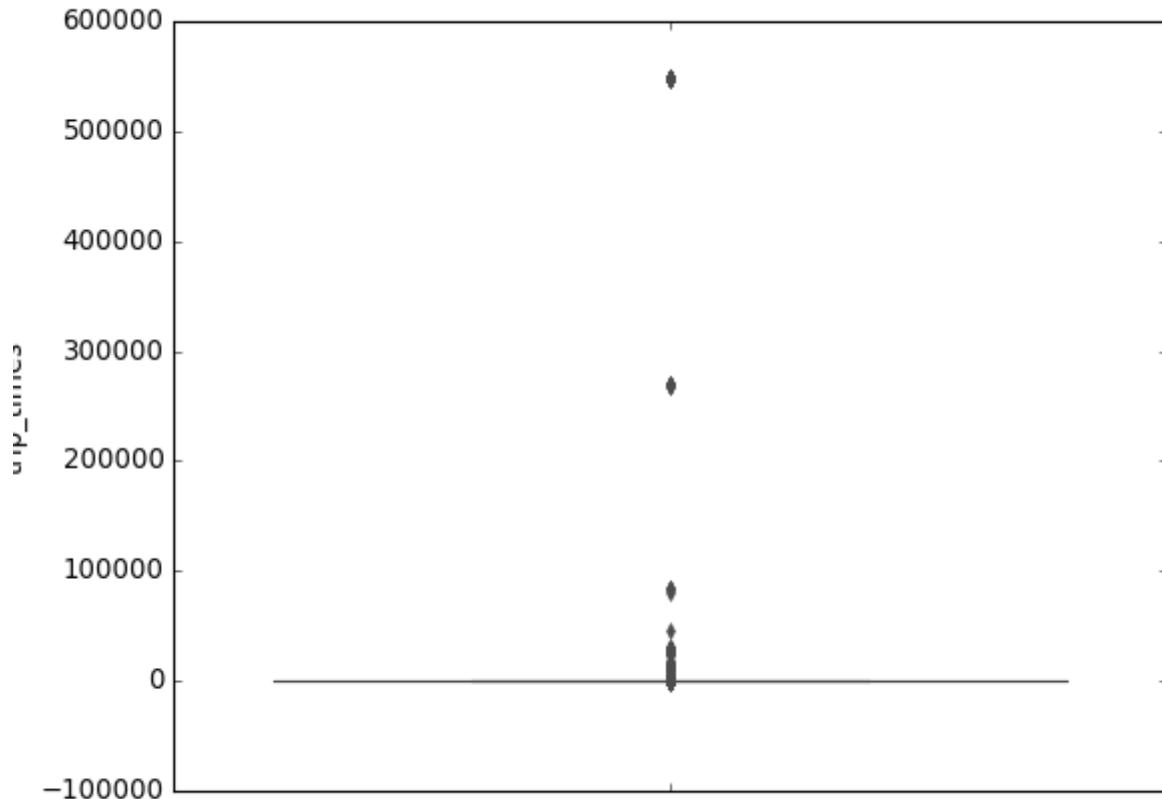
    return new_frame

# print(frame_with_durations.head())
#   passenger_count  trip_distance  pickup_longitude  pickup_latitude  dropoff_longitude  dropoff_latitude
# 1           1.59      -73.993896        40.750111       -73.974785        40.750618
# 1           3.30      -74.001648        40.724243       -73.994415        40.759109
# 1           1.80      -73.963341        40.802788       -73.951820        40.824413
# 1           0.50      -74.009087        40.713818       -74.004326        40.719986
# 1           3.00      -73.971176        40.762428       -74.004181        40.742653
frame_with_durations = return_with_trip_times(month)

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

```





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

 0 percentile value is -1211.0166666666667  
 10 percentile value is 3.8333333333333335  
 20 percentile value is 5.383333333333334  
 30 percentile value is 6.816666666666666  
 40 percentile value is 8.3  
 50 percentile value is 9.95  
 60 percentile value is 11.866666666666667  
 70 percentile value is 14.28333333333333  
 80 percentile value is 17.63333333333333  
 90 percentile value is 23.45  
 100 percentile value is 548555.6333333333

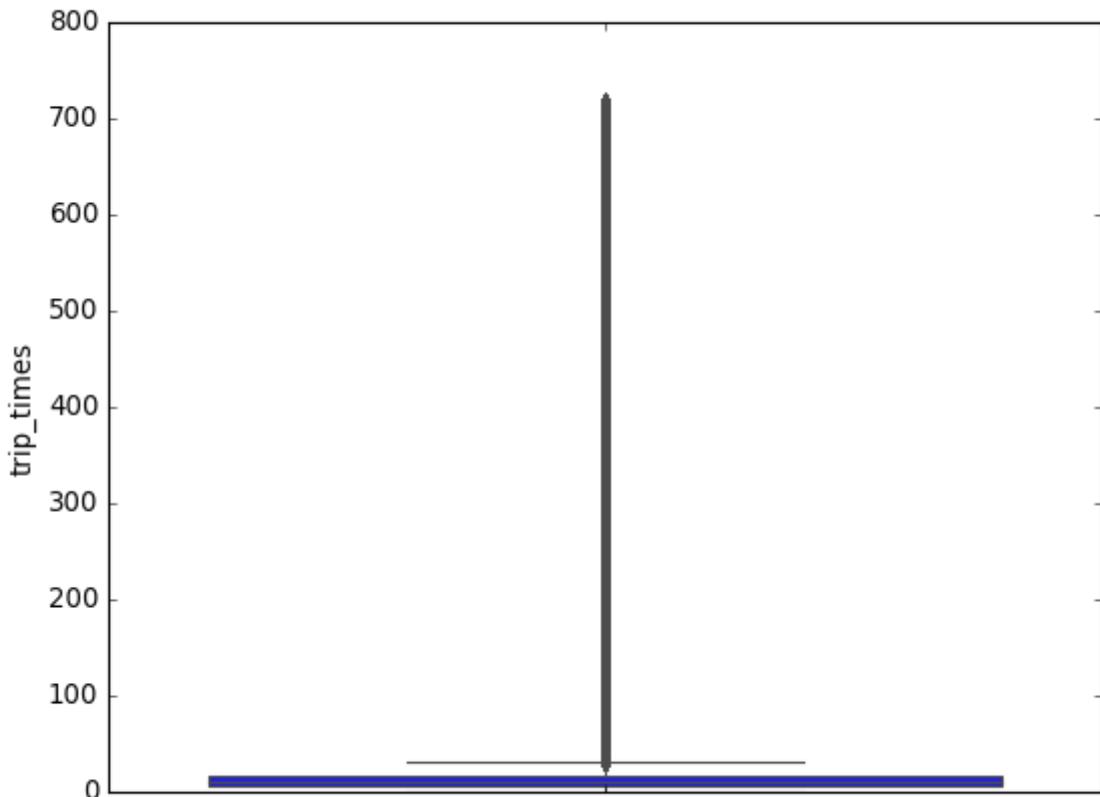
```
#looking further from the 99th percecntile
for i in range(90,100):
    var = frame_with_durations["trip_times"].values
    var = np.sort(var, axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100))]))
print ("100 percentile value is ",var[-1])
```



```
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 27.93333333333334
95 percentile value is 29.58333333333332
96 percentile value is 31.68333333333334
97 percentile value is 34.46666666666667
98 percentile value is 38.71666666666667
99 percentile value is 46.75
100 percentile value is 548555.6333333333
```

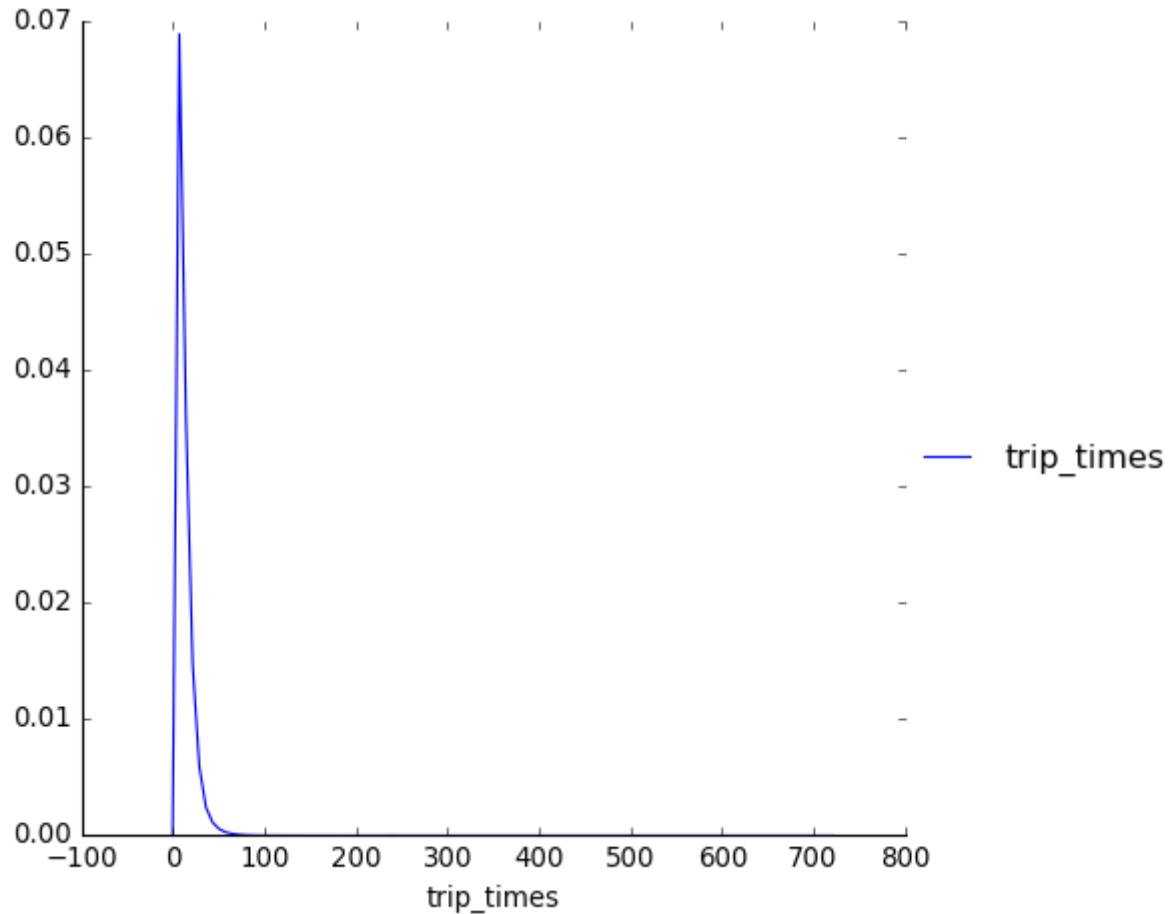
```
#removing data based on our analysis and TLC regulations
frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip_times>1) & (frame_w
```

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



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

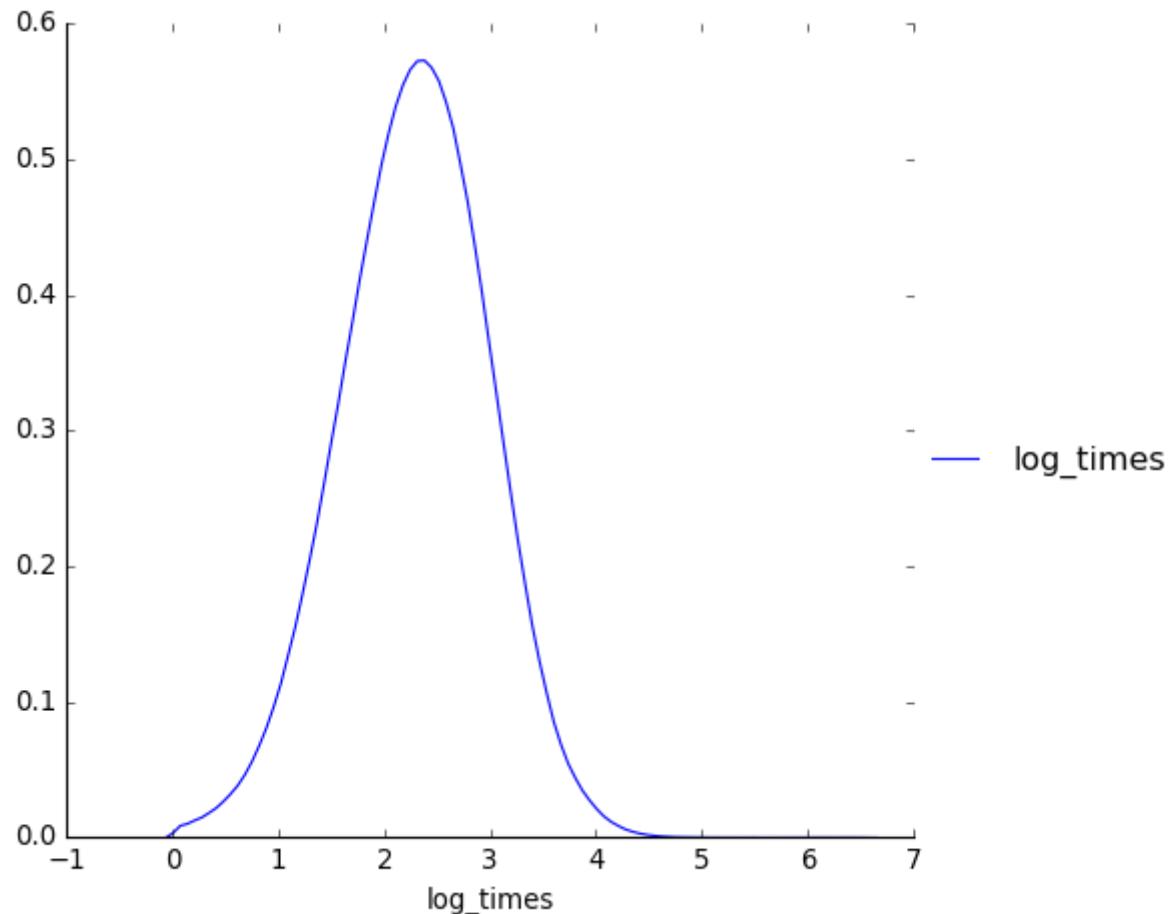




```
#converting the values to log-values to chec for log-normal
import math
frame_with_durations_modified['log_times']=[math.log(i) for i in frame_with_durations_modified['t']]

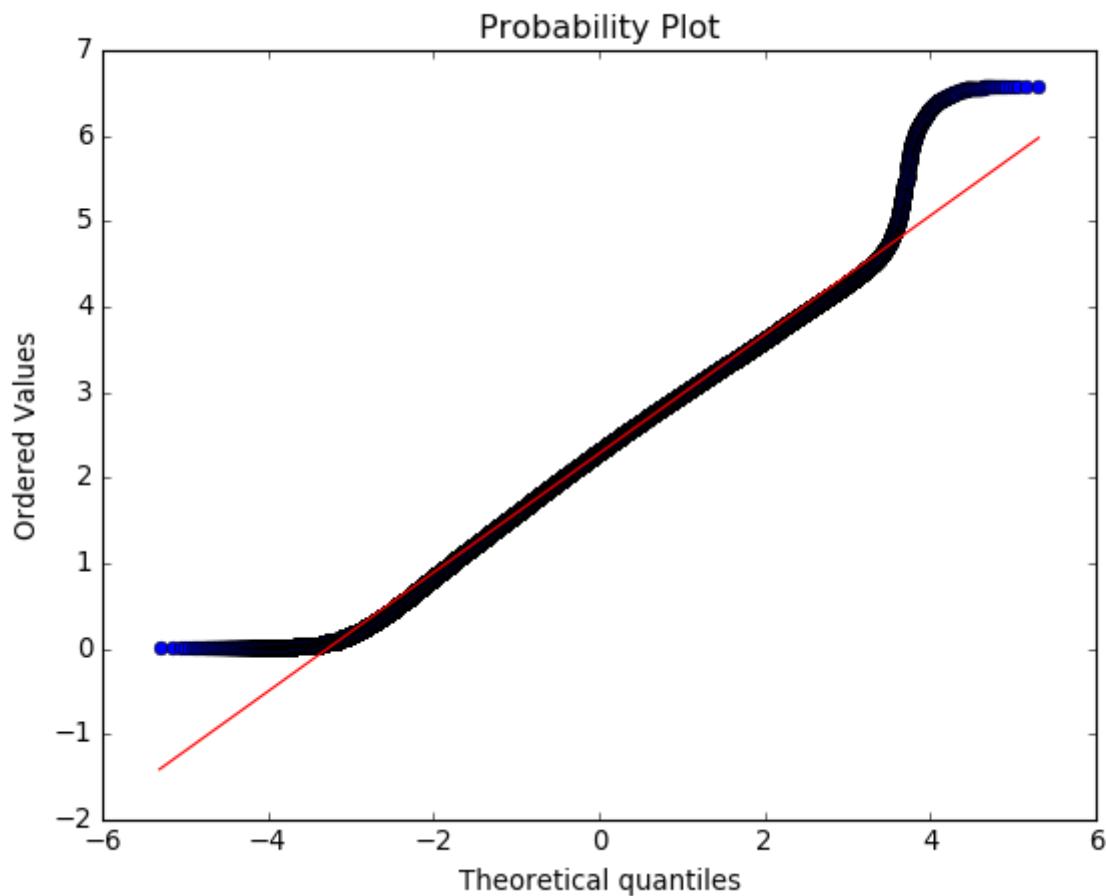
#pdf of log-values
sns.FacetGrid(frame_with_durations_modified,size=6) \
    .map(sns.kdeplot,"log_times") \
    .add_legend();
plt.show();
```





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

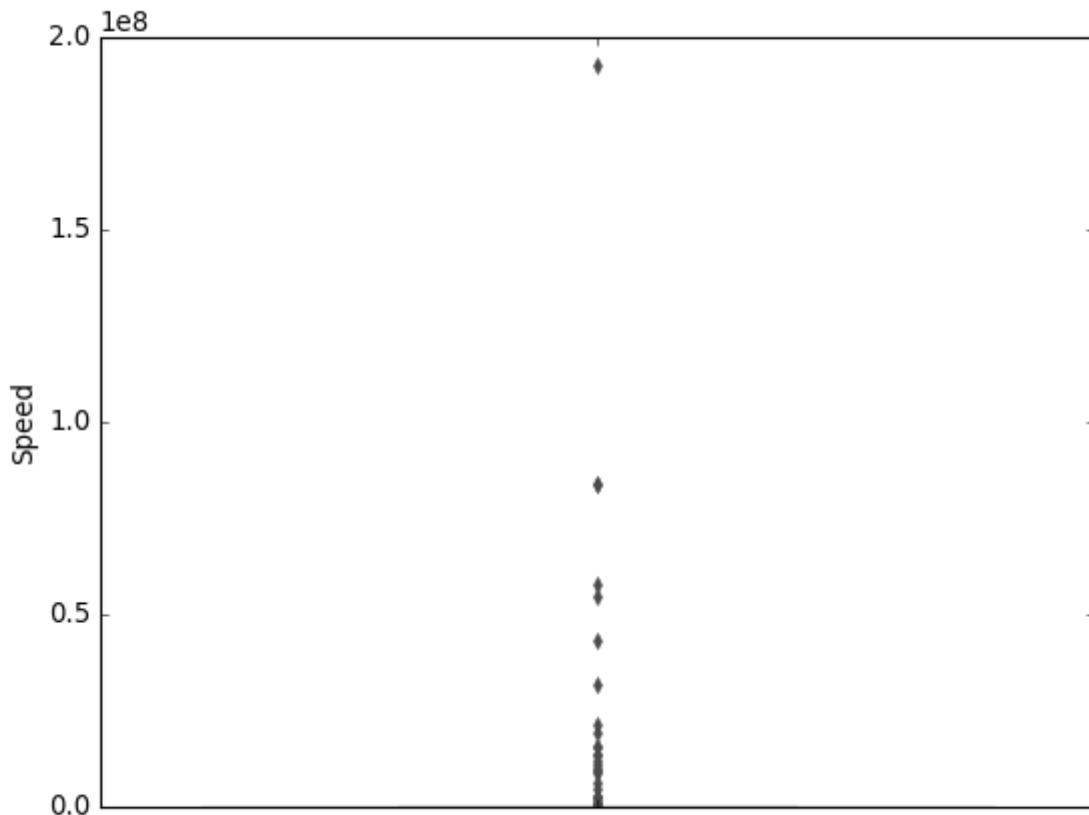




## ▼ 4. Speed

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





```
#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

```
#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
```

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

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

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

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

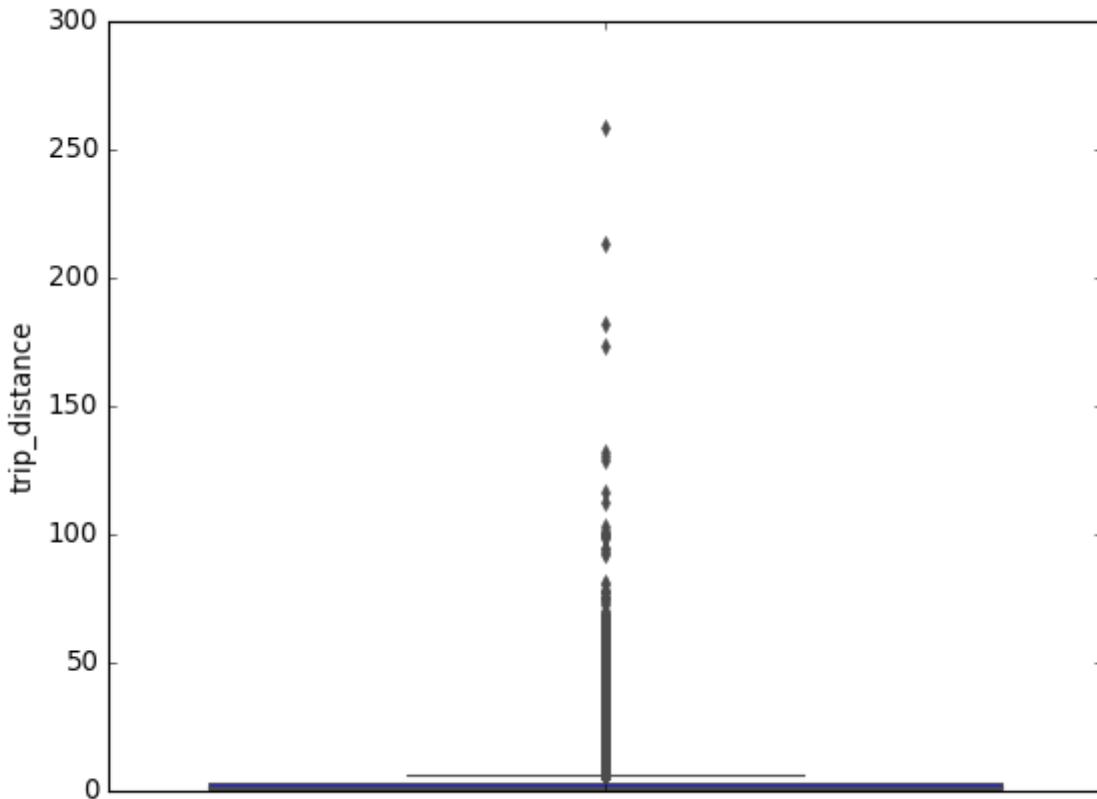
 12.450173996027528

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

```
# 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()
```





```
#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

```
#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
```

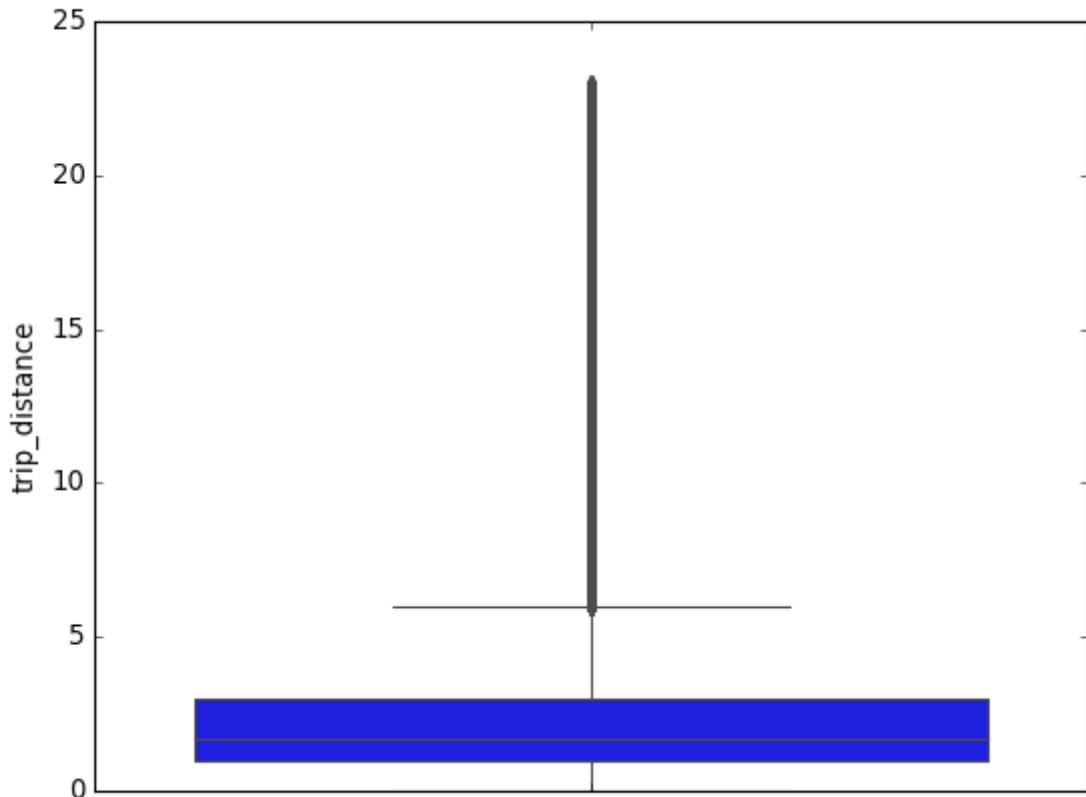
```
#calculating trip distance values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,
for i in np.arange(0.0, 1.0, 0.1):
    var = frame_with_durations_modified["trip_distance"].values
    var = np.sort(var, axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))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

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

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

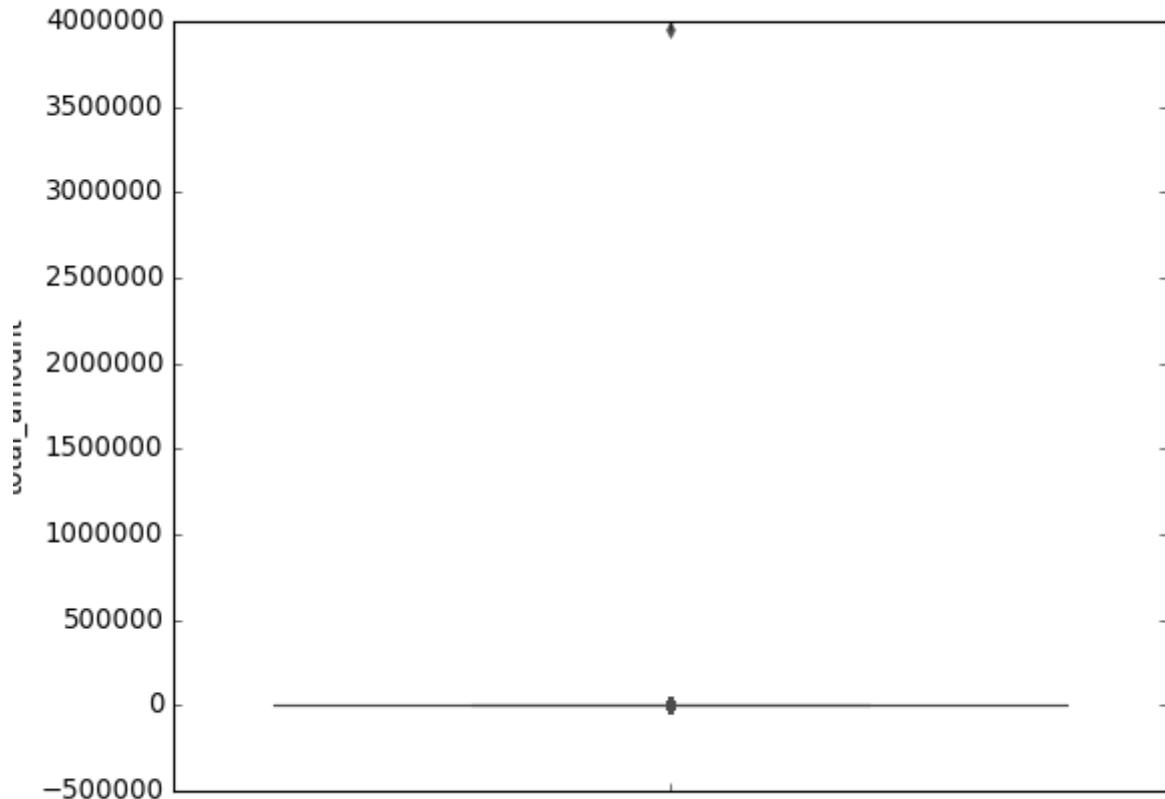




## ▼ 5. Total Fare

```
# 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()
```





```
#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

```
#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
```

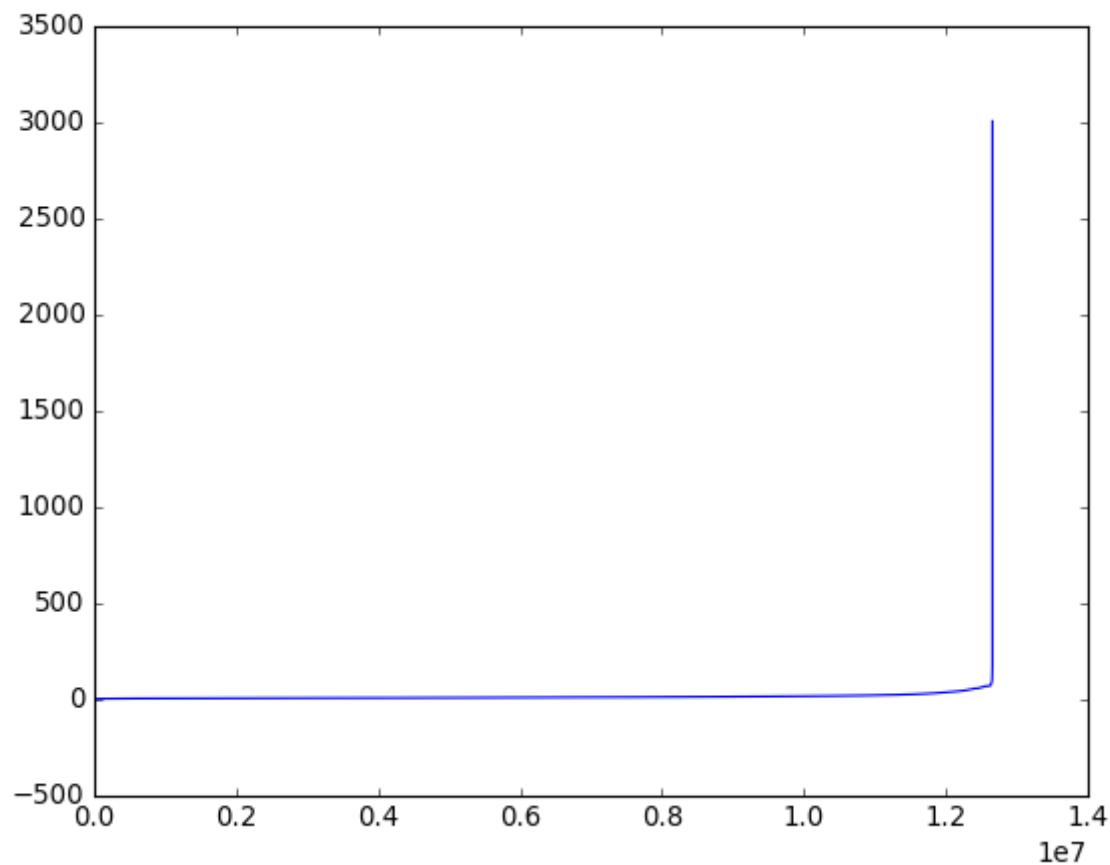
```
#calculating total fare amount values at each percentile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9
for i in np.arange(0.0, 1.0, 0.1):
    var = frame_with_durations_modified["total_amount"].values
    var = np.sort(var, axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100))]))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

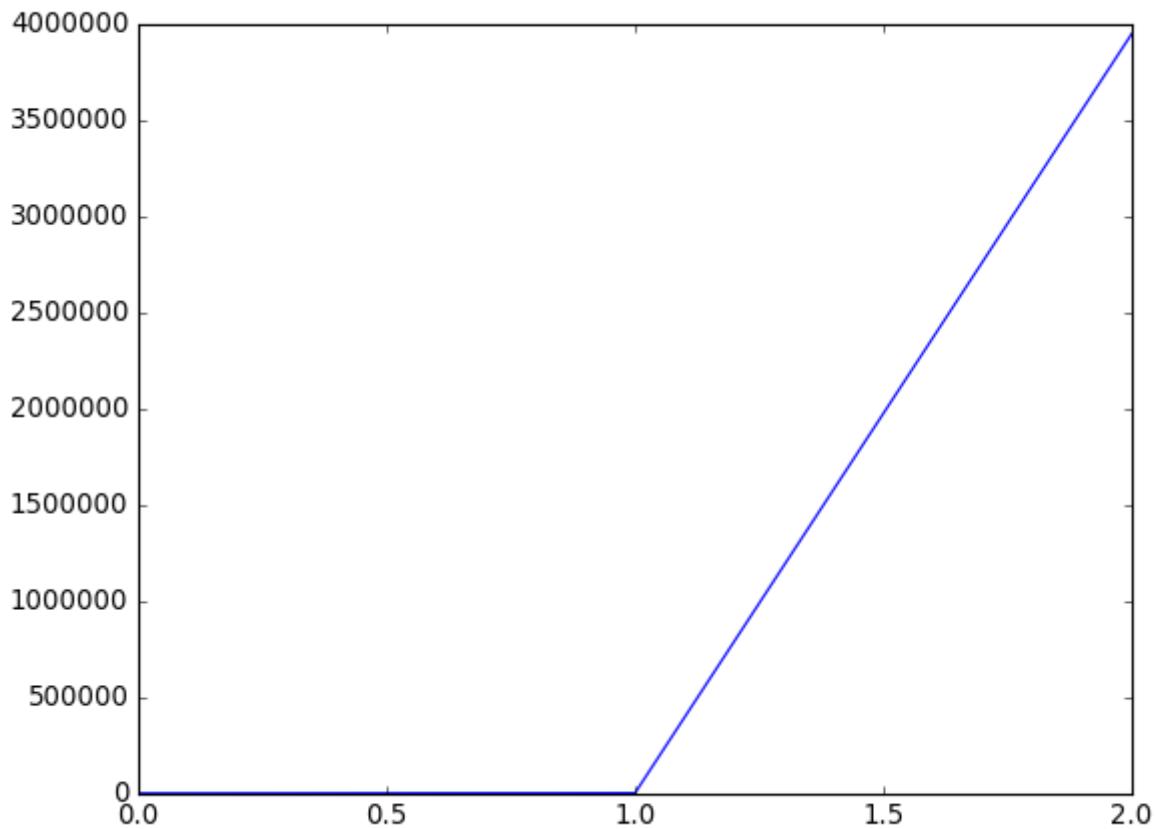
```
#below plot shows us the fare values(sorted) to find a sharp increase to remove those values as outliers
# plot the fare amount excluding last two values in sorted data
plt.plot(var[:-2])
plt.show()
```





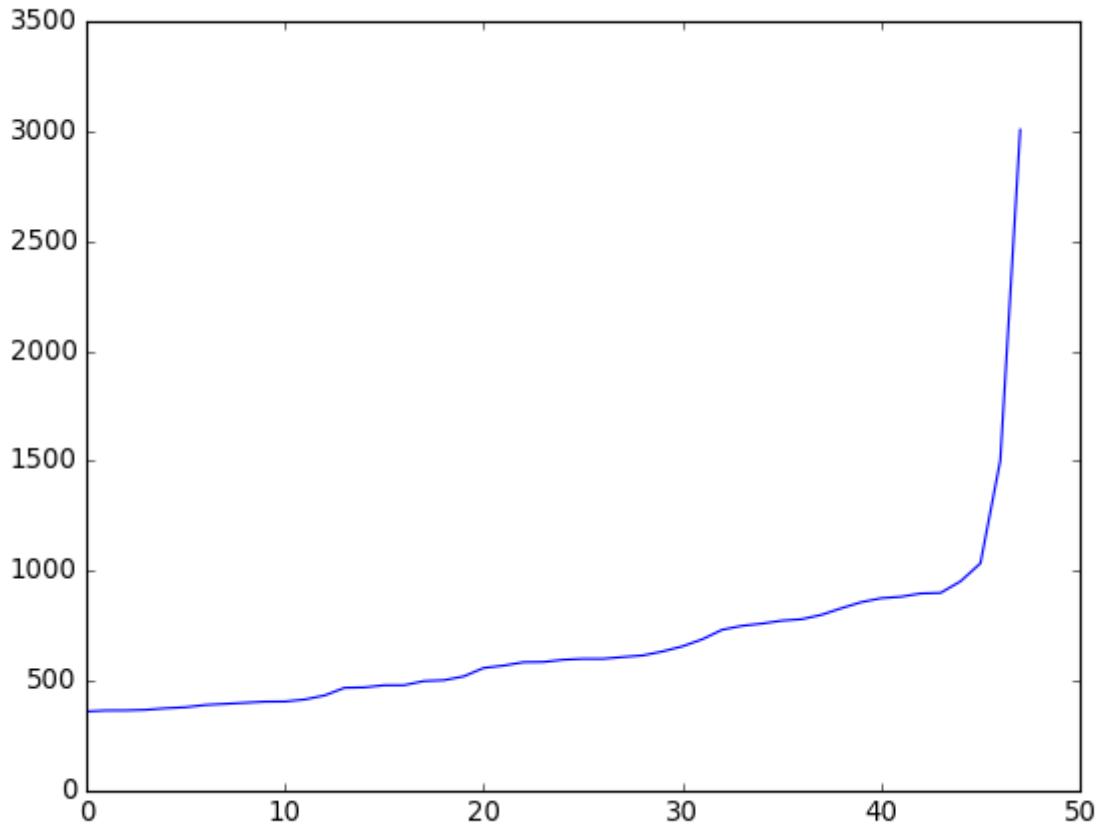
```
# a very sharp increase in fare values can be seen  
# plotting last three total fare values, and we can observe there is share increase in the values  
plt.plot(var[-3:])  
plt.show()
```





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





## ▼ Remove all outliers/errorous points.

```
#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 <= 40.5774) & (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitude <= 41.5) & ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_latitude >= 40.5774) & (new_frame.pickup_longitude <= -73.7004) & (new_frame.pickup_latitude <= 41.5))]
    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 <= -74.05) & (new_frame.dropoff_latitude >= 40.5774) & (new_frame.dropoff_latitude <= 40.6074) & ((new_frame.pickup_longitude >= -74.15) & (new_frame.pickup_longitude <= -74.05) & (new_frame.pickup_latitude >= 40.5774) & (new_frame.pickup_latitude <= 40.6074))

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

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

```

#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 = gpixpy.geo.haversine_distance(cluster_centers[i][0], cluster_centers[i][1], cluster_centers[j][0], cluster_centers[j][1])
                min_dist = min(min_dist,distance/(1.60934*1000))
                if (distance/(1.60934*1000)) <= 2:
                    nice_points +=1
                else:
                    wrong_points += 1
            less2.append(nice_points)
            more2.append(wrong_points)
        neighbours.append(less2)
    print ("On choosing a cluster size of ",cluster_len,"Avg. Number of Clusters within the vicinity",min_dist)

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

```

```

cluster_centers = kmeans.cluster_centers_
cluster_len = len(cluster_centers)
return cluster_centers, cluster_len

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

```



On choosing a cluster size of 10

Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 2.0  
 Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 8.0  
 Min inter-cluster distance = 1.0945442325142543

---

On choosing a cluster size of 20

Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 4.0  
 Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 16.0  
 Min inter-cluster distance = 0.7131298007387813

---

On choosing a cluster size of 30

Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0  
 Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 22.0  
 Min inter-cluster distance = 0.5185088176172206

---

On choosing a cluster size of 40

Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 8.0  
 Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2): 32.0  
 Min inter-cluster distance = 0.5069768450363973

---

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:

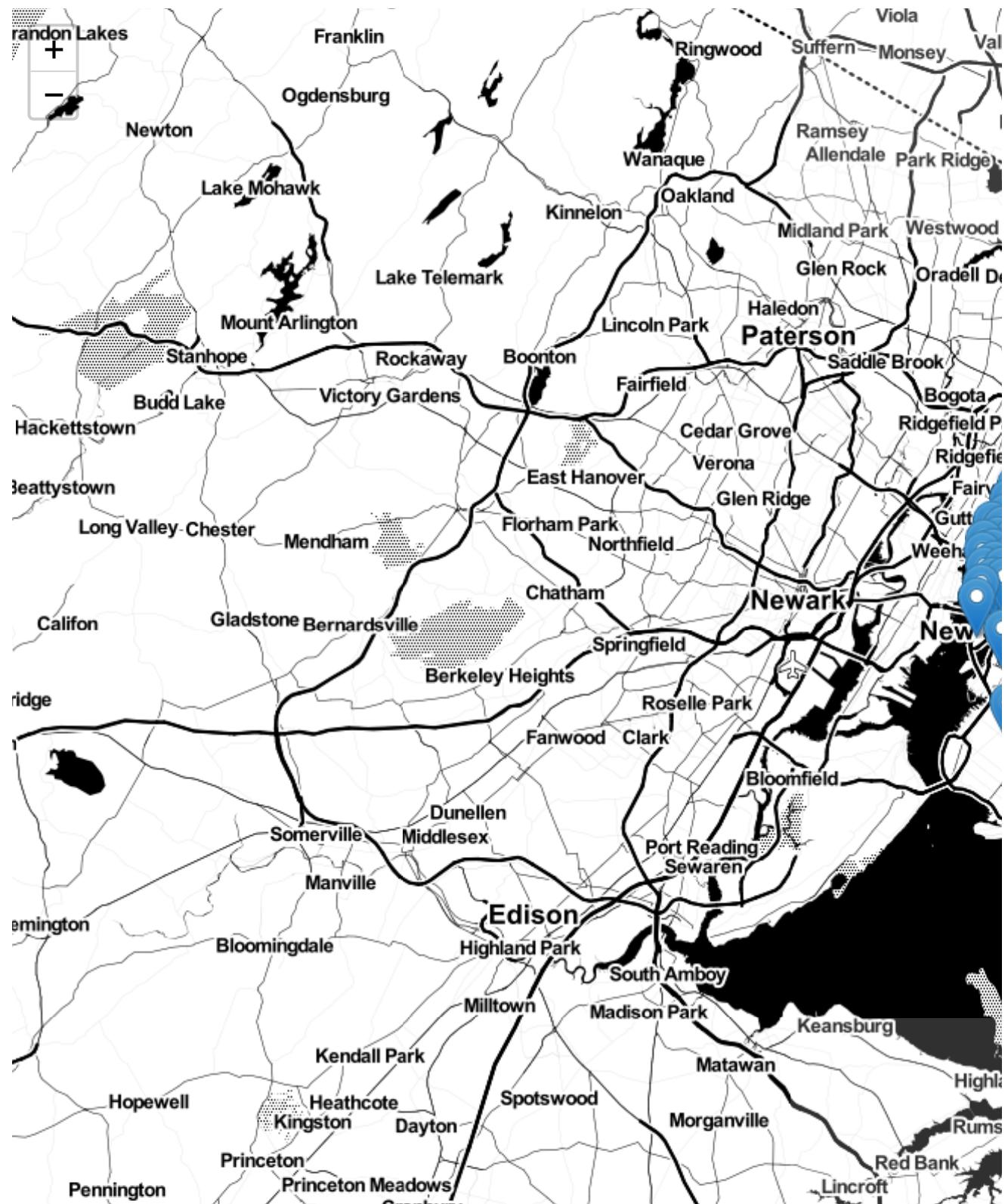
- The main objective was to find a optimal min. distance (Which roughly estimates to the radius of a cluster)

```
# if check for the 50 clusters you can observe that there are two clusters with only 0.3 miles ap  
# 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_out
```

## ▼ Plotting the cluster centers:

```
# 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
```





## ▼ Plotting the clusters:

```
#Visualising the clusters on a map
def plot_clusters(frame):
    city_long_border = (-74.03, -73.75)
    city_lat_border = (40.63, 40.85)
    fig, ax = plt.subplots(ncols=1, nrows=1)
    ax.scatter(frame.pickup_longitude.values[:100000], frame.pickup_latitude.values[:100000], s=1
               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

```
#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 es
    tenminutewise_binned_unix_pickup_times=[(int((i-start_pickup_unix)/600)+33) for i in unix_pic
frame['pickup_bins'] = np.array(tenminutewise_binned_unix_pickup_times)
return frame

# clustering, making pickup bins and grouping by pickup cluster and pickup bins
frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(frame_with_durations_out
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(['pic

# 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()
```

	passenger_count	trip_distance	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
0	1	1.59	-73.993896	40.750111	-73.9747	40.750111
1	1	3.30	-74.001648	40.724243	-73.9944	40.724243
2	1	1.80	-73.963341	40.802788	-73.9518	40.802788
3	1	0.50	-74.009087	40.713818	-74.0043	40.713818
4	1	3.00	-73.971176	40.762428	-74.0041	40.762428

```
# hear the trip_distance represents the number of pickups that are happennd in that particular 10m
# this data frame has two indices
# primary index: pickup_cluster (cluster number)
# secondary index : pickup_bins (we devide whole months time into 10min intravels 24*31*60/10 =446
jan_2015_groupby.head()
```



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

```
# 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)
#frame_with_durations_outliers_removed_2016['pickup_cluster'] = kmeans.predict(frame_with_durations)

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

    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

```

# 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
# 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

# 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)  
  
# 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(j  
    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:  $\lfloor \frac{x}{4} \rfloor, \lceil \frac{x}{4} \rceil, \lceil \frac{x}{4} \rceil, \lceil \frac{x}{4} \rceil$   
 Ex2:  $\lfloor \frac{x}{3} \rfloor, \lceil \frac{x}{3} \rceil, \lceil \frac{x}{3} \rceil$
  - Case 2:(values missing in middle)
 

Ex1:  $\lfloor \frac{x+y}{4} \rfloor, \lceil \frac{x+y}{4} \rceil, \lceil \frac{x+y}{4} \rceil, \lceil \frac{x+y}{4} \rceil$   
 Ex2:  $\lfloor \frac{x+y}{5} \rfloor, \lceil \frac{x+y}{5} \rceil, \lceil \frac{x+y}{5} \rceil, \lceil \frac{x+y}{5} \rceil, \lceil \frac{x+y}{5} \rceil$
  - Case 3:(values missing at the end)
 

Ex1:  $\lfloor \frac{x}{4} \rfloor, \lceil \frac{x}{4} \rceil, \lceil \frac{x}{4} \rceil, \lceil \frac{x}{4} \rceil$   
 Ex2:  $\lfloor \frac{x}{2} \rfloor, \lceil \frac{x}{2} \rceil$

```
# 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
```

```
# 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 smoothed data (which is calculated based on the methods that are discussed in the
# we finally return smoothed data
def smoothing(count_values,values):
    smoothed_regions=[] # stores list of final smoothed values of each region
    ind=0
    repeat=0
    smoothed_value=0
    for r in range(0,40):
        smoothed_bins=[] #stores the final smoothed values
        repeat=0
        for i in range(4464):
            if repeat!=0: # prevents iteration for a value which is already visited/resolved
                repeat-=1
                continue
            if i in values[r]: #checks if the pickup-bin exists
                smoothed_bins.append(count_values[ind]) # appends the value of the pickup bin if
            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 va
                            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
                        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_l
                            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,henc
                            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

```

```

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

```

```

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

```

```

# 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

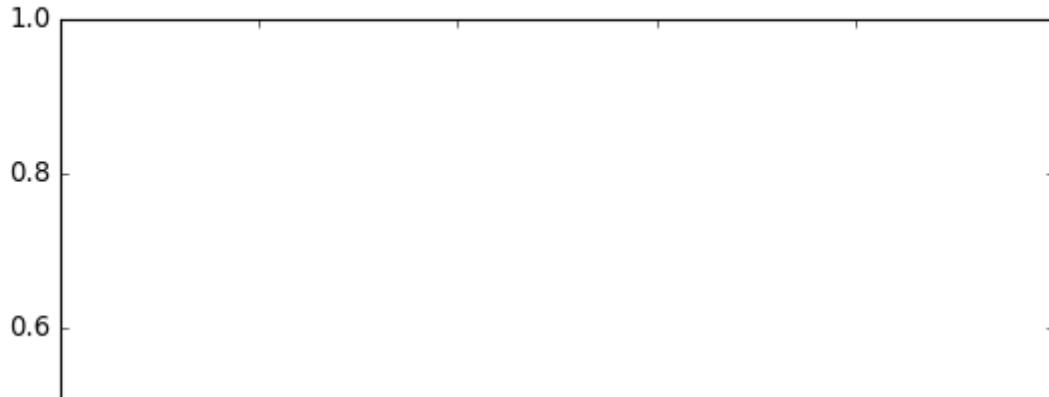
```

```
# 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

```
# 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()
```





```
# 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, 0, 20, i.e there are 10 pickups
# 10st 10min intravel, 0 pickups happened in 2nd 10mins intravel, 0 pickups happened in 3rd 10min
# 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
# that are happened in the first 40min are same in both cases, but if you can observe that we loo
# wheen you are using smoothing we are looking at the future number of pickups which might cause

# 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.
```

```
# 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
regions_cum = []

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

# number of 10min indices for jan 2015= 24*31*60/10 = 4464
# number of 10min indices for jan 2016 = 24*31*60/10 = 4464
# number of 10min indices for feb 2016 = 24*29*60/10 = 4176
# number of 10min indices for march 2016 = 24*31*60/10 = 4464
# regions_cum: it will contain 40 lists, each list will contain 4464+4176+4464 values which repre
# 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_
```

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

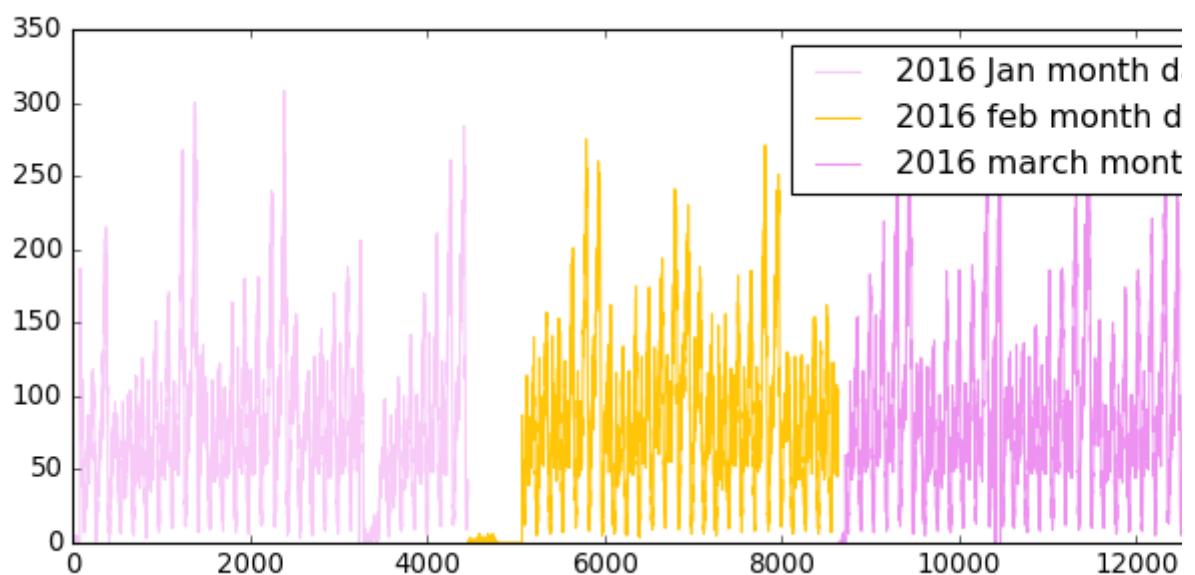
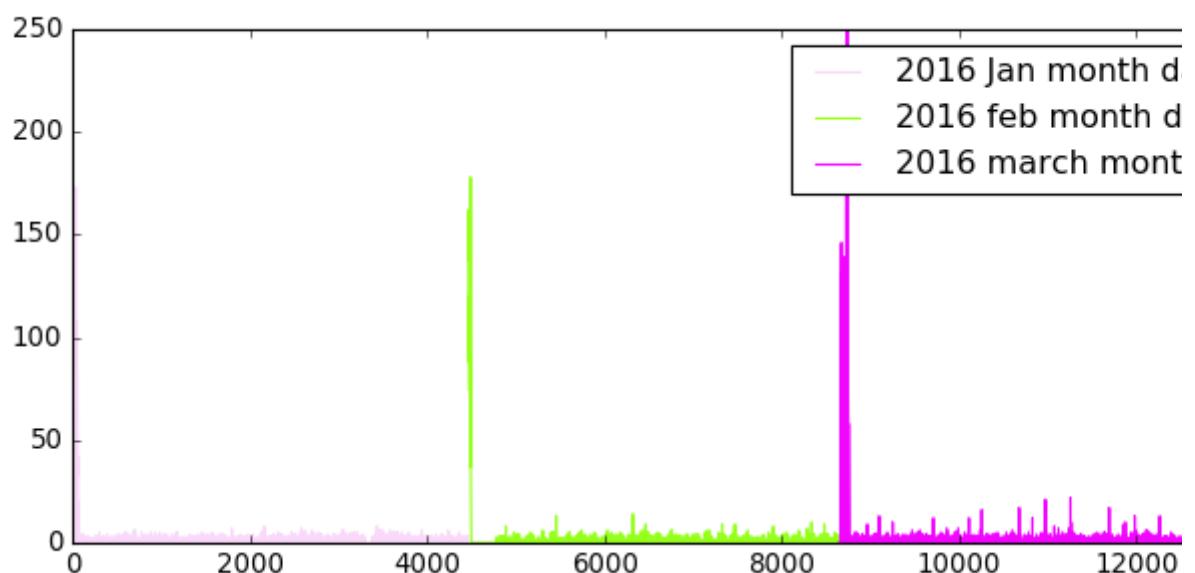
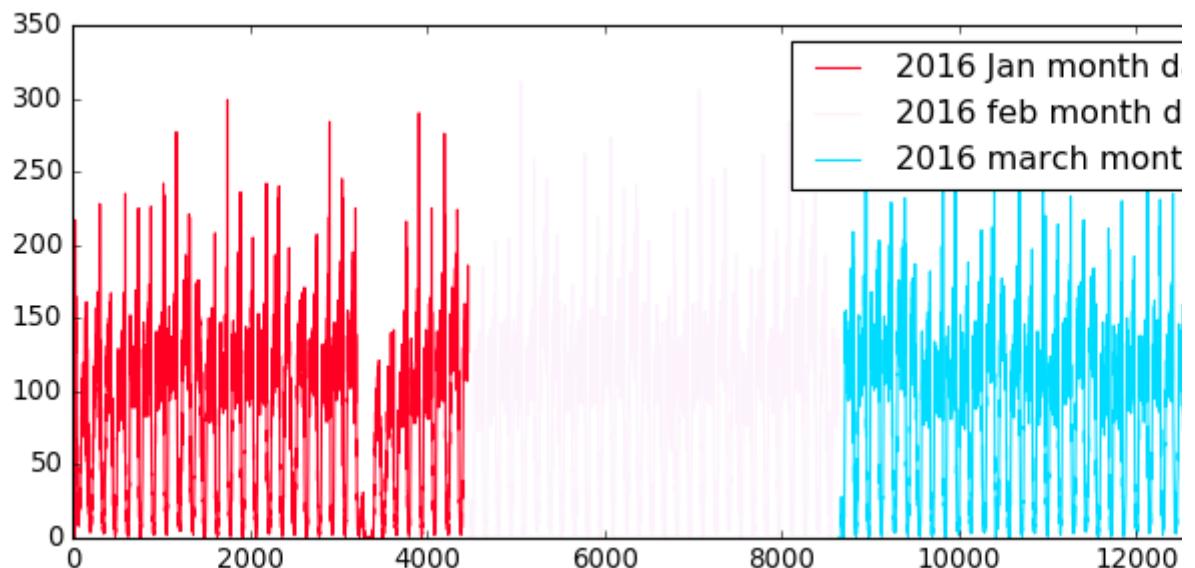


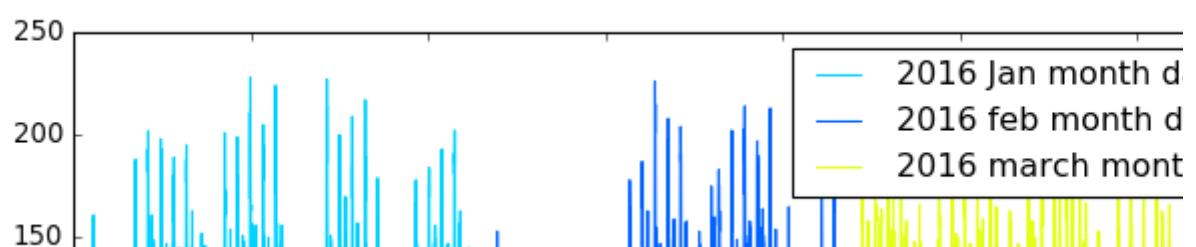
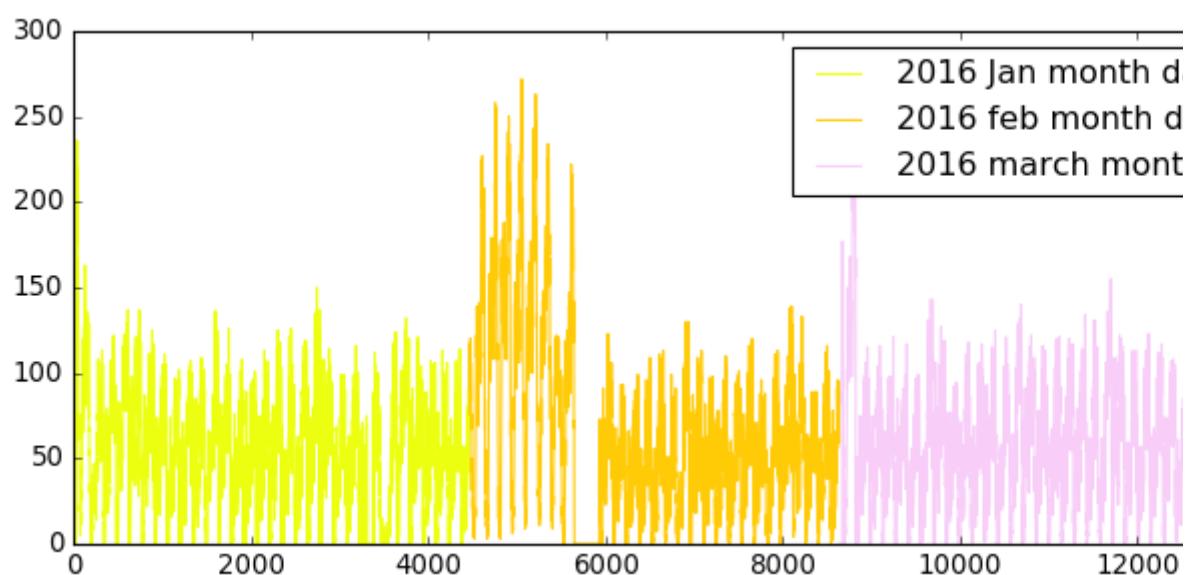
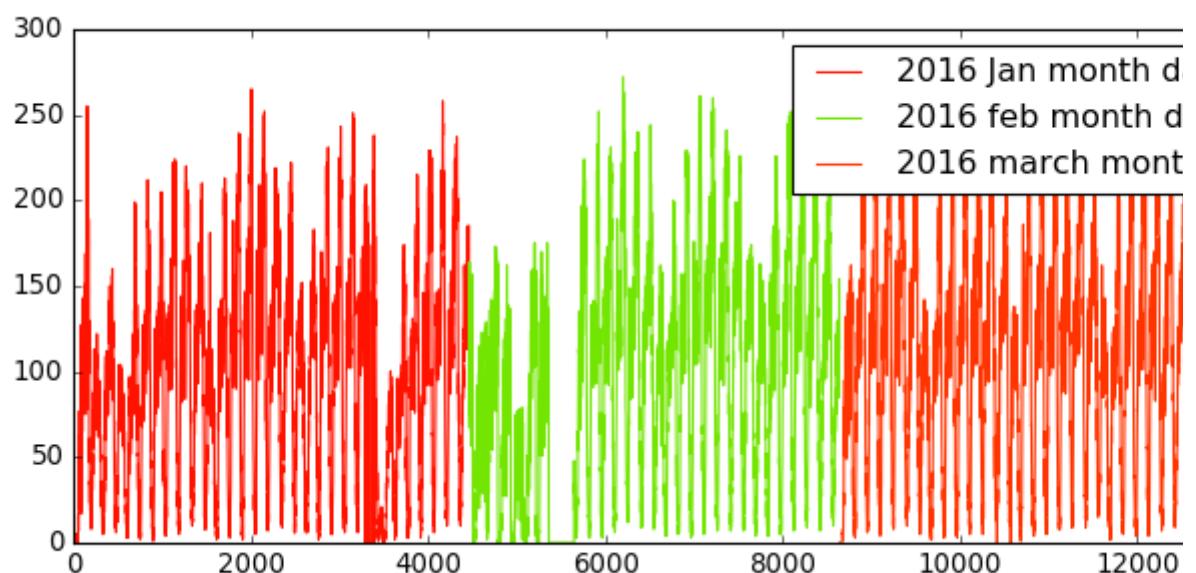
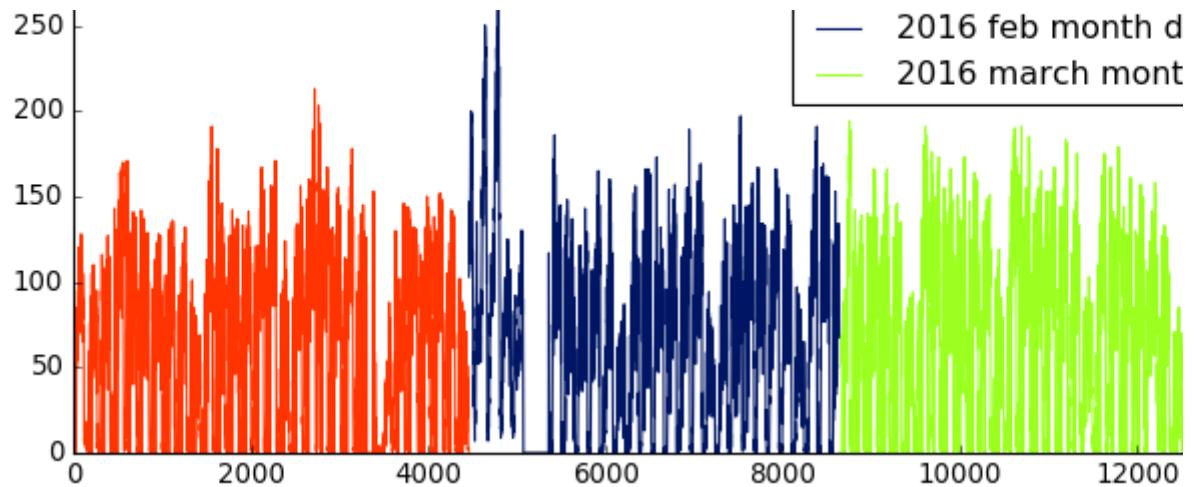
## ▼ Time series and Fourier Transforms

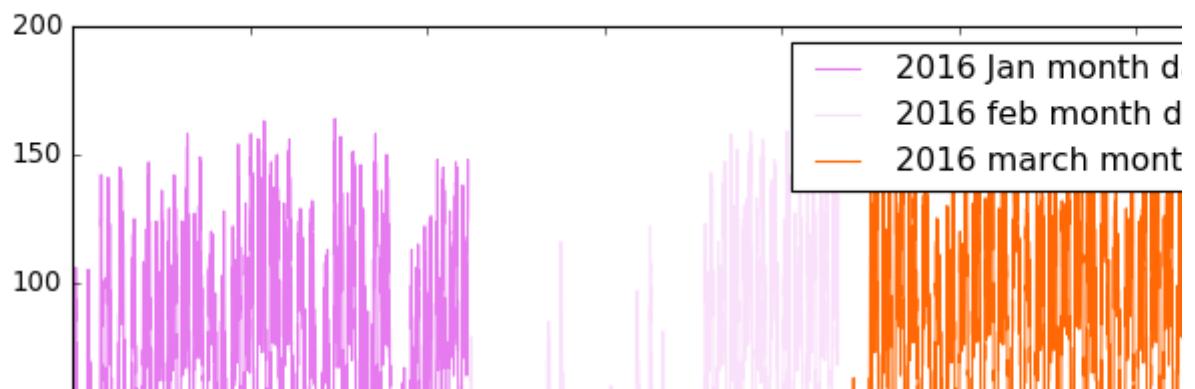
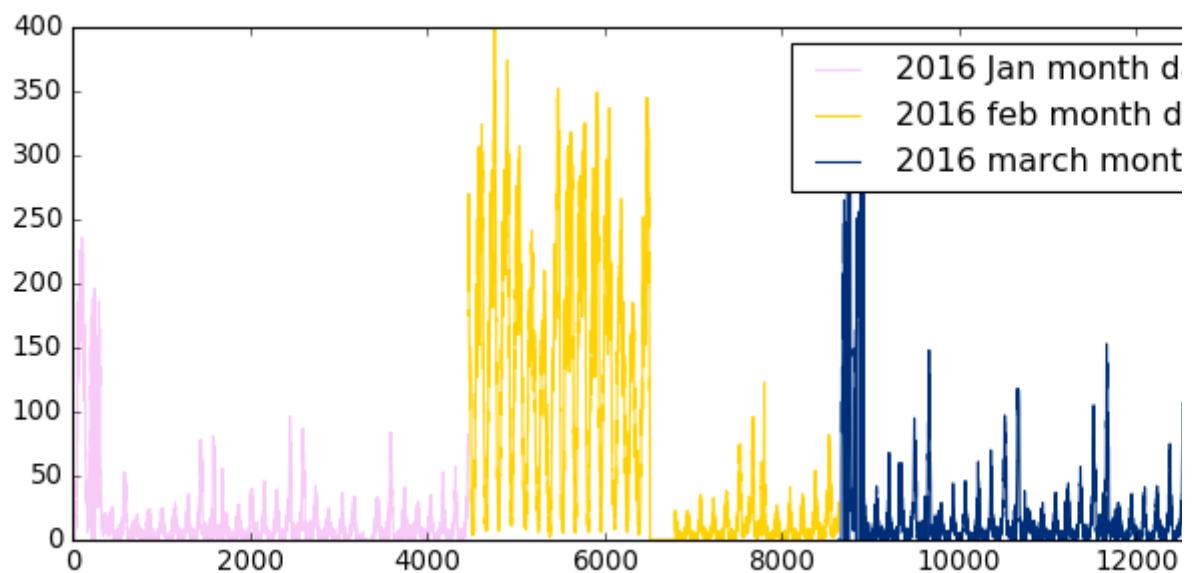
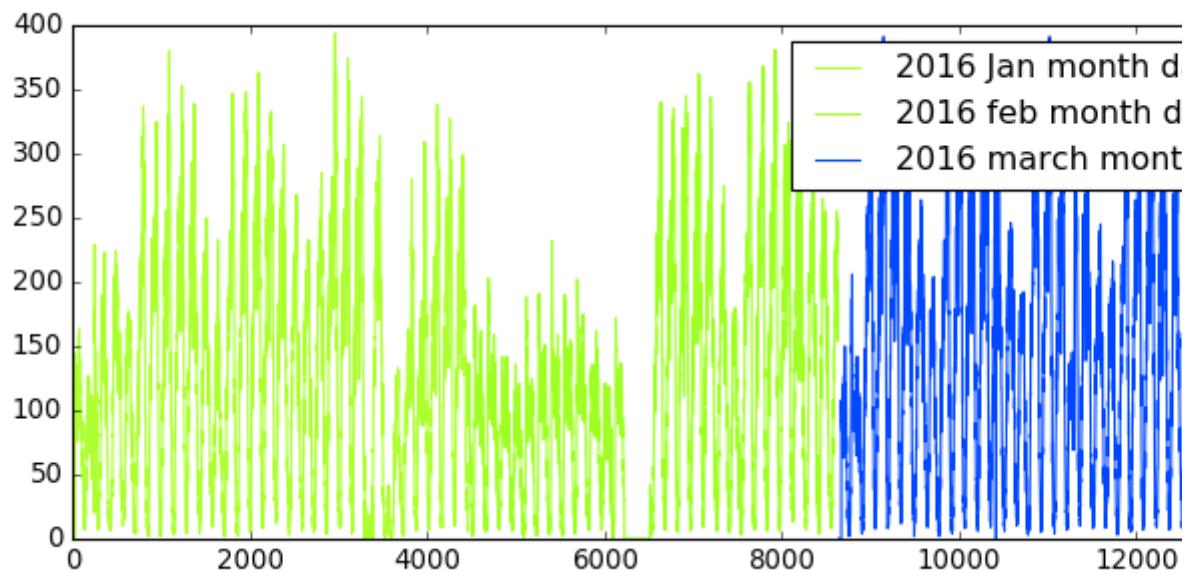
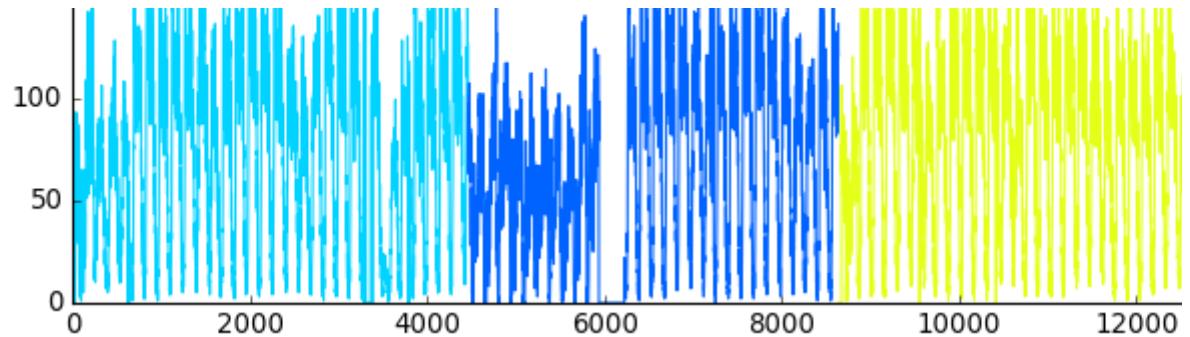
```
def uniqueish_color():
    """There're better ways to generate unique colors, but this isn't awful."""
    return plt.cm.gist_ncar(np.random.random())
first_x = list(range(0,4464))
second_x = list(range(4464,8640))
third_x = list(range(8640,13104))
for i in range(40):
```

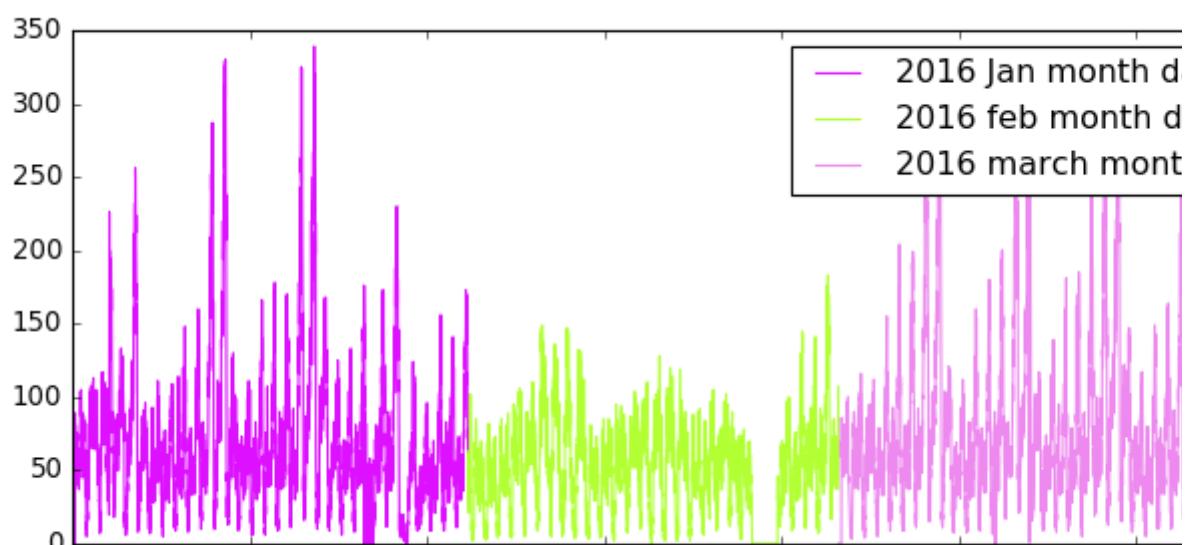
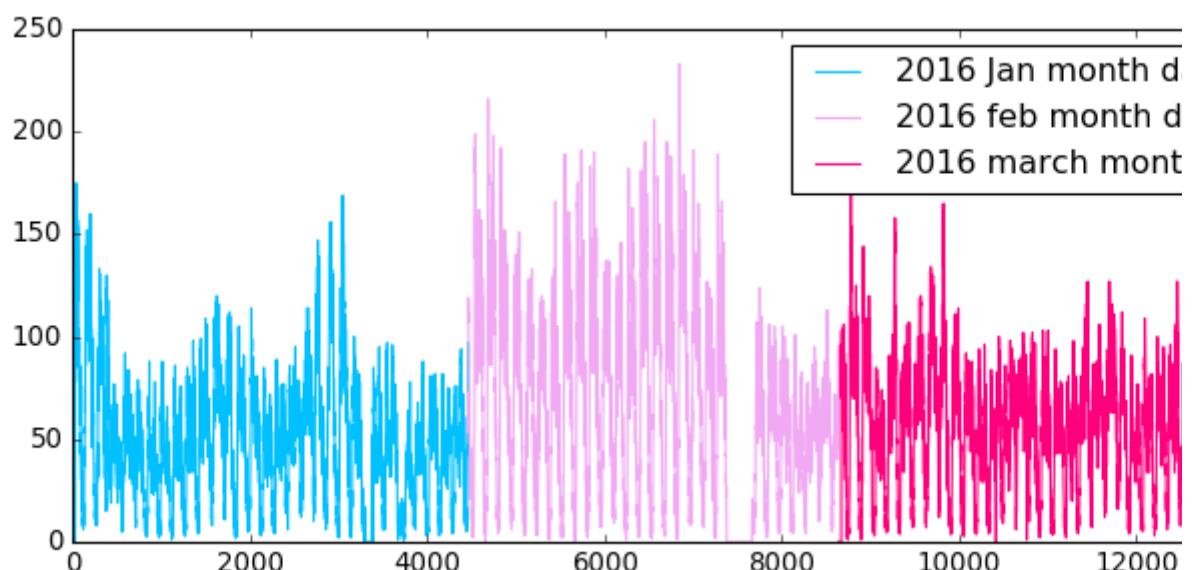
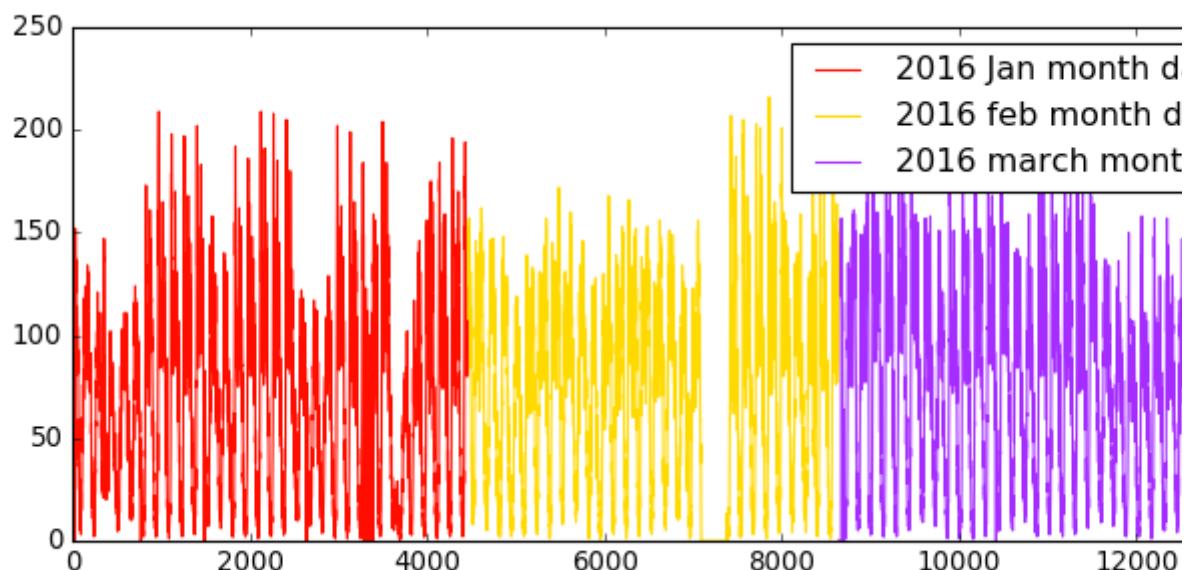
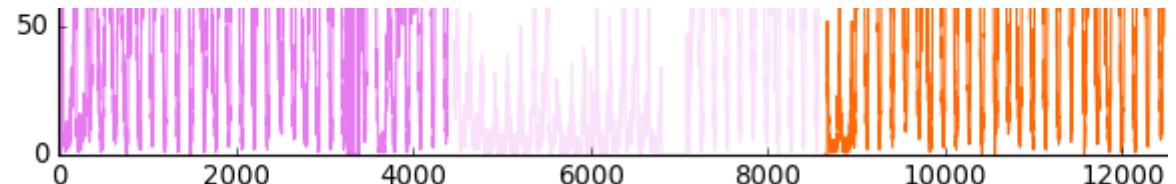
```
plt.figure(figsize=(10,4))
plt.plot(first_x,regions_cum[i][:4464], color=uniqueish_color(), label='2016 Jan month data')
plt.plot(second_x,regions_cum[i][4464:8640], color=uniqueish_color(), label='2016 feb month d
plt.plot(third_x,regions_cum[i][8640:], color=uniqueish_color(), label='2016 march month data
plt.legend()
plt.show()
```

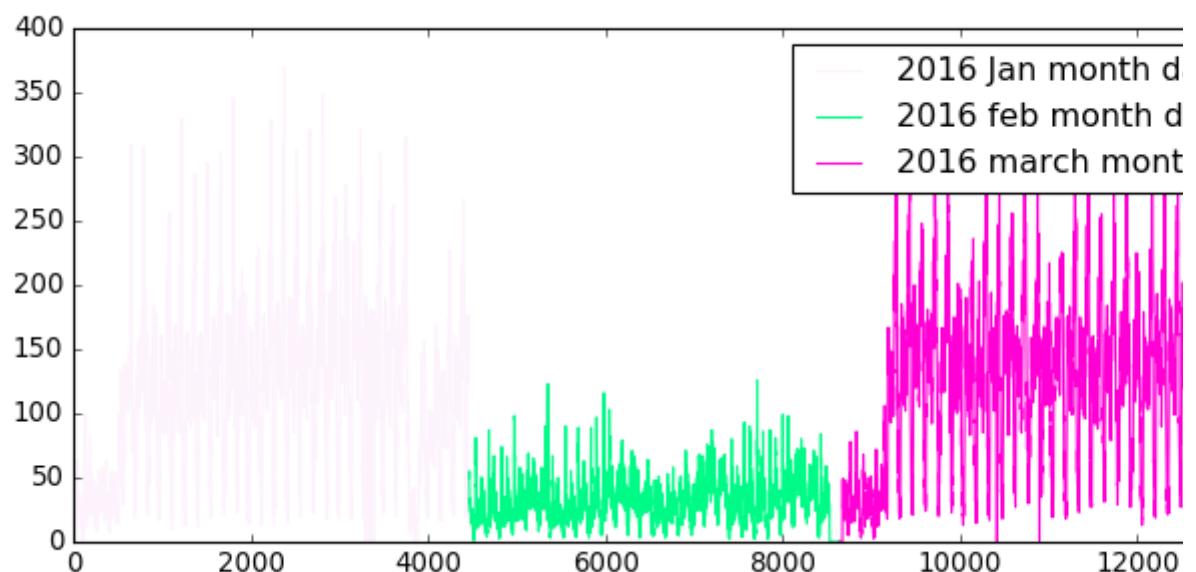
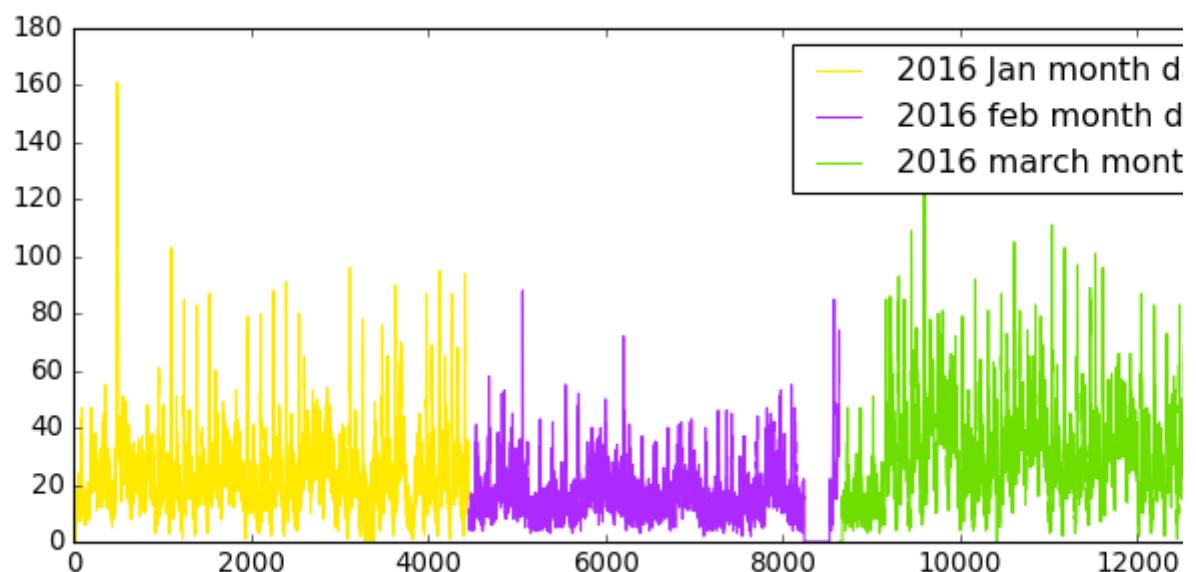
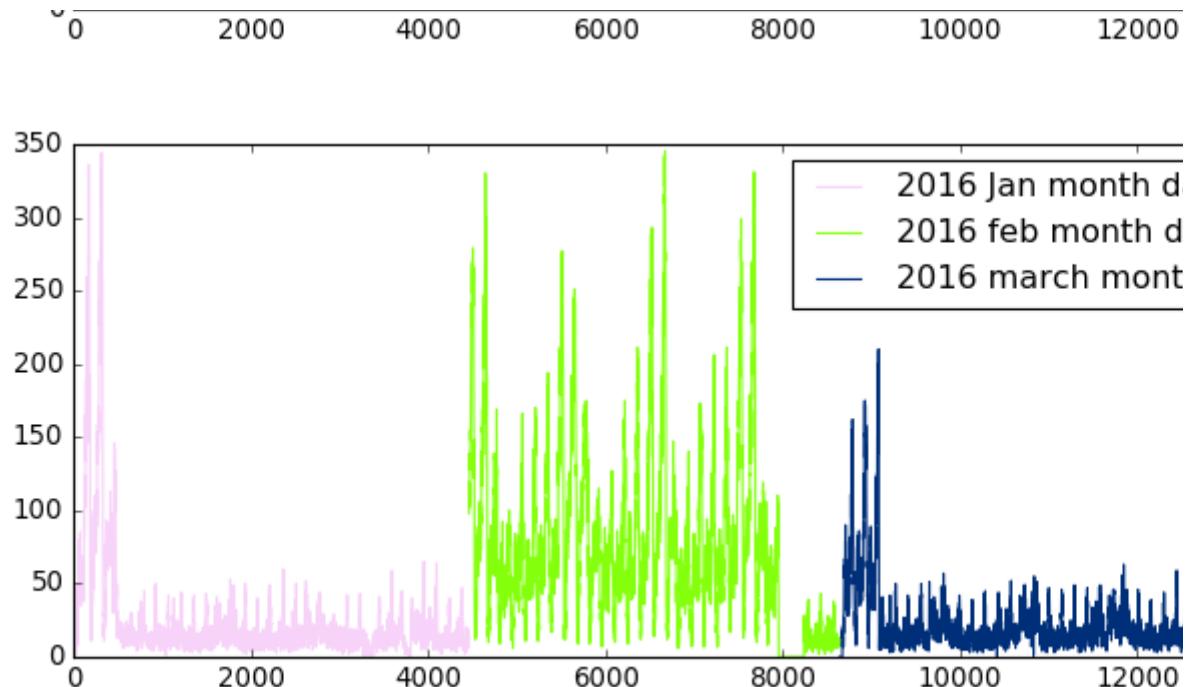


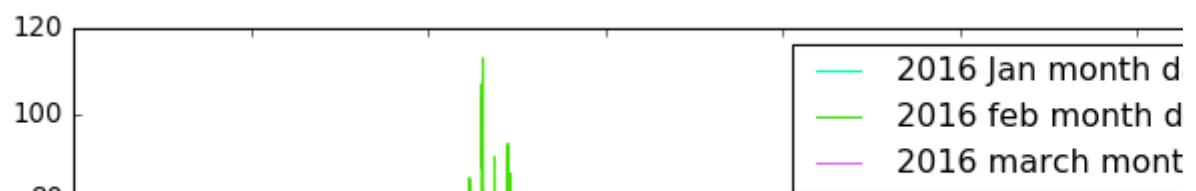
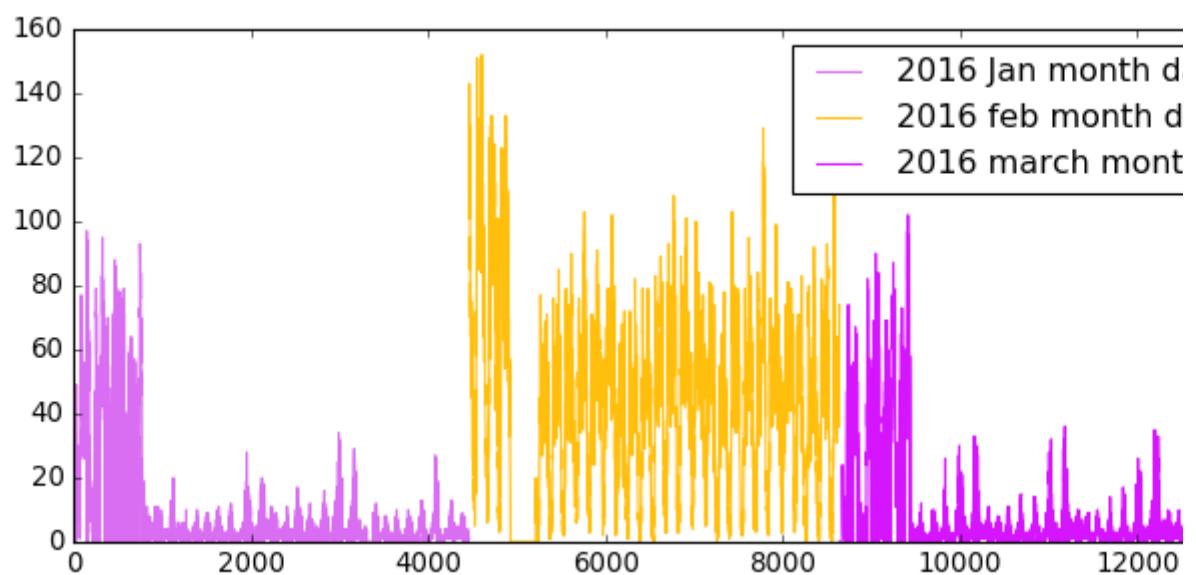
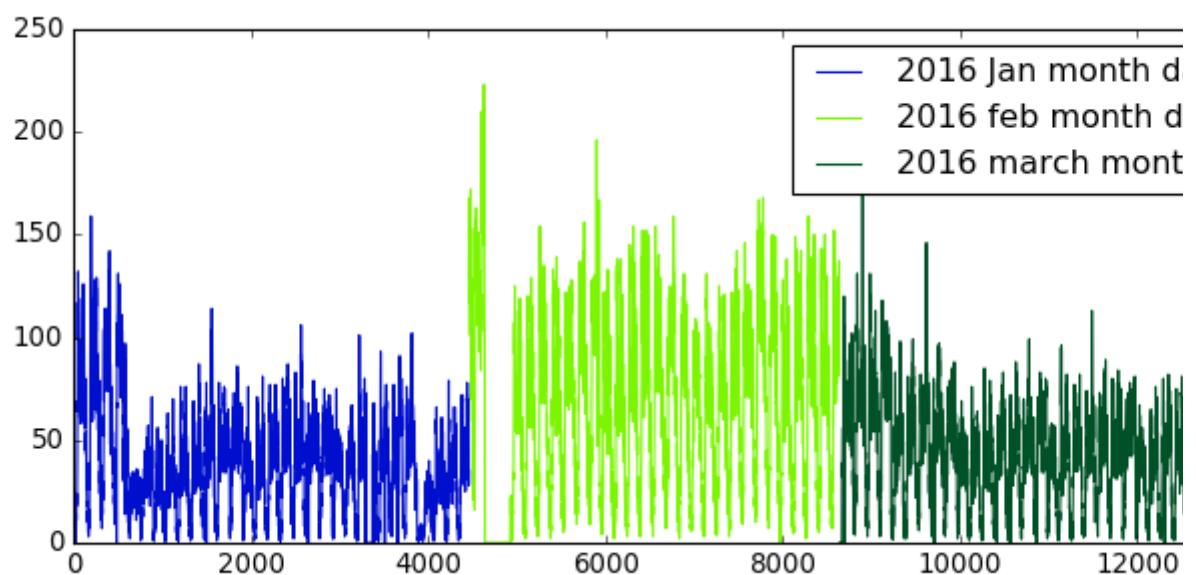
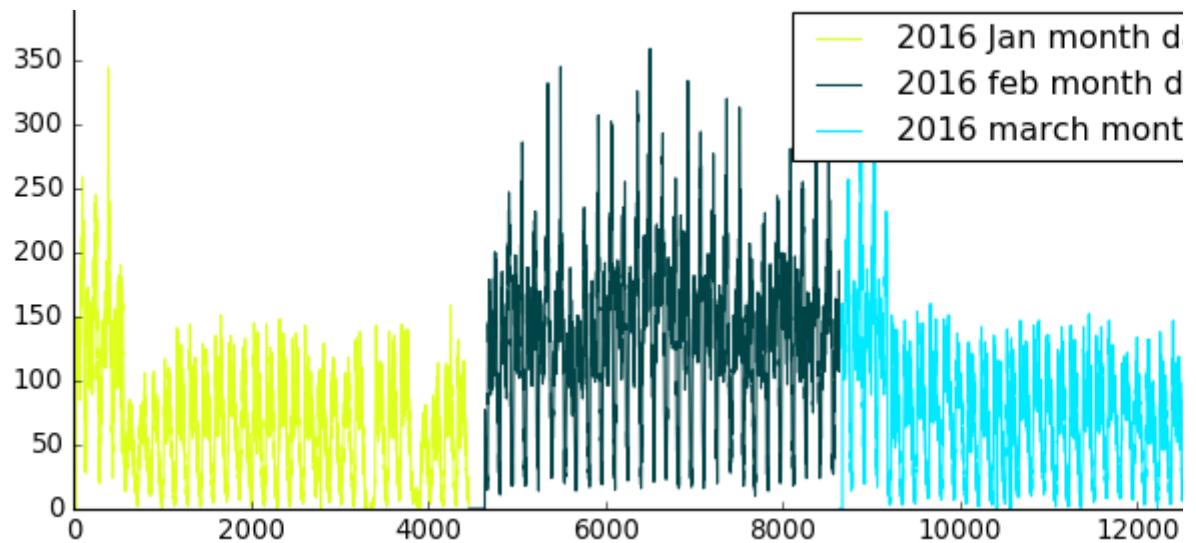


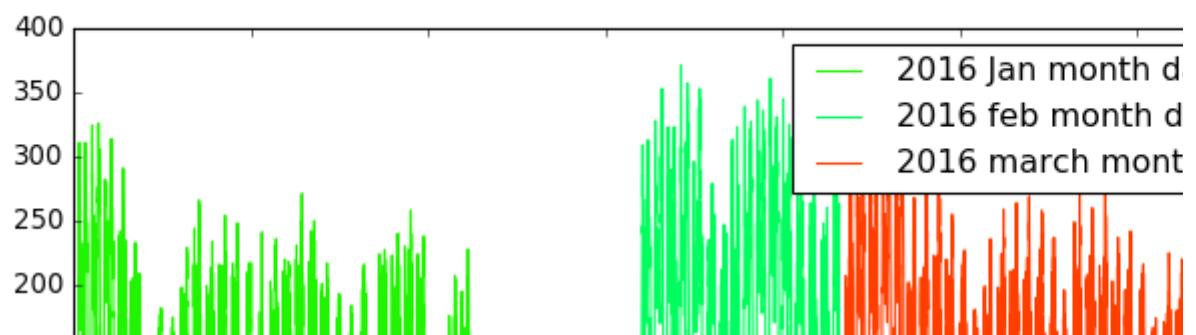
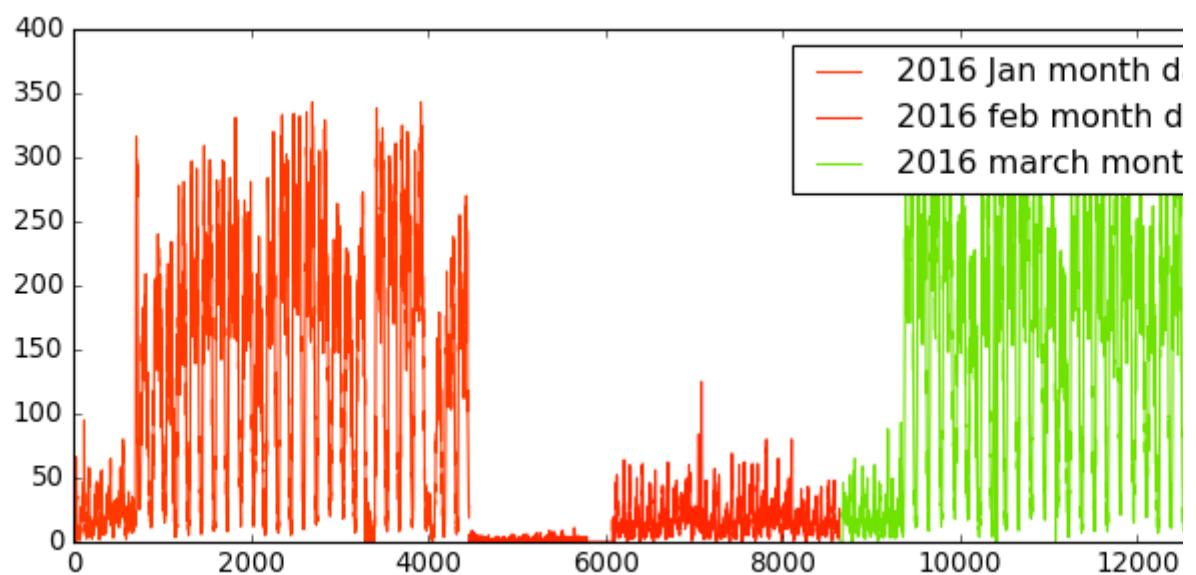
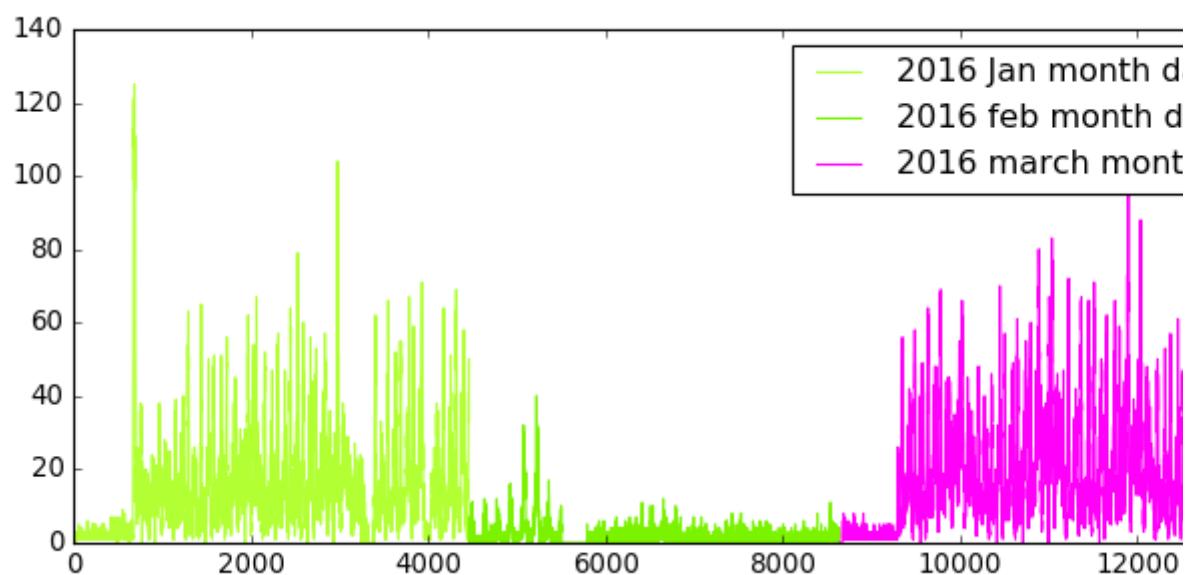
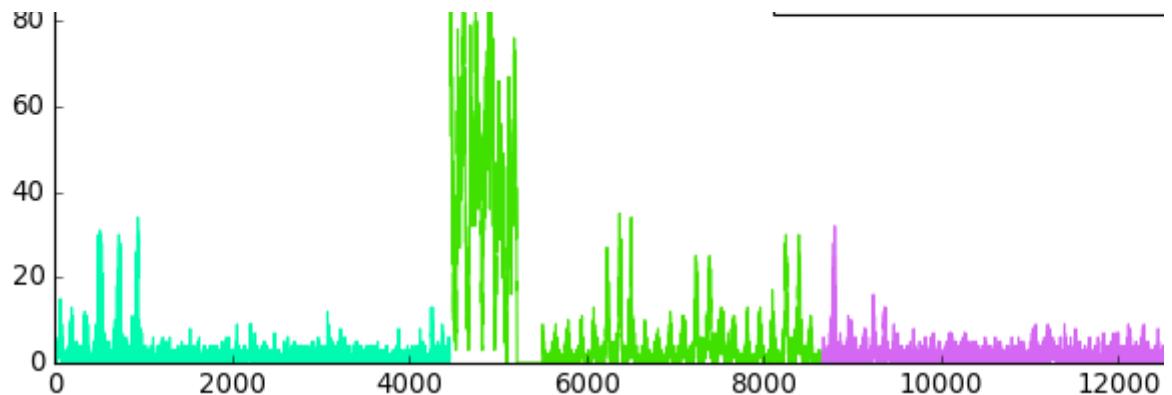


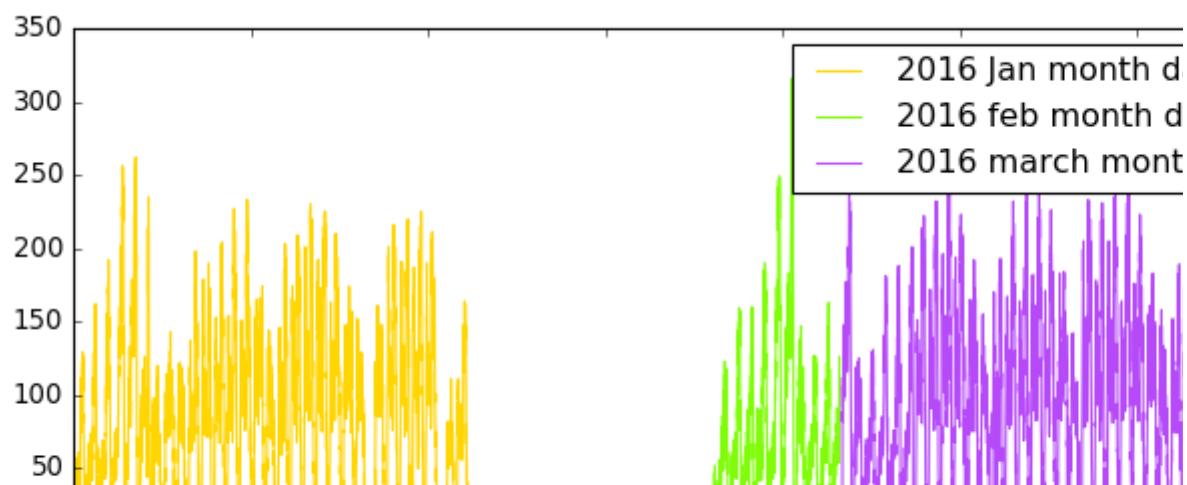
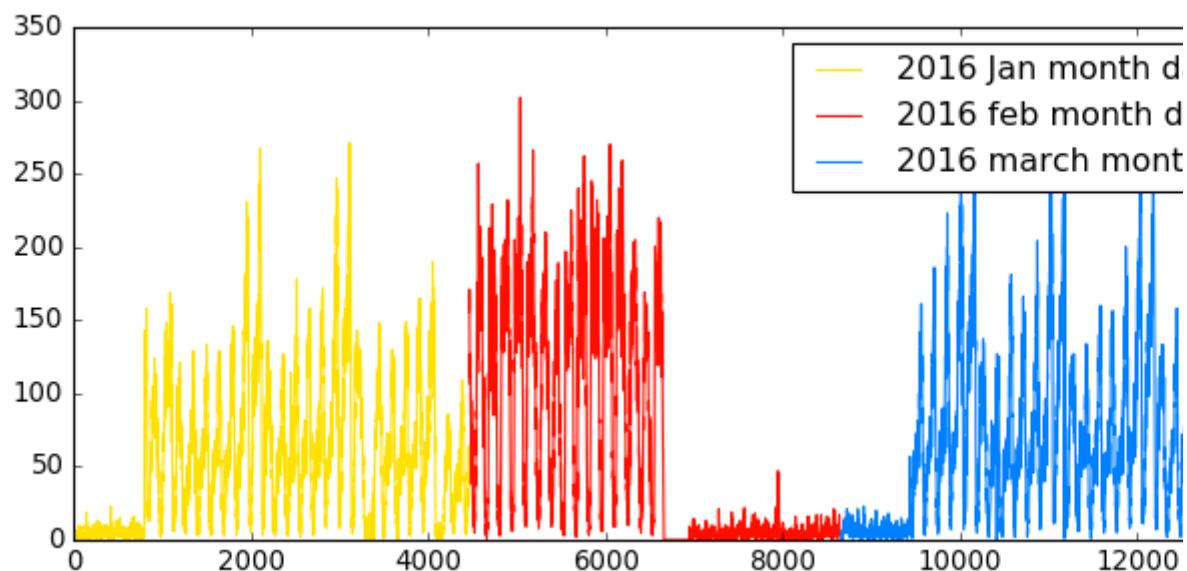
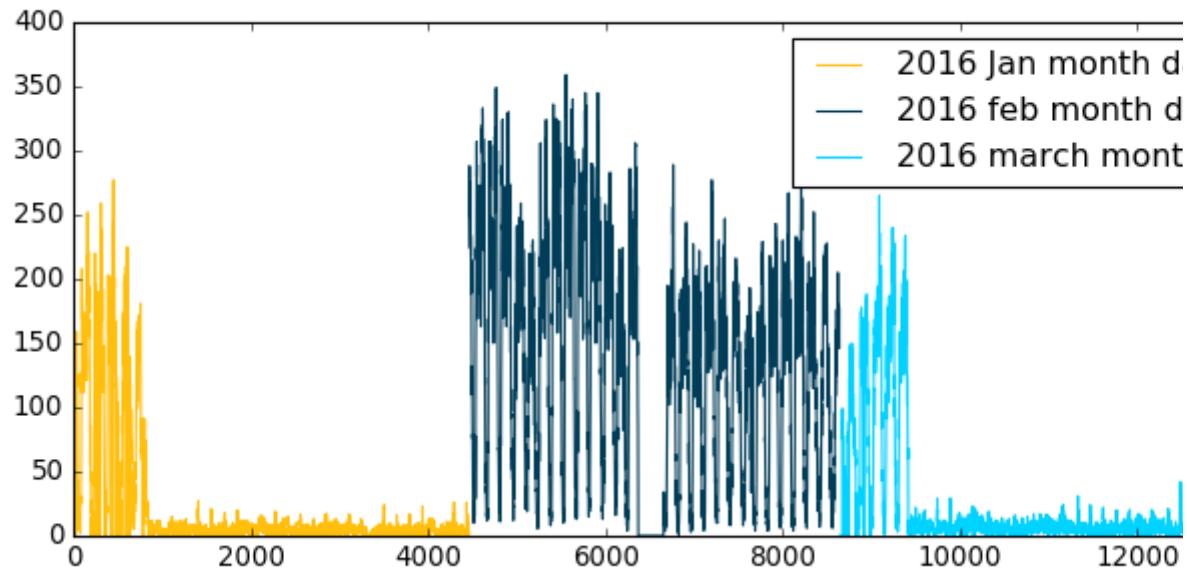
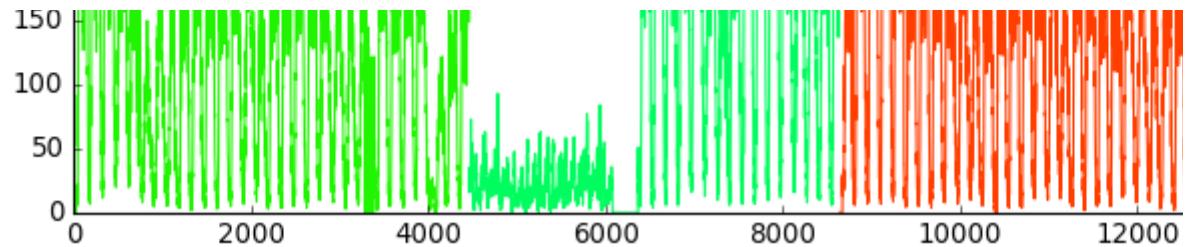


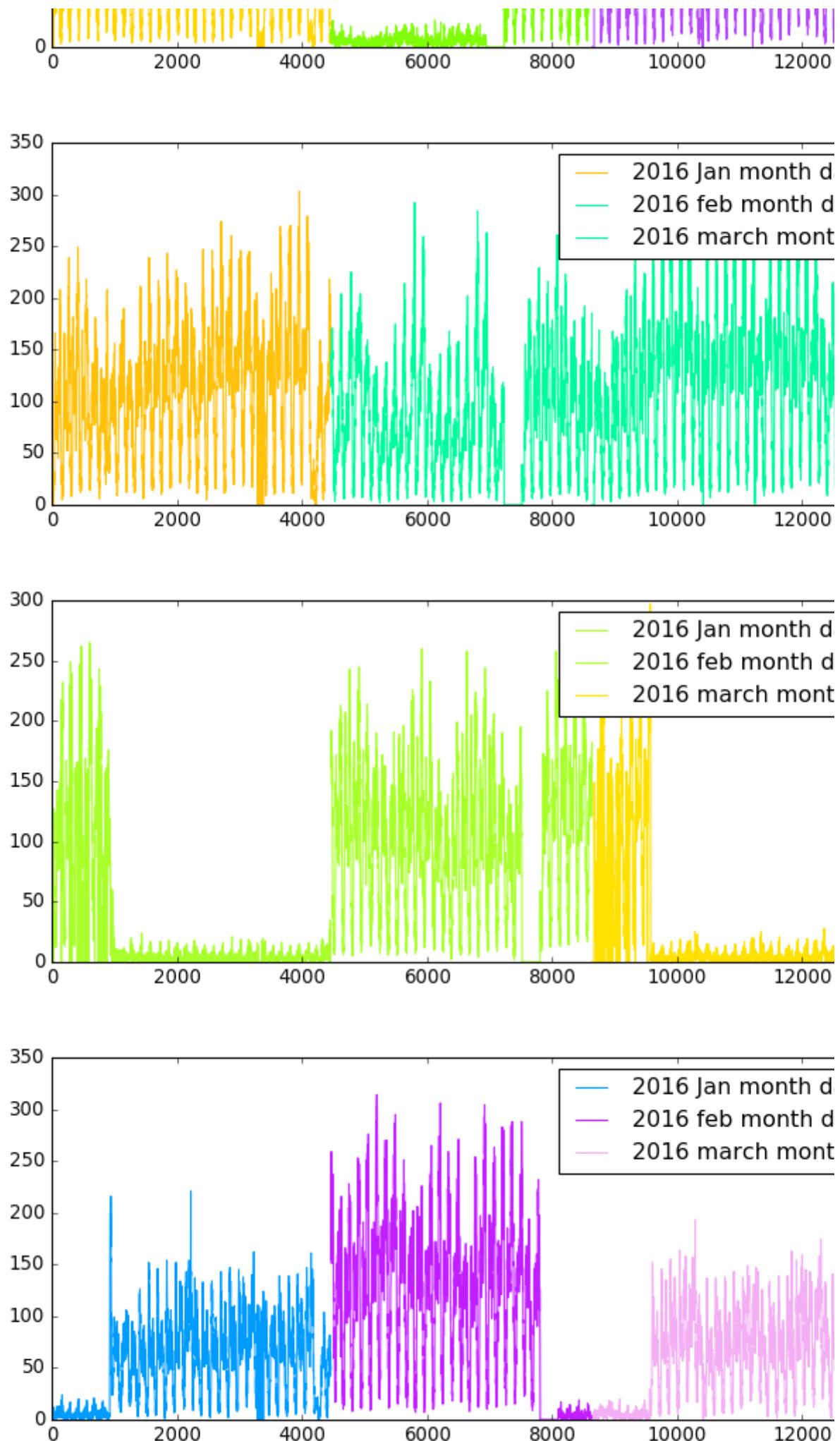


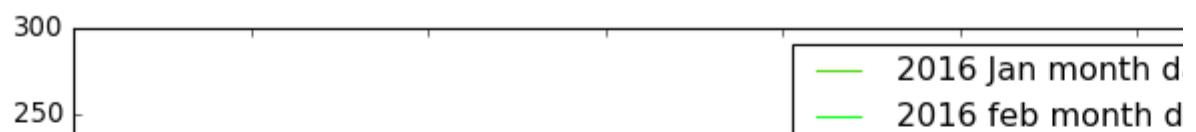
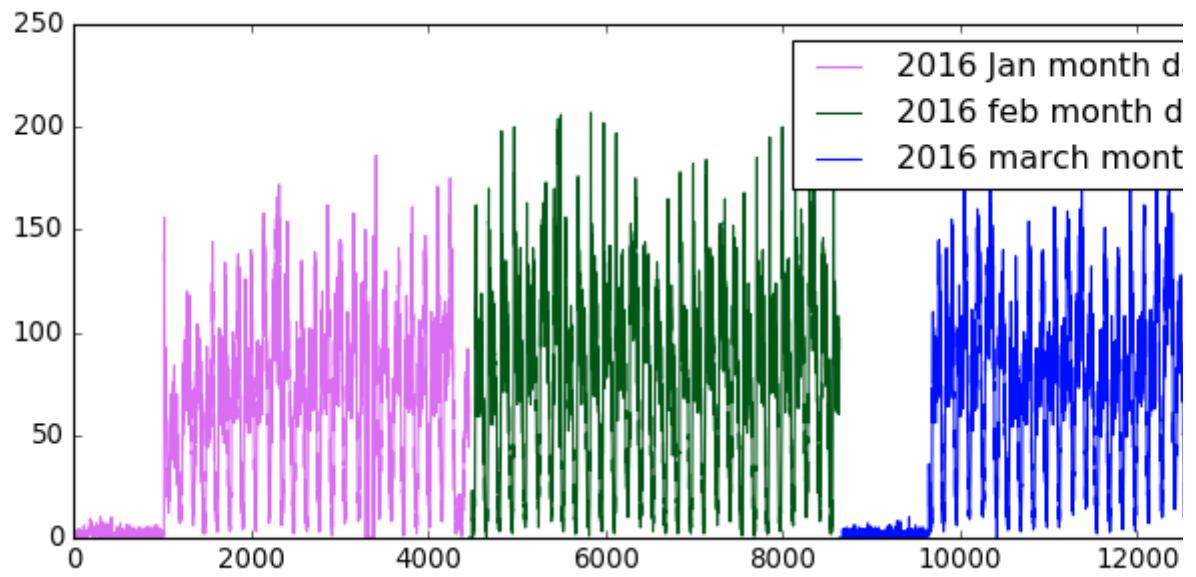
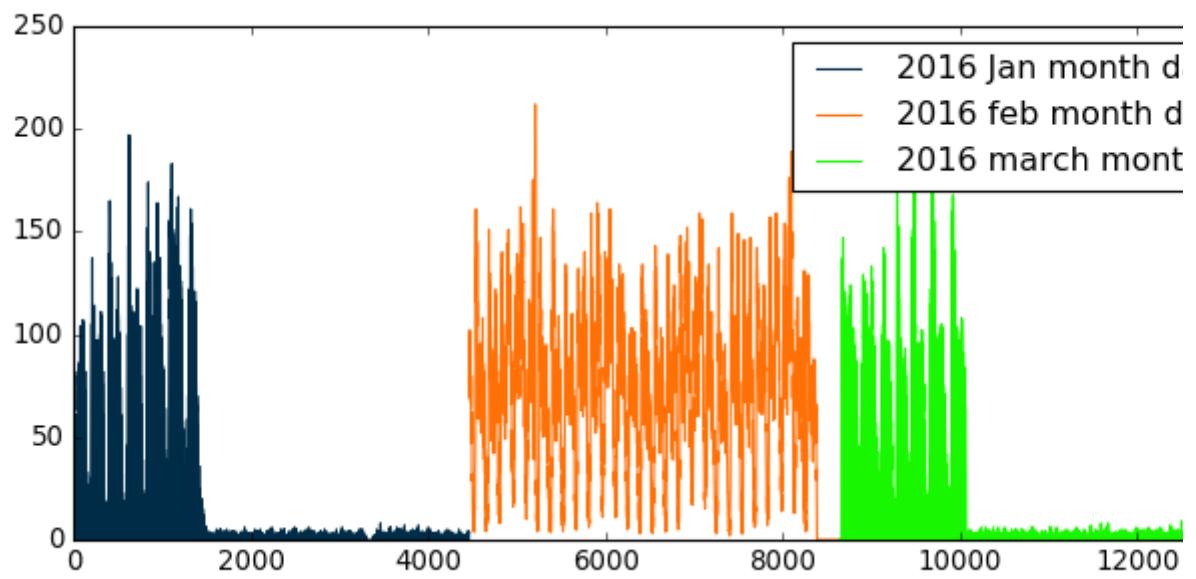
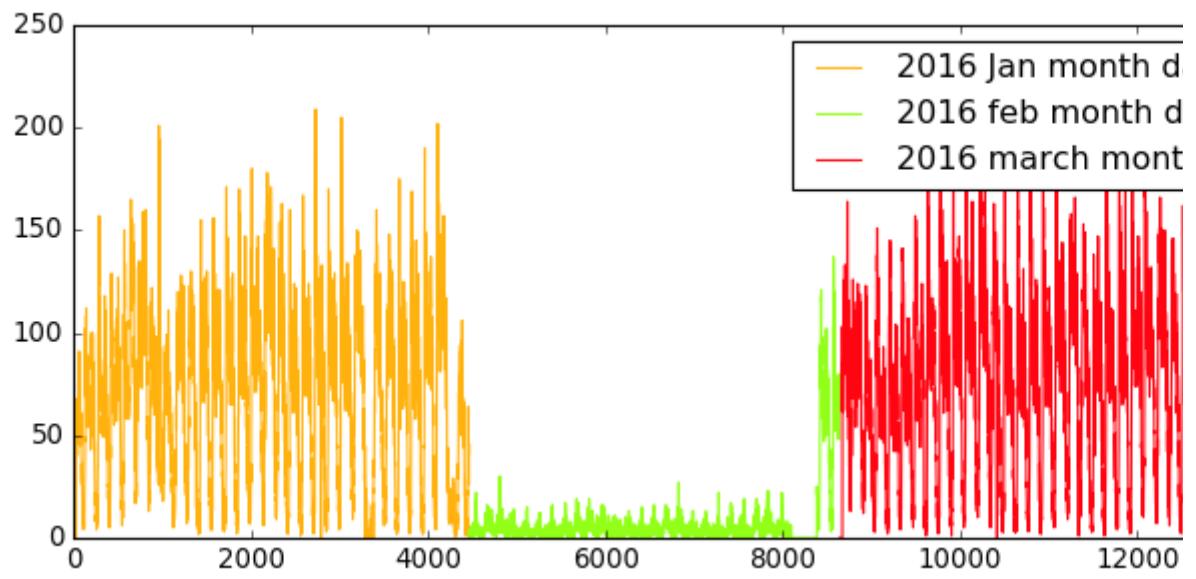


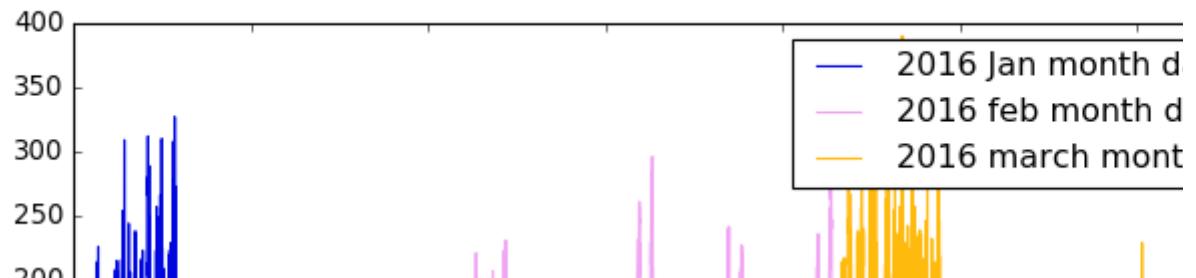
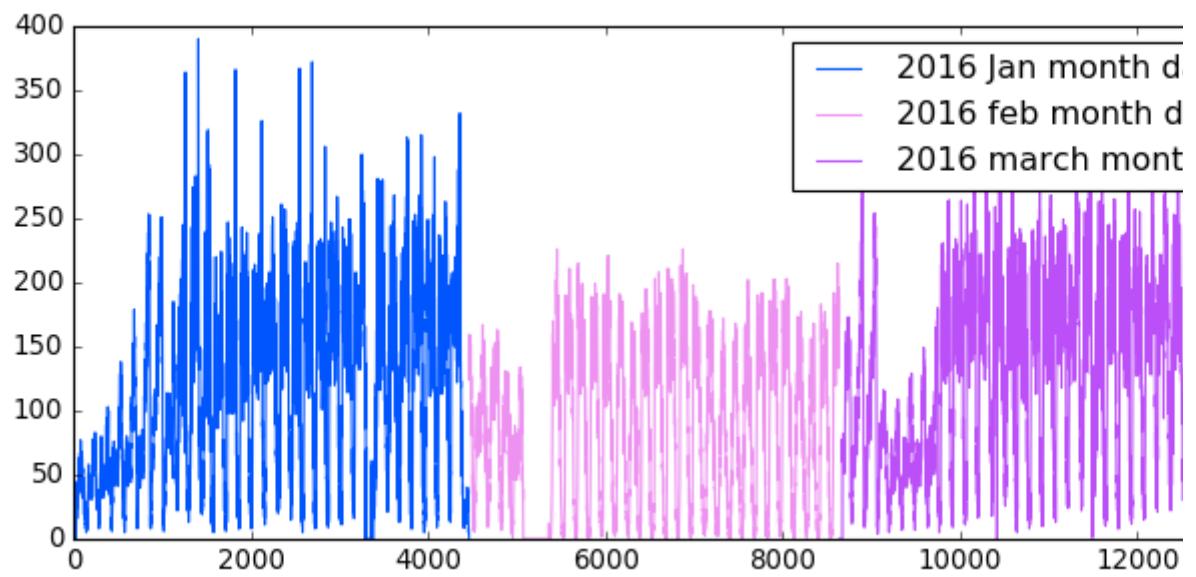
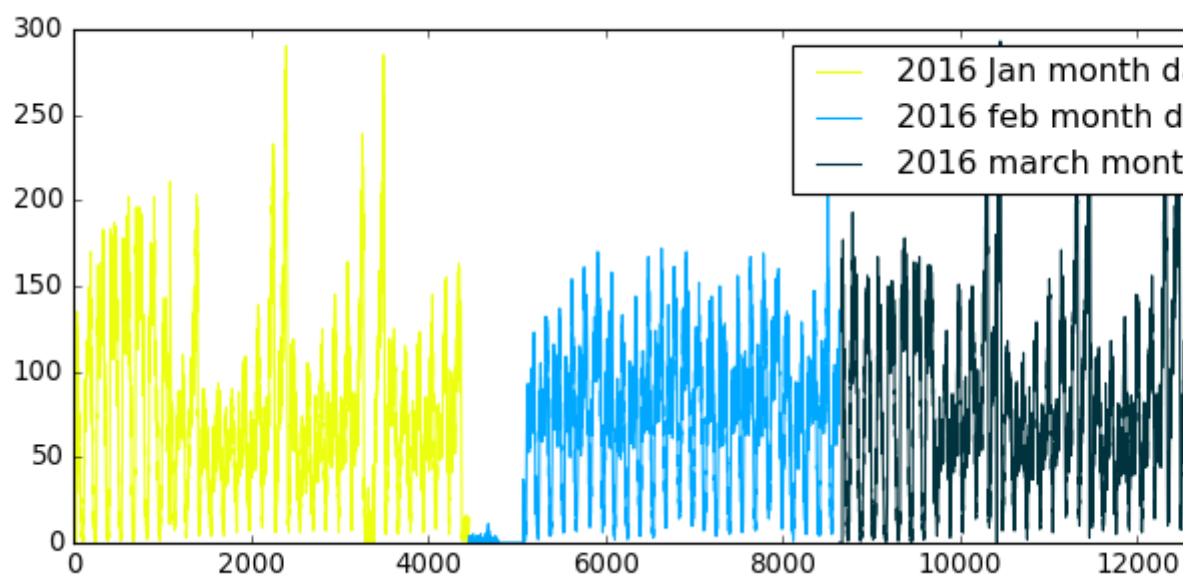
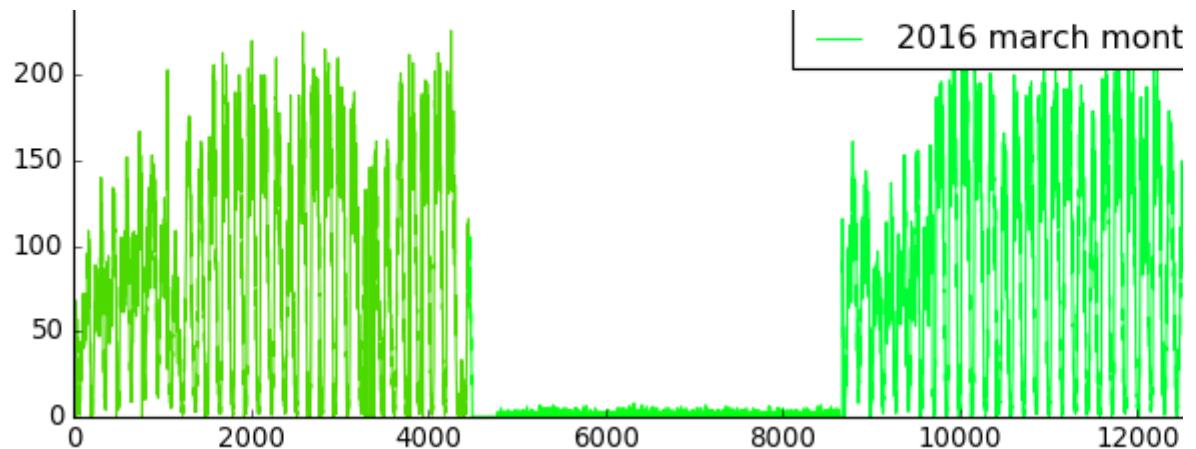


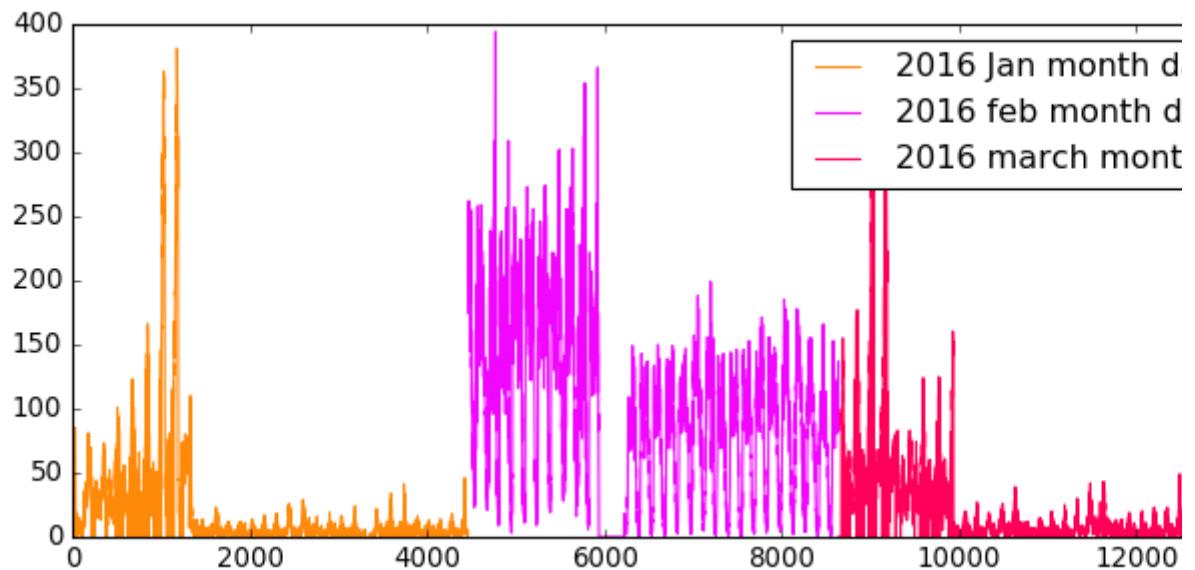
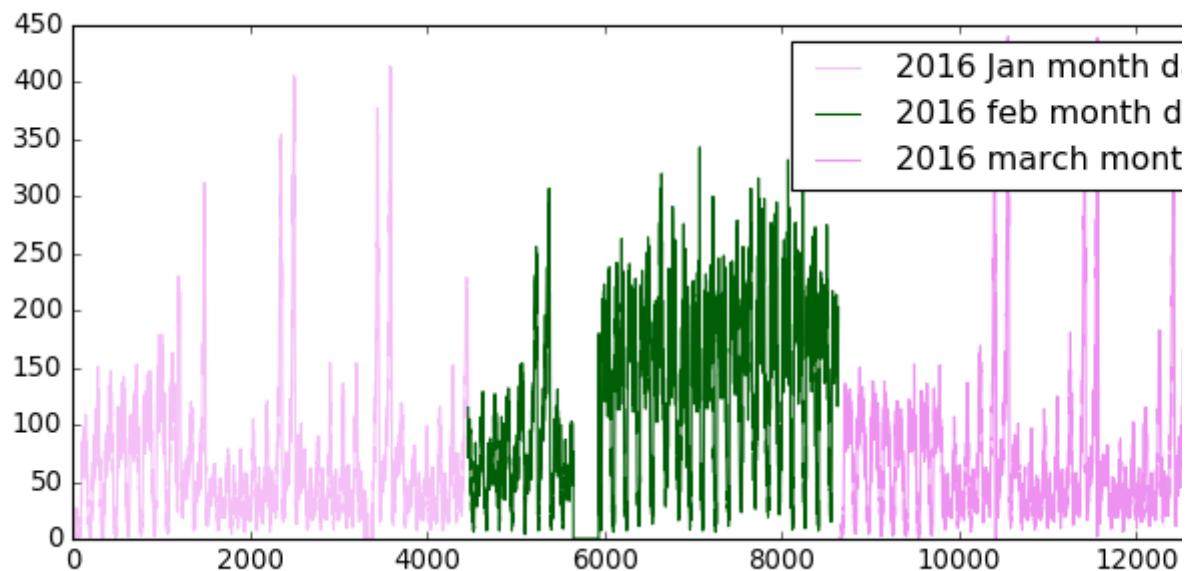
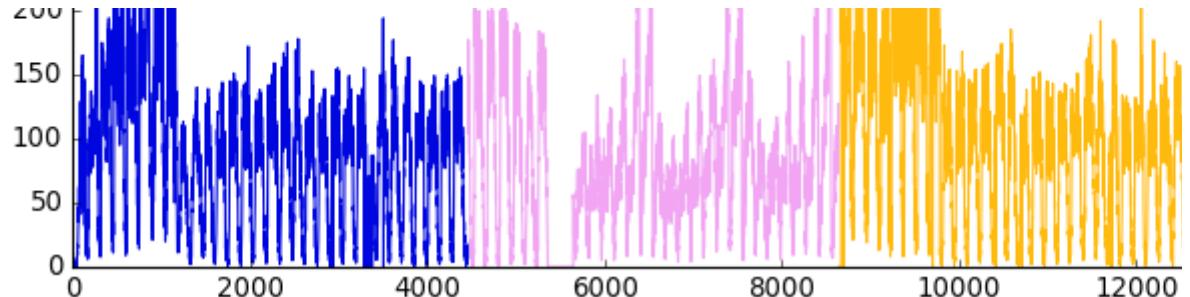






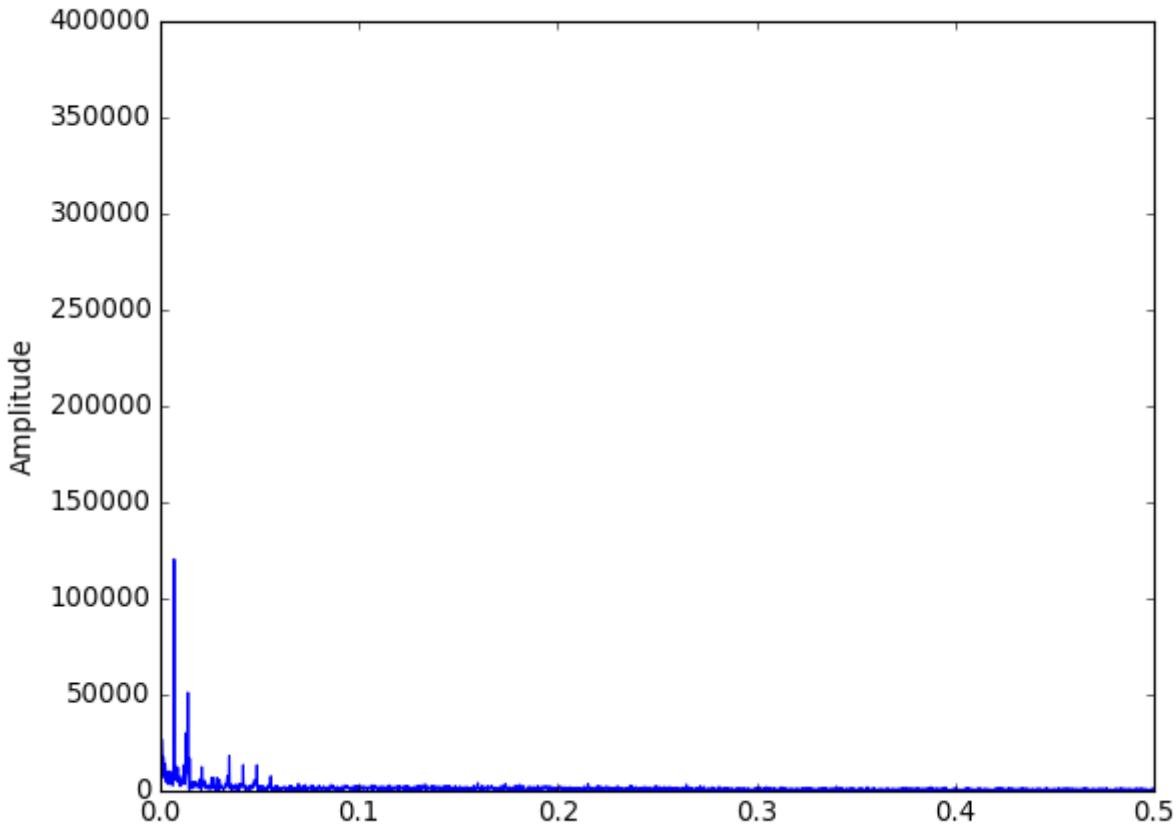






```
# getting peaks: https://blog.ytotech.com/2015/11/01/findpeaks-in-python/
# read more about fft function : https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.f
Y      = np.fft.fft(np.array(jan_2016_smooth)[0:4460])
# read more about the fftfreq: https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fft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





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

## ▼ Modelling: Baseline Models

Now we get into modelling in order to forecast the pickup densities for the months of Jan, Feb and March of 2016 variations

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

## ▼ Simple Moving Averages

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

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

```
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['Predicted'][i]),2)))
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']))
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 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  $P_t = (P_{t-1} + P_{t-2} + \dots + P_{t-3})/3$

```

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)])/windo
        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']))
    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 using Moving Averages using previous 2016 values therefore we get  $P_t = P_{t-1}$

## ▼ Weighted Moving Averages

The Moving Averages Model used gave equal importance to all the values in the window used, but we know intuitively the latest values are more similar to the older values. Weighted Averages converts this analogy into a mathematical model by computing the averages to the latest previous value and decreasing weights to the subsequent older ones.

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

```

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['Pr
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']))
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 using Weighted Moving Averages using previous Ratio values therefore we get  $R_t = (5 * R_{t-1} + 4 * R_{t-2} +$

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

```

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']))
    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 using Weighted Moving Averages using previous 2016 values therefore we get  $P_t = (2 * P_{t-1} + P_{t-2})/3$