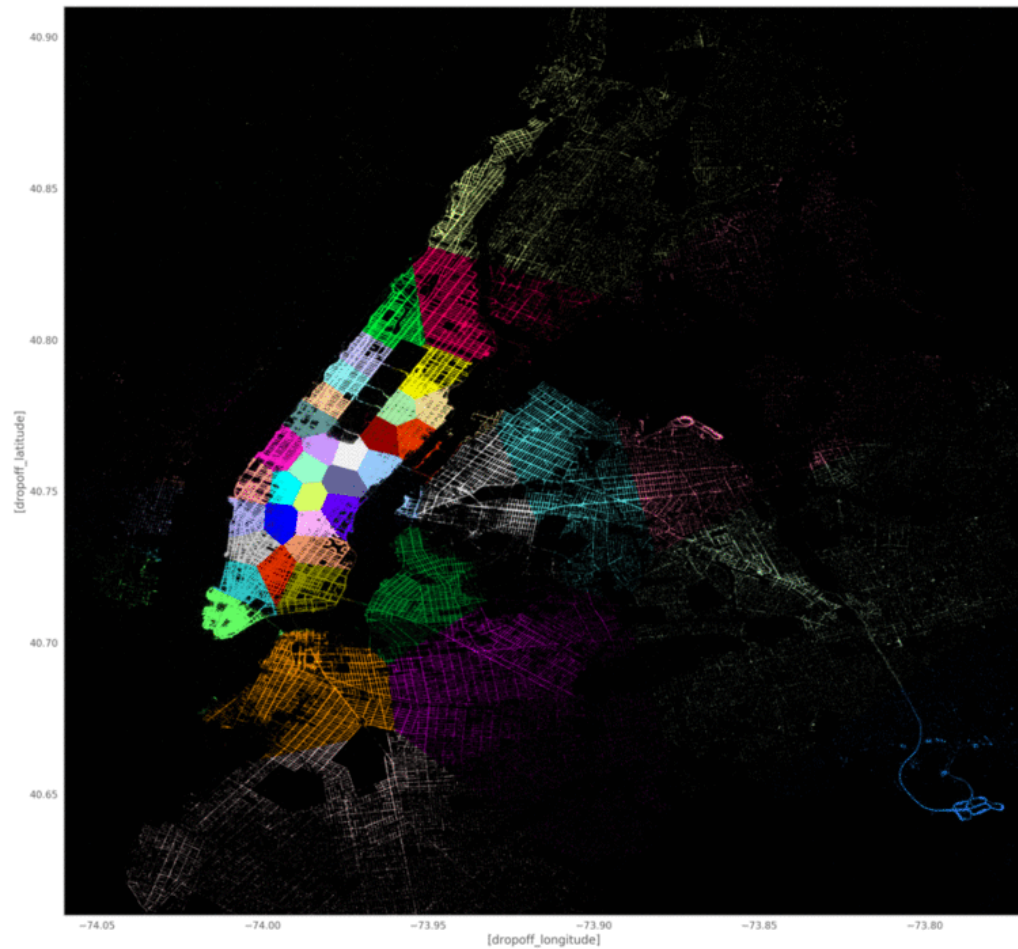


# Taxi demand prediction in New York City



```
In [2]: 1 !pip install gpxpy
```

Collecting gpxpy

Downloading <https://files.pythonhosted.org/packages/9e/96/d9fdd462f17ef54fd64fe3f5e870c88e35686c740dcaed6460d048b43566/gpxpy-1.4.0.tar.gz> (<https://files.pythonhosted.org/packages/9e/96/d9fdd462f17ef54fd64fe3f5e870c88e35686c740dcaed6460d048b43566/gpxpy-1.4.0.tar.gz>) (105kB)

|██| 112kB 2.8MB/s eta 0:00:01

Building wheels for collected packages: gpxpy

Building wheel for gpxpy (setup.py) ... done

Created wheel for gpxpy: filename=gpxpy-1.4.0-cp36-none-any.whl size=42813 sha256=5a911a384a97e99e02e7f24dfae7552e3e5b393f41d4f55ff596828ce6ff11c8

Stored in directory: /root/.cache/pip/wheels/77/d7/ee/cb4d7a151ce924c35e681377fb90a0b882f55bfd3c2c586739

Successfully built gpxpy

Installing collected packages: gpxpy

Successfully installed gpxpy-1.4.0

```

In [0]: 1 #Importing Libraries
2 # pip3 install graphviz
3 #pip3 install dask
4 #pip3 install toolz
5 #pip3 install cloudpickle
6 # https://www.youtube.com/watch?v=ieW3G7ZzRZ0
7 # https://github.com/dask/dask-tutorial
8 # please do go through this python notebook: https://github.com/dask/dask-tut
9 import dask.dataframe as dd#similar to pandas
10
11 import pandas as pd#pandas to create small dataframes
12
13 # pip3 install folium
14 # if this doesnt work refere install_folium.JPG in drive
15 import folium #open street map
16
17 # unix time: https://www.unixtimestamp.com/
18 import datetime #Convert to unix time
19
20 import time #Convert to unix time
21
22 # if numpy is not installed already : pip3 install numpy
23 import numpy as np#Do arithmetic operations on arrays
24
25 # matplotlib: used to plot graphs
26 import matplotlib
27 # matplotlib.use('nbagg') : matplotlib uses this protocol which makes plots
28 #matplotlib.use('nbagg')
29 import matplotlib.pyplot as plt
30 import seaborn as sns#Plots
31 from matplotlib import rcParams#Size of plots
32
33 # this lib is used while we calculate the stight line distance between two (L
34 import gpxpy.geo #Get the haversine distance
35
36 from sklearn.cluster import MiniBatchKMeans, KMeans#Clustering
37 import math
38 import pickle
39 import os
40
41 # download mingwin: https://mingw-w64.org/doku.php/download/mingw-builds
42 # install it in your system and keep the path, mingw_path = 'installed path'
43 mingw_path = 'C:\\Program Files\\mingw-w64\\x86_64-5.3.0-posix-seh-rt_v4-rev0
44 os.environ['PATH'] = mingw_path + ';' + os.environ['PATH']
45
46 # to install xgboost: pip3 install xgboost
47 # if it didnt happen check install_xgboost.JPG
48 import xgboost as xgb
49
50 # to install sklearn: pip install -U scikit-learn
51 from sklearn.ensemble import RandomForestRegressor
52 from sklearn.metrics import mean_squared_error
53 from sklearn.metrics import mean_absolute_error
54 import warnings
55 warnings.filterwarnings("ignore")

```

```
In [0]: 1 from google.colab import drive
        2 drive.mount('/gdrive')
```

Mounted at /gdrive

## 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)  
([http://www.nyc.gov/html/tlc/html/about/trip\\_record\\_data.shtml](http://www.nyc.gov/html/tlc/html/about/trip_record_data.shtml)) (2016 data) The data used in the attached datasets were collected and provided to the NYC Taxi and Limousine Commission (TLC)

## Information on taxis:

### ***Yellow Taxi: Yellow Medallion Taxicabs***

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

### ***For Hire Vehicles (FHVs)***

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

### ***Green Taxi: Street Hail Livery (SHL)***

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

Credits: Quora

### ***Footnote:***

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

## Data Collection

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

	file name	file name size	number of records	number of features
	yellow_tripdata_2016-01	1. 59G	10906858	19
	yellow_tripdata_2016-02	1. 66G	11382049	19

yellow_tripdata_2016-03	1. 78G	12210952	19
yellow_tripdata_2016-04	1. 74G	11934338	19
yellow_tripdata_2016-05	1. 73G	11836853	19
yellow_tripdata_2016-06	1. 62G	11135470	19
yellow_tripdata_2016-07	884Mb	10294080	17
yellow_tripdata_2016-08	854Mb	9942263	17
yellow_tripdata_2016-09	870Mb	10116018	17
yellow_tripdata_2016-10	933Mb	10854626	17
yellow_tripdata_2016-11	868Mb	10102128	17
yellow_tripdata_2016-12	897Mb	10449408	17
yellow_tripdata_2015-01	1.84Gb	12748986	19
yellow_tripdata_2015-02	1.81Gb	12450521	19
yellow_tripdata_2015-03	1.94Gb	13351609	19
yellow_tripdata_2015-04	1.90Gb	13071789	19
yellow_tripdata_2015-05	1.91Gb	13158262	19
yellow_tripdata_2015-06	1.79Gb	12324935	19
yellow_tripdata_2015-07	1.68Gb	11562783	19
yellow_tripdata_2015-08	1.62Gb	11130304	19
yellow_tripdata_2015-09	1.63Gb	11225063	19
yellow_tripdata_2015-10	1.79Gb	12315488	19
yellow_tripdata_2015-11	1.65Gb	11312676	19
yellow_tripdata_2015-12	1.67Gb	11460573	19

In [4]:

```

1 from google.colab import drive
2 drive.mount('/content/drive')
```

Go to this URL in a browser: [https://accounts.google.com/o/oauth2/auth?client\\_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect\\_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response\\_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly](https://accounts.google.com/o/oauth2/auth?client_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly) (http s://accounts.google.com/o/oauth2/auth?client\_id=947318989803-6bn6qk8qdgf4n4g3pfee6491hc0brc4i.apps.googleusercontent.com&redirect\_uri=urn%3aietf%3awg%3aoauth%3a2.0%3aob&response\_type=code&scope=email%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdocs.test%20https%3a%2f%2fwww.googleapis.com%2fauth%2fdrive.photos.readonly%20https%3a%2f%2fwww.googleapis.com%2fauth%2fpeopleapi.readonly)

Enter your authorization code:

.....

Mounted at /content/drive

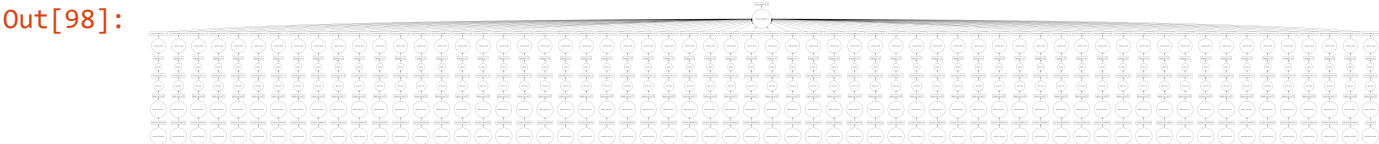
```
In [0]: 1 #Looking at the features
        2 # dask dataframe : # https://github.com/dask/dask-tutorial/blob/master/07_da
        3 month = dd.read_csv('drive/My Drive/NYTD/yellow_tripdata_2015-01.csv')
        4 print(month.columns)
```

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

```
In [0]: 1 # However unlike Pandas, operations on dask.dataframes don't trigger immediat
        2 # instead they add key-value pairs to an underlying Dask graph. Recall that i
        3 # circles are operations and rectangles are results.
        4
        5 # to see the visulaization you need to install graphviz
        6 # pip3 install graphviz if this doesnt work please check the install_graphviz
        7 month.visualize()
```



```
In [0]: 1 month.fare_amount.sum().visualize()
```



Features in the dataset:

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

		The final rate code in effect at the end of the trip.
	1.	Standard rate
	2.	JFK
RateCodeID	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.	
		A numeric code signifying how the passenger paid for the trip.
	1.	Credit card
	2.	Cash
Payment_type	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 [0]: 1 #table below shows few datapoints along with all our features
        2 start = datetime.datetime.now()
        3 print(month.head(5))
        4 print('time taken:',datetime.datetime.now() - start)
```

	VendorID	tpep_pickup_datetime	...	improvement_surcharge	total_amount
0	2	2015-01-15 19:05:39	...	0.3	17.05
1	1	2015-01-10 20:33:38	...	0.3	17.80
2	1	2015-01-10 20:33:38	...	0.3	10.80
3	1	2015-01-10 20:33:39	...	0.3	4.80
4	1	2015-01-10 20:33:39	...	0.3	16.30

[5 rows x 19 columns]

time taken: 0:00:03.256042

## 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 [0]: 1 # Plotting pickup coordinates which are outside the bounding box of New-York
2 # we will collect all the points outside the bounding box of newyork city to
3 outlier_locations = month[((month.pickup_longitude <= -74.15) | (month.pickup
4 (month.pickup_longitude >= -73.7004) | (month.pickup_latit
5 #outlier_locations.compute()
6
7 # creating a map with the a base location
8 # read more about the folium here: http://folium.readthedocs.io/en/latest/qui
9
10 # note: you dont need to remember any of these, you dont need indeepth knowle
11
12 map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')
13
14 # we will spot only first 100 outliers on the map, plotting all the outliers
15 sample_locations = outlier_locations.head(10000)
16 for i,j in sample_locations.iterrows():
17     if int(j['pickup_latitude']) != 0:
18         folium.Marker(list((j['pickup_latitude'],j['pickup_longitude']))).add
19 map_osm

```

Out[10]:

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

## 2. Dropoff Latitude & Dropoff Longitude

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

```

In [0]: 1 # Plotting dropoff coordinates which are outside the bounding box of New-York
2 # we will collect all the points outside the bounding box of newyork city to
3 outlier_locations = month[((month.dropoff_longitude <= -74.15) | (month.dropo
4 (month.dropoff_longitude >= -73.7004) | (month.dropoff_lat
5
6 # creating a map with the a base location
7 # read more about the folium here: http://folium.readthedocs.io/en/latest/qui
8
9 # note: you dont need to remember any of these, you dont need indepth knowle
10
11 map_osm = folium.Map(location=[40.734695, -73.990372], tiles='Stamen Toner')
12
13 # we will spot only first 100 outliers on the map, plotting all the outliers
14 sample_locations = outlier_locations.head(10000)
15 for i,j in sample_locations.iterrows():
16     if int(j['pickup_latitude']) != 0:
17         folium.Marker(list((j['dropoff_latitude'],j['dropoff_longitude']))).a
18 map_osm

```

Out[11]:

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

### 3. Trip Durations:

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

```

In [0]: 1 #The timestamps are converted to unix so as to get duration(trip-time) & speed
2
3 # in out data we have time in the formate "YYYY-MM-DD HH:MM:SS" we convert th
4 # https://stackoverflow.com/a/27914405
5 def convert_to_unix(s):
6     return time.mktime(datetime.datetime.strptime(s, "%Y-%m-%d %H:%M:%S").tim
7
8
9
10 # we return a data frame which contains the columns
11 # 1.'passenger_count' : self explanatory
12 # 2.'trip_distance' : self explanatory
13 # 3.'pickup_longitude' : self explanatory
14 # 4.'pickup_latitude' : self explanatory
15 # 5.'dropoff_longitude' : self explanatory
16 # 6.'dropoff_latitude' : self explanatory
17 # 7.'total_amount' : total fair that was paid
18 # 8.'trip_times' : duration of each trip
19 # 9.'pickup_times' : pickup time converted into unix time
20 # 10.'Speed' : velocity of each trip
21 def return_with_trip_times(month):
22     duration = month[['tpep_pickup_datetime', 'tpep_dropoff_datetime']].compute
23     #pickups and dropoffs to unix time
24     duration_pickup = [convert_to_unix(x) for x in duration['tpep_pickup_date
25     duration_drop = [convert_to_unix(x) for x in duration['tpep_dropoff_datet
26     #calculate duration of trips
27     durations = (np.array(duration_drop) - np.array(duration_pickup))/float(6
28
29     #append durations of trips and speed in miles/hr to a new dataframe
30     new_frame = month[['passenger_count', 'trip_distance', 'pickup_longitude', '
31
32     new_frame['trip_times'] = durations
33     new_frame['pickup_times'] = duration_pickup
34     new_frame['Speed'] = 60*(new_frame['trip_distance']/new_frame['trip_times
35
36     return new_frame
37
38 # print(frame_with_durations.head())
39 # passenger_count trip_distance pickup_longitude pickup_latitude dropo
40 # 1 1.59 -73.993896 40.750111 -73.9
41 # 1 3.30 -74.001648 40.724243 -73.9
42 # 1 1.80 -73.963341 40.802788 -73.9
43 # 1 0.50 -74.009087 40.713818 -74.0
44 # 1 3.00 -73.971176 40.762428 -74.0
45 frame_with_durations = return_with_trip_times(month)

```

```

In [0]: 1 # the skewed box plot shows us the presence of outliers
2 sns.boxplot(y="trip_times", data =frame_with_durations)
3 plt.show()

```

<IPython.core.display.Javascript object>

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

```
0 percentile value is -1211.0166666666667
10 percentile value is 3.8333333333333335
20 percentile value is 5.383333333333334
30 percentile value is 6.816666666666666
40 percentile value is 8.3
50 percentile value is 9.95
60 percentile value is 11.866666666666667
70 percentile value is 14.283333333333333
80 percentile value is 17.633333333333333
90 percentile value is 23.45
100 percentile value is 548555.633333
```

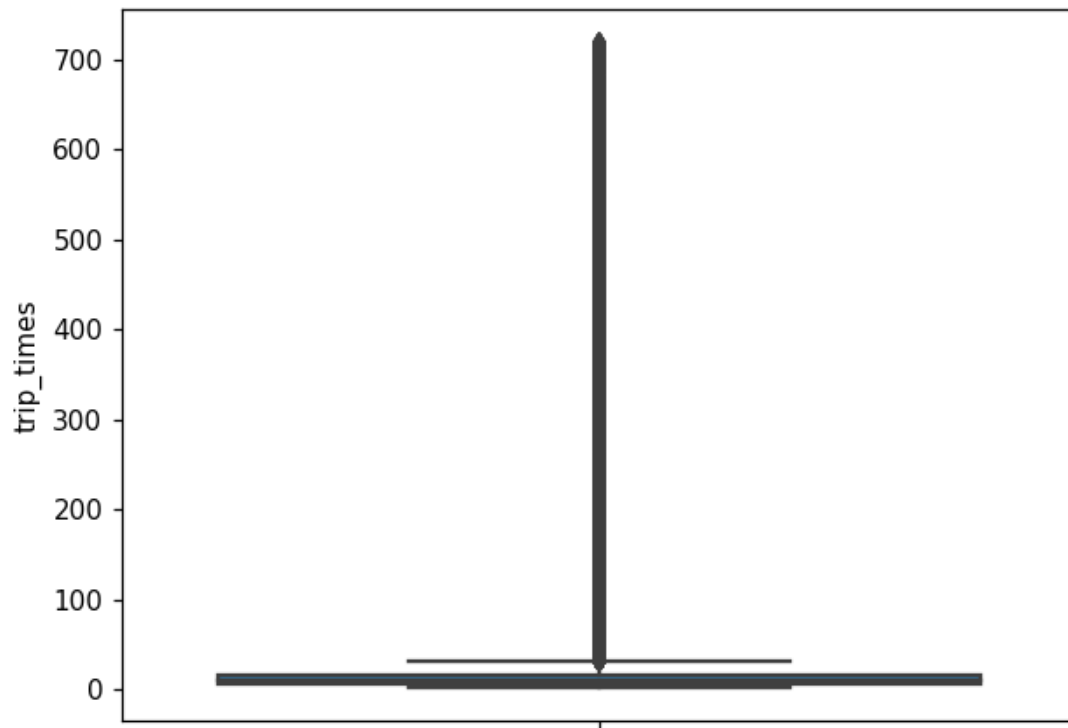
```
In [0]: 1 #Looking further from the 99th percetntile
2 for i in range(90,100):
3     var =frame_with_durations["trip_times"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100)
6     print ("100 percentile value is ",var[-1]))
```

```
90 percentile value is 23.45
91 percentile value is 24.35
92 percentile value is 25.383333333333333
93 percentile value is 26.55
94 percentile value is 27.933333333333334
95 percentile value is 29.583333333333332
96 percentile value is 31.683333333333334
97 percentile value is 34.466666666666667
98 percentile value is 38.716666666666667
99 percentile value is 46.75
100 percentile value is 548555.633333
```

```
In [0]: 1 #removing data based on our analysis and TLC regulations
2 frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip
```

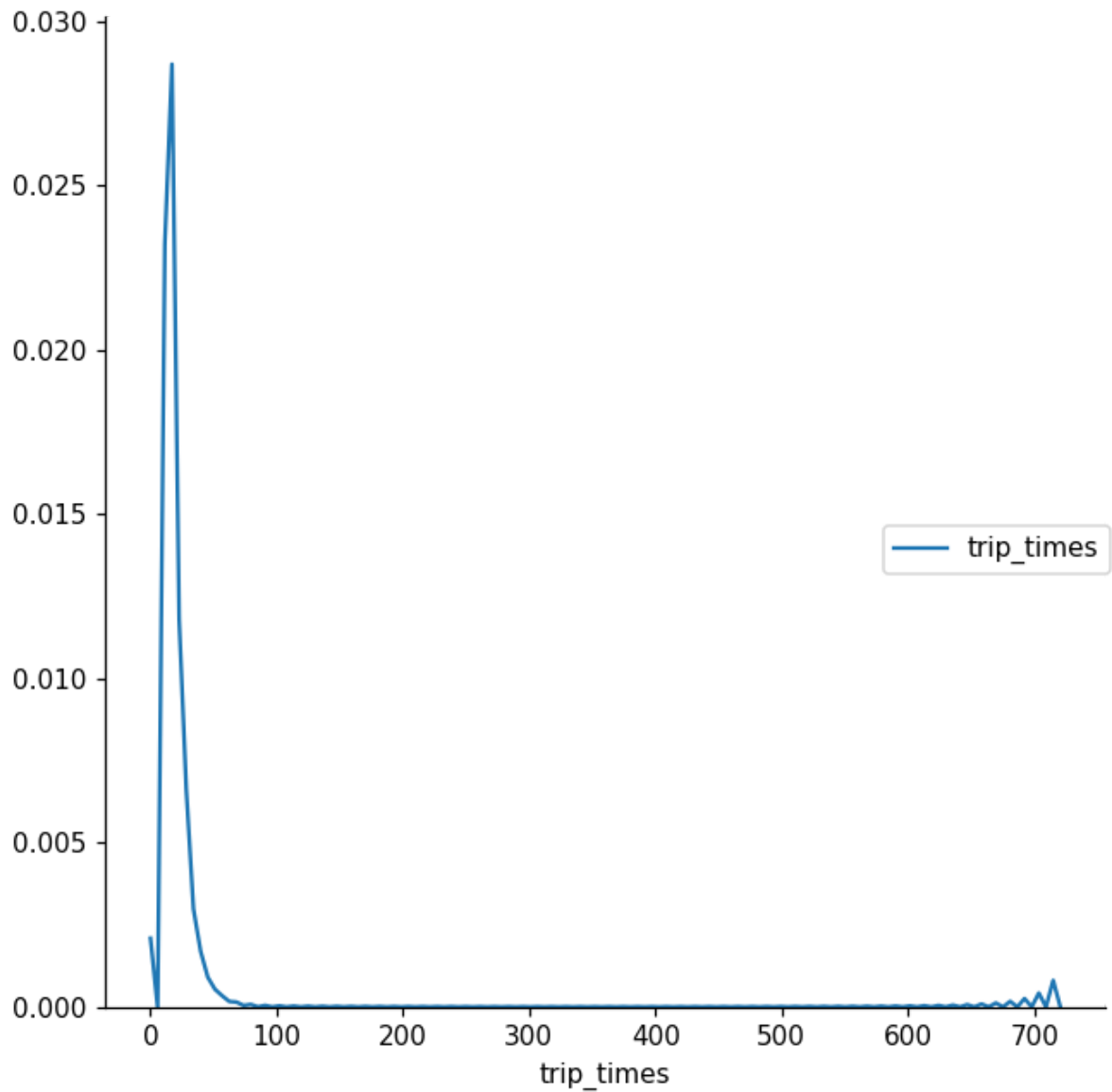
```
In [0]: 1 #box-plot after removal of outliers  
2 sns.boxplot(y="trip_times", data =frame_with_durations_modified)  
3 plt.show()
```

<IPython.core.display.Javascript object>



```
In [0]: 1 #pdf of trip-times after removing the outliers
2 sns.FacetGrid(frame_with_durations_modified,size=6) \
3     .map(sns.kdeplot,"trip_times") \
4     .add_legend();
5 plt.show();
```

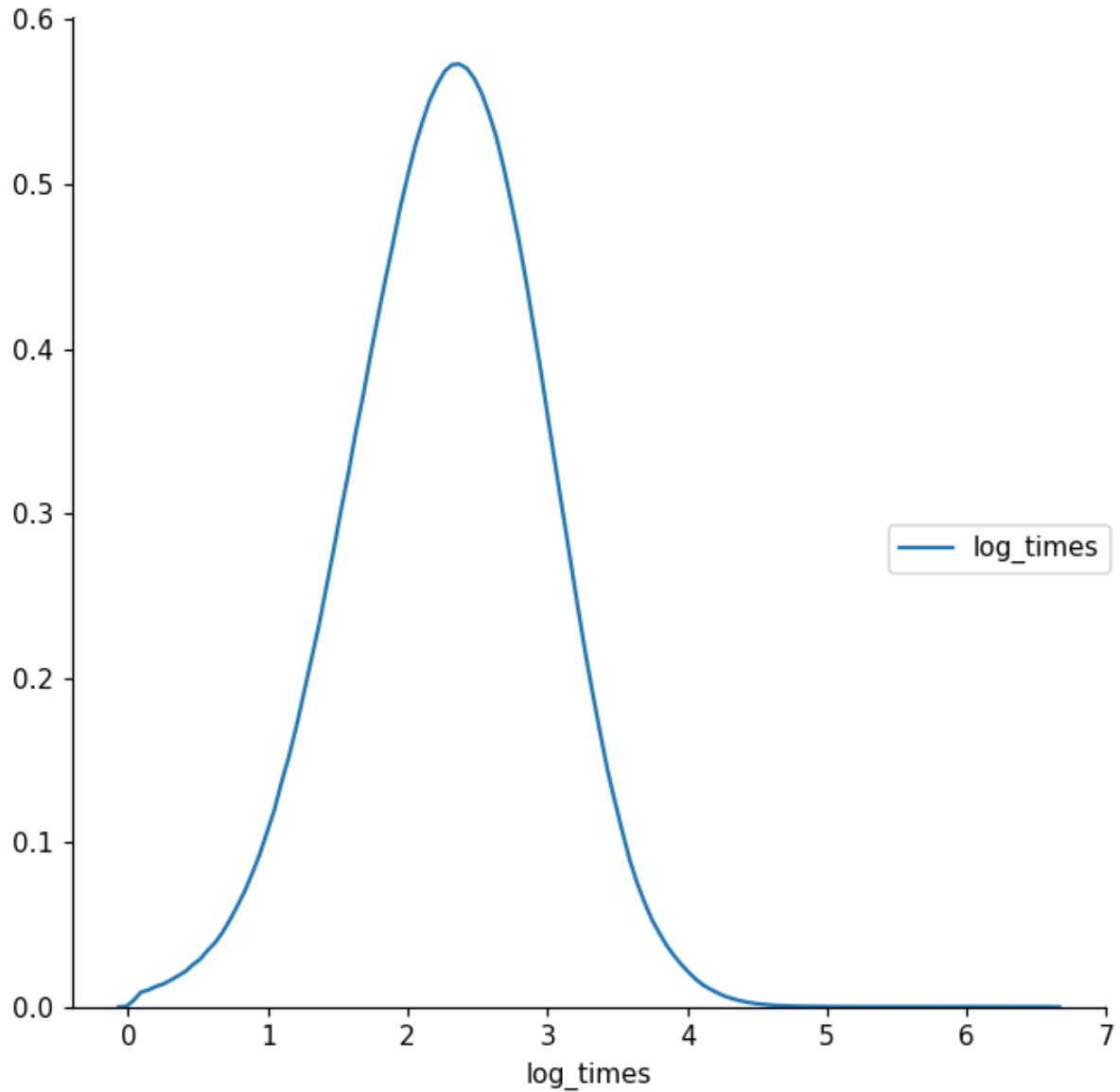
<IPython.core.display.Javascript object>



```
In [0]: 1 #converting the values to log-values to chec for log-normal
        2 import math
        3 frame_with_durations_modified['log_times']=[math.log(i) for i in frame_with_d
```

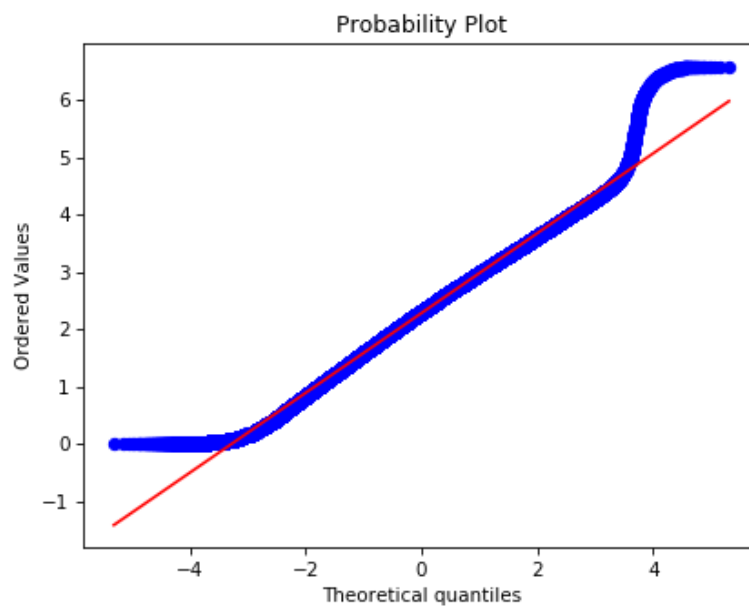
```
In [0]: 1 #pdf of log-values
        2 sns.FacetGrid(frame_with_durations_modified,size=6) \
        3     .map(sns.kdeplot,"log_times") \
        4     .add_legend();
        5 plt.show();
```

<IPython.core.display.Javascript object>



```
In [0]: 1 #Q-Q plot for checking if trip-times is log-normal  
2 scipy.stats.probplot(frame_with_durations_modified['log_times'].values, plot=  
3 plt.show())
```

<IPython.core.display.Javascript object>

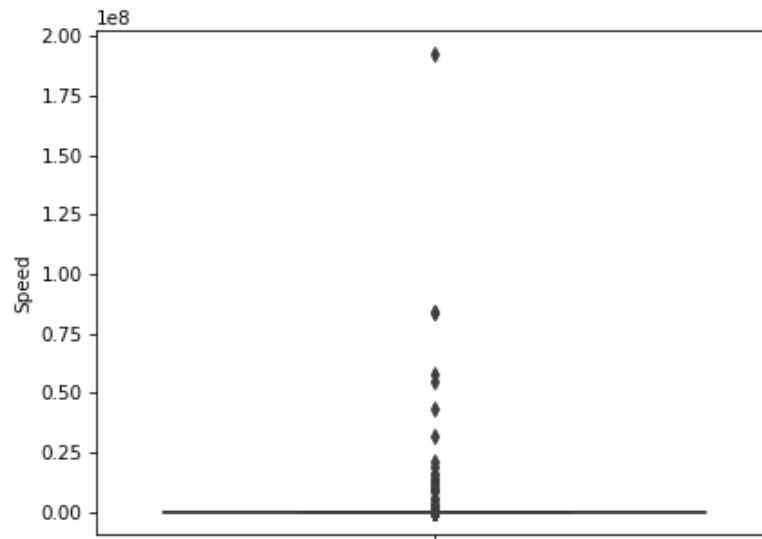


## 4. Speed



```
In [0]: 1 # check for any outliers in the data after trip duration outliers removed  
2 # box-plot for speeds with outliers  
3 frame_with_durations_modified['Speed'] = 60*(frame_with_durations_modified['t  
4 sns.boxplot(y="Speed", data =frame_with_durations_modified)  
5 plt.show()
```

<IPython.core.display.Javascript object>



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

```
0 percentile value is 0.0
10 percentile value is 6.409495548961425
20 percentile value is 7.80952380952381
30 percentile value is 8.929133858267717
40 percentile value is 9.98019801980198
50 percentile value is 11.06865671641791
60 percentile value is 12.286689419795222
70 percentile value is 13.796407185628745
80 percentile value is 15.963224893917962
90 percentile value is 20.186915887850468
100 percentile value is 192857142.857
```

```
In [0]: 1 #calculating speed values at each percentile 90,91,92,93,94,95,96,97,98,99,100
2 for i in range(90,100):
3     var =frame_with_durations_modified["Speed"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100)]))
6     print("100 percentile value is ",var[-1])
```

```
90 percentile value is 20.186915887850468
91 percentile value is 20.91645569620253
92 percentile value is 21.752988047808763
93 percentile value is 22.721893491124263
94 percentile value is 23.844155844155843
95 percentile value is 25.182552504038775
96 percentile value is 26.80851063829787
97 percentile value is 28.84304932735426
98 percentile value is 31.591128254580514
99 percentile value is 35.7513566847558
100 percentile value is 192857142.857
```

```
In [0]: 1 #calculating speed values at each percentile 99.0,99.1,99.2,99.3,99.4,99.5,99.
2 for i in np.arange(0.0, 1.0, 0.1):
3     var =frame_with_durations_modified["Speed"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+
6     print("100 percentile value is ",var[-1]))
```

```
99.0 percentile value is 35.7513566847558
99.1 percentile value is 36.31084727468969
99.2 percentile value is 36.91470054446461
99.3 percentile value is 37.588235294117645
99.4 percentile value is 38.33035714285714
99.5 percentile value is 39.17580340264651
99.6 percentile value is 40.15384615384615
99.7 percentile value is 41.338301043219076
99.8 percentile value is 42.86631016042781
99.9 percentile value is 45.3107822410148
100 percentile value is 192857142.857
```

```
In [0]: 1 #removing further outliers based on the 99.9th percentile value
2 frame_with_durations_modified=frame_with_durations[(frame_with_durations.Speed
```

```
In [0]: 1 #avg.speed of cabs in New-York
2 sum(frame_with_durations_modified["Speed"]) / float(len(frame_with_durations_
```

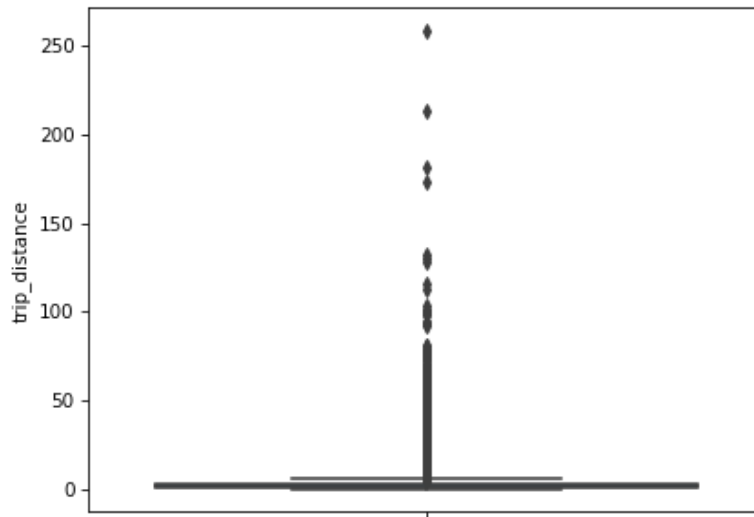
```
Out[17]: 12.450173996027528
```

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 [0]: 1 # up to now we have removed the outliers based on trip durations and cab speed
        2 # Lets try if there are any outliers in trip distances
        3 # box-plot showing outliers in trip-distance values
        4 sns.boxplot(y="trip_distance", data =frame_with_durations_modified)
        5 plt.show()
```

<IPython.core.display.Javascript object>



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

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

```
In [0]: 1 #calculating trip distance values at each percentile 90,91,92,93,94,95,96,97,98,99,100
2 for i in range(90,100):
3     var =frame_with_durations_modified["trip_distance"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100)]))
6     print("100 percentile value is ",var[-1])
```

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

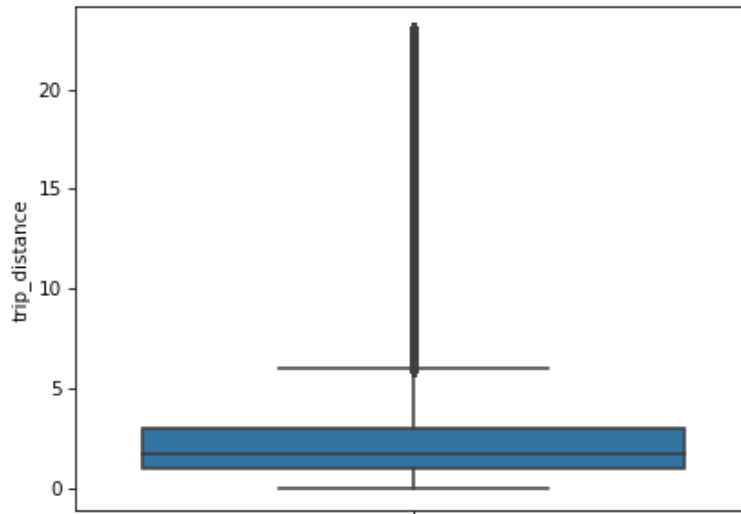
```
In [0]: 1 #calculating trip distance values at each percentile 99.0,99.1,99.2,99.3,99.4,99.5,99.6,99.7,99.8,99.9,100
2 for i in np.arange(0.0, 1.0, 0.1):
3     var =frame_with_durations_modified["trip_distance"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/100)]))
6     print("100 percentile value is ",var[-1])
```

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

```
In [0]: 1 #removing further outliers based on the 99.9th percentile value
2 frame_with_durations_modified=frame_with_durations[(frame_with_durations.trip
```

```
In [0]: 1 #box-plot after removal of outliers  
2 sns.boxplot(y="trip_distance", data = frame_with_durations_modified)  
3 plt.show()
```

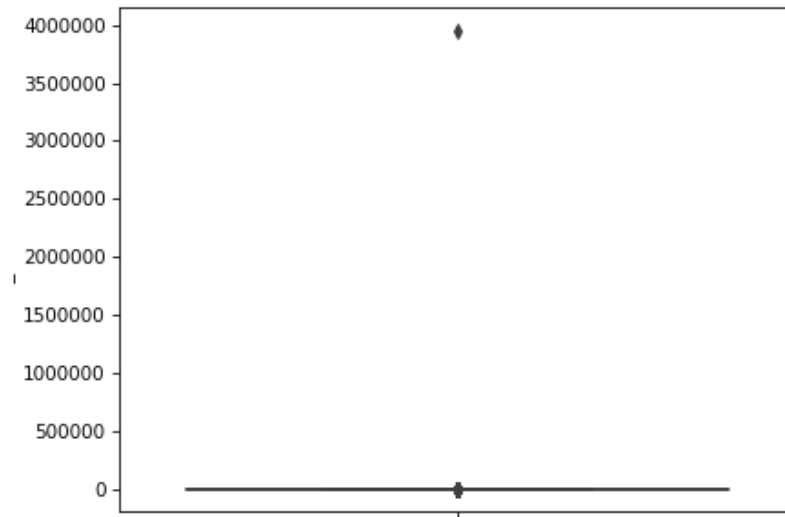
<IPython.core.display.Javascript object>



## 5. Total Fare

```
In [0]: 1 # up to now we have removed the outliers based on trip durations, cab speeds,  
2 # Lets try if there are any outliers in based on the total_amount  
3 # box-plot showing outliers in fare  
4 sns.boxplot(y="total_amount", data =frame_with_durations_modified)  
5 plt.show()
```

<IPython.core.display.Javascript object>



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

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

```
In [0]: 1 #calculating total fare amount values at each percentile 90,91,92,93,94,95,96,97,98,99,100
2 for i in range(90,100):
3     var = frame_with_durations_modified["total_amount"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/100)]))
6     print("100 percentile value is ",var[-1])
```

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



```
In [0]: 1 #calculating total fare amount values at each percentile 99.0,99.1,99.2,99.3,9
2 for i in np.arange(0.0, 1.0, 0.1):
3     var = frame_with_durations_modified["total_amount"].values
4     var = np.sort(var,axis = None)
5     print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+
6     print("100 percentile value is ",var[-1])
```

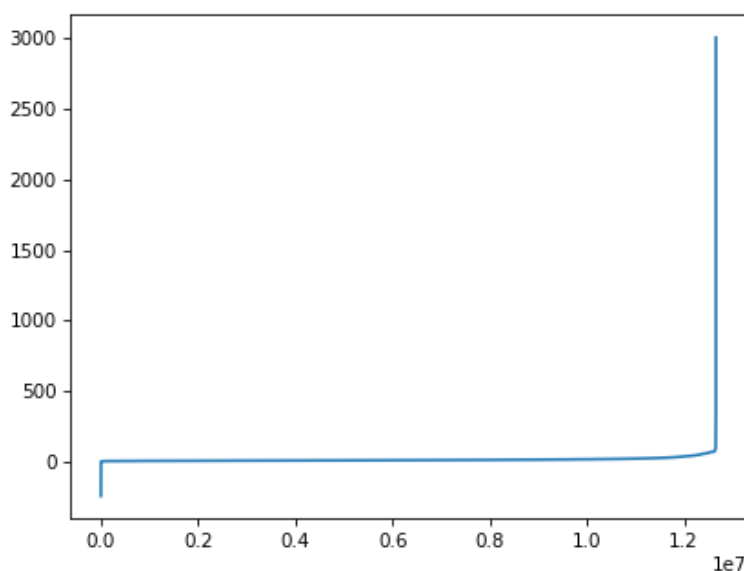
```
99.0 percentile value is 68.13
99.1 percentile value is 69.13
99.2 percentile value is 69.6
99.3 percentile value is 69.73
99.4 percentile value is 69.73
99.5 percentile value is 69.76
99.6 percentile value is 72.46
99.7 percentile value is 72.73
99.8 percentile value is 80.05
99.9 percentile value is 95.55
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 [0]: 1 del frame_with_durations_modified
```

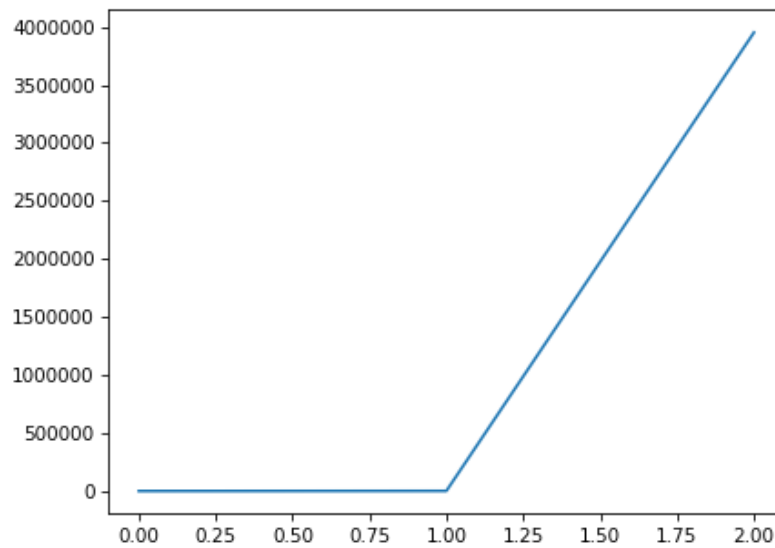
```
In [0]: 1 #below plot shows us the fare values(sorted) to find a sharp increase to remo
2 # plot the fare amount excluding last two values in sorted data
3 plt.plot(var[:-2])
4 plt.show()
```

<IPython.core.display.Javascript object>



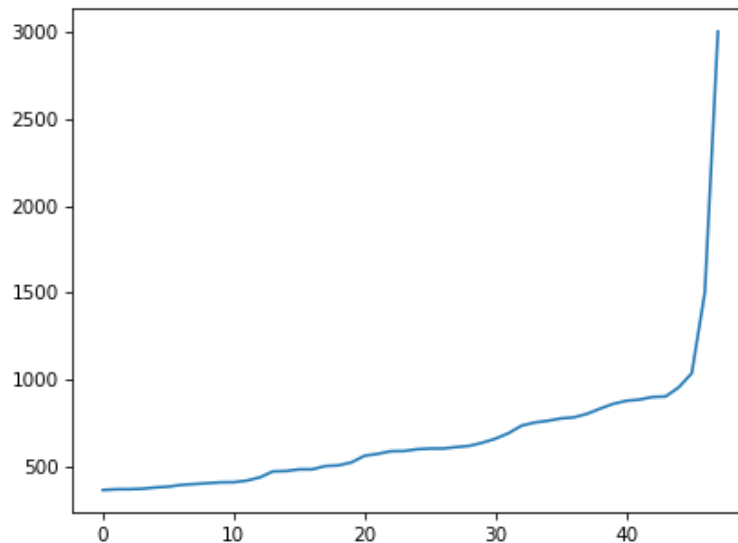
```
In [0]: 1 # a very sharp increase in fare values can be seen  
2 # plotting last three total fare values, and we can observe there is share in  
3 plt.plot(var[-3:])  
4 plt.show()
```

<IPython.core.display.Javascript object>



```
In [0]: 1 #now looking at values not including the last two points we again find a dras  
2 # we plot last 50 values excluding last two values  
3 plt.plot(var[-50:-2])  
4 plt.show()
```

<IPython.core.display.Javascript object>



**Remove all outliers/erronous points.**

```

In [0]: 1 #removing all outliers based on our univariate analysis above
2 def remove_outliers(new_frame):
3
4
5     a = new_frame.shape[0]
6     print ("Number of pickup records = ",a)
7     temp_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_fr
8                             (new_frame.dropoff_latitude >= 40.5774) & (new_frame.d
9                             ((new_frame.pickup_longitude >= -74.15) & (new_frame.p
10                            (new_frame.pickup_longitude <= -73.7004) & (new_frame.
11     b = temp_frame.shape[0]
12     print ("Number of outlier coordinates lying outside NY boundaries:",(a-b))
13
14
15     temp_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times
16     c = temp_frame.shape[0]
17     print ("Number of outliers from trip times analysis:",(a-c))
18
19
20     temp_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_di
21     d = temp_frame.shape[0]
22     print ("Number of outliers from trip distance analysis:",(a-d))
23
24     temp_frame = new_frame[(new_frame.Speed <= 65) & (new_frame.Speed >= 0)]
25     e = temp_frame.shape[0]
26     print ("Number of outliers from speed analysis:",(a-e))
27
28     temp_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_
29     f = temp_frame.shape[0]
30     print ("Number of outliers from fare analysis:",(a-f))
31
32
33     new_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) & (new_fr
34                            (new_frame.dropoff_latitude >= 40.5774) & (new_frame.d
35                            ((new_frame.pickup_longitude >= -74.15) & (new_frame.p
36                            (new_frame.pickup_longitude <= -73.7004) & (new_frame.
37
38     new_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times
39     new_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.trip_dis
40     new_frame = new_frame[(new_frame.Speed < 45.31) & (new_frame.Speed > 0)]
41     new_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.total_a
42
43     print ("Total outliers removed",a - new_frame.shape[0])
44     print ("---")
45     return new_frame

```

```
In [0]: 1 print ("Removing outliers in the month of Jan-2015")
        2 print ("----")
        3 frame_with_durations_outliers_removed = remove_outliers(frame_with_durations)
        4 print("fraction of data points that remain after removing outliers", float(le
```

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 [0]: 1 #trying different cluster sizes to choose the right K in K-means
2 coords = frame_with_durations_outliers_removed[['pickup_latitude', 'pickup_lo
3 neighbours=[]
4
5 def find_min_distance(cluster_centers, cluster_len):
6     nice_points = 0
7     wrong_points = 0
8     less2 = []
9     more2 = []
10    min_dist=1000
11    for i in range(0, cluster_len):
12        nice_points = 0
13        wrong_points = 0
14        for j in range(0, cluster_len):
15            if j!=i:
16                distance = gpxpy.geo.haversine_distance(cluster_centers[i][0]
17                min_dist = min(min_dist,distance/(1.60934*1000))
18                if (distance/(1.60934*1000)) <= 2:
19                    nice_points +=1
20                else:
21                    wrong_points += 1
22            less2.append(nice_points)
23            more2.append(wrong_points)
24    neighbours.append(less2)
25    print ("On choosing a cluster size of ",cluster_len,"\nAvg. Number of Clu
26
27 def find_clusters(increment):
28     kmeans = MiniBatchKMeans(n_clusters=increment, batch_size=10000,random_st
29     frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(
30     cluster_centers = kmeans.cluster_centers_
31     cluster_len = len(cluster_centers)
32     return cluster_centers, cluster_len
33
34 # we need to choose number of clusters so that, there are more number of clus
35 #that are close to any cluster center
36 # and make sure that the minimum inter cluster should not be very less
37 for increment in range(10, 100, 10):
38     cluster_centers, cluster_len = find_clusters(increment)
39     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 = 1.0933194607372518
---
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 = 0.7123318236197774
---
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 = 0.5179286172497254
---
On choosing a cluster size of 40
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2):
9.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):
31.0
Min inter-cluster distance = 0.5064095487015858
---
On choosing a cluster size of 50
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 1
2.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):
38.0
Min inter-cluster distance = 0.36495419250817024
---
On choosing a cluster size of 60
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 1
4.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):
46.0
Min inter-cluster distance = 0.346654501371586
---
On choosing a cluster size of 70
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 1
6.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):
54.0
Min inter-cluster distance = 0.30468071844965394
---
On choosing a cluster size of 80
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 1
8.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):
62.0
Min inter-cluster distance = 0.29187627608454664
---
On choosing a cluster size of 90
Avg. Number of Clusters within the vicinity (i.e. intercluster-distance < 2): 2
1.0
Avg. Number of Clusters outside the vicinity (i.e. intercluster-distance > 2):
69.0
Min inter-cluster distance = 0.18237562550345013
---
```

## 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 [0]: 1 coords = jan_2015_frame[['pickup_latitude', 'pickup_longitude']].values
```

```
In [0]: 1 # if check for the 50 clusters you can observe that there are two clusters wi  
2 # so we choose 40 clusters for solve the further problem  
3  
4 # Getting 40 clusters using the kmeans  
5 kmeans = MiniBatchKMeans(n_clusters=40, batch_size=10000, random_state=0).fit(  
6 #frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(fra
```

## Plotting the cluster centers:

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

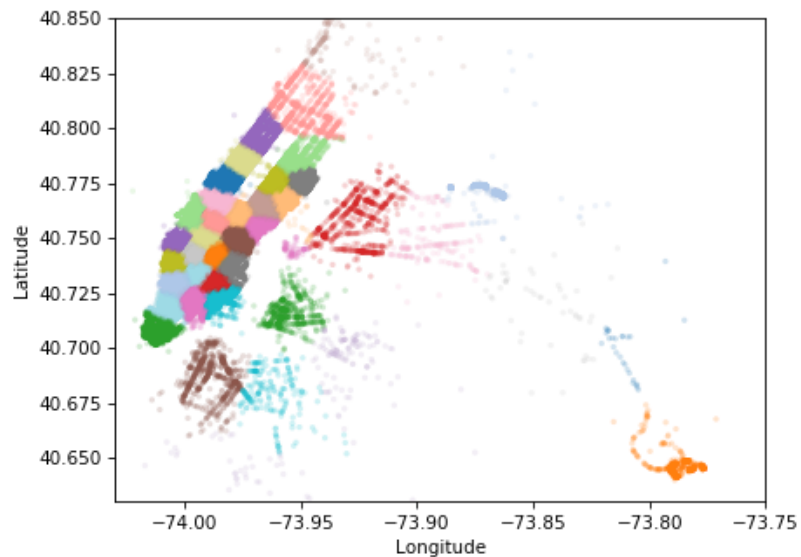
Out[48]:

## Plotting the clusters:



```
In [0]: 1 #Visualising the clusters on a map
2 def plot_clusters(frame):
3     city_long_border = (-74.03, -73.75)
4     city_lat_border = (40.63, 40.85)
5     fig, ax = plt.subplots(ncols=1, nrows=1)
6     ax.scatter(frame.pickup_longitude.values[:100000], frame.pickup_latitude.
7               c=frame.pickup_cluster.values[:100000], cmap='tab20', alpha=0.
8               ax.set_xlim(city_long_border)
9               ax.set_ylim(city_lat_border)
10              ax.set_xlabel('Longitude')
11              ax.set_ylabel('Latitude')
12              plt.show()
13
14 plot_clusters(frame_with_durations_outliers_removed)
```

<IPython.core.display.Javascript object>



## Time-binning

```

In [0]: 1 #Refer:https://www.unixtimestamp.com/
2 # 1420070400 : 2015-01-01 00:00:00
3 # 1422748800 : 2015-02-01 00:00:00
4 # 1425168000 : 2015-03-01 00:00:00
5 # 1427846400 : 2015-04-01 00:00:00
6 # 1430438400 : 2015-05-01 00:00:00
7 # 1433116800 : 2015-06-01 00:00:00
8
9 # 1451606400 : 2016-01-01 00:00:00
10 # 1454284800 : 2016-02-01 00:00:00
11 # 1456790400 : 2016-03-01 00:00:00
12 # 1459468800 : 2016-04-01 00:00:00
13 # 1462060800 : 2016-05-01 00:00:00
14 # 1464739200 : 2016-06-01 00:00:00
15
16 def add_pickup_bins(frame,month,year):
17     unix_pickup_times=[i for i in frame['pickup_times'].values]
18     unix_times = [[1420070400,1422748800,1425168000,1427846400,1430438400,143
19                    [1451606400,1454284800,1456790400,1459468800,1462060800,1
20
21     start_pickup_unix=unix_times[year-2015][month-1]
22     # https://www.timeanddate.com/time/zones/est
23     # (int((i-start_pickup_unix)/600)+33): our unix time is in gmt to we are
24     tenminutewise_binned_unix_pickup_times=[(int((i-start_pickup_unix)/600)+3
25     frame['pickup_bins'] = np.array(tenminutewise_binned_unix_pickup_times)
26     return frame

```

```

In [0]: 1 # clustering, making pickup bins and grouping by pickup cluster and pickup bi
2 frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(fram
3 jan_2015_frame = add_pickup_bins(frame_with_durations_outliers_removed,1,2015
4 jan_2015_groupby = jan_2015_frame[['pickup_cluster','pickup_bins','trip_dista

```

```

In [0]: 1 del jan_2015_frame
2 del jan_2015_groupby

```

```

In [0]: 1 jan_2015_frame.to_csv('jan_2015_frame.csv',index = False)
2 jan_2015_groupby.to_csv('jan_2015_groupby.csv',index = False)

```

```

In [0]: 1 #del frame_with_durations
2 del frame_with_durations_outliers_removed
3 del frame_with_durations

```

```
In [0]: 1 # we add two more columns 'pickup_cluster'(to which cluster it belongs to)
        2 # and 'pickup_bins' (to which 10min intravel the trip belongs to)
        3 jan_2015_frame.head()
```

Out[14]:

	passenger_count	trip_distance	pickup_longitude	pickup_latitude	dropoff_longitude	dropoff_latitude
0	1	1.59	-73.993896	40.750111	-73.974785	40.750111
1	1	3.30	-74.001648	40.724243	-73.994415	40.750111
2	1	1.80	-73.963341	40.802788	-73.951820	40.824243
3	1	0.50	-74.009087	40.713818	-74.004326	40.713818
4	1	3.00	-73.971176	40.762428	-74.004181	40.742428



```
In [0]: 1 # hear the trip_distance represents the number of pickups that are happend in
        2 # this data frame has two indices
        3 # primary index: pickup_cluster (cluster number)
        4 # secondary index : pickup_bins (we devid whole months time into 10min intravel)
        5 jan_2015_groupby.head()
```

Out[31]:

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

```

In [0]: 1 # upto now we cleaned data and prepared data for the month 2015,
2
3 # now do the same operations for months Jan, Feb, March of 2016
4 # 1. get the dataframe which includes only required columns
5 # 2. adding trip times, speed, unix time stamp of pickup_time
6 # 4. remove the outliers based on trip_times, speed, trip_duration, total_amo
7 # 5. add pickup_cluster to each data point
8 # 6. add pickup_bin (index of 10min intravel to which that trip belongs to)
9 # 7. group by data, based on 'pickup_cluster' and 'pickup_bin'
10
11 # Data Preparation for the months of Jan, Feb and March 2016
12 def datapreparation(month, kmeans, month_no, year_no):
13
14     print ("Return with trip times..")
15
16     frame_with_durations = return_with_trip_times(month)
17
18     print ("Remove outliers..")
19     frame_with_durations_outliers_removed = remove_outliers(frame_with_durati
20
21     print ("Estimating clusters..")
22     frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.predict(
23     #frame_with_durations_outliers_removed_2016['pickup_cluster'] = kmeans.pr
24
25     print ("Final groupbying..")
26     final_updated_frame = add_pickup_bins(frame_with_durations_outliers_remov
27     final_groupby_frame = final_updated_frame[['pickup_cluster', 'pickup_bins'
28
29     return final_updated_frame, final_groupby_frame
30
31 month_jan_2016 = dd.read_csv('drive/My Drive/NYTD/yellow_tripdata_2016-01.csv
32 month_feb_2016 = dd.read_csv('drive/My Drive/NYTD/yellow_tripdata_2016-02.csv
33 month_mar_2016 = dd.read_csv('drive/My Drive/NYTD/yellow_tripdata_2016-03.csv

```

Type Markdown and LaTeX:  $\alpha^2$

```
In [0]: 1 jan_2016_frame,jan_2016_groupby = datapreparation(month_jan_2016,kmeans,1,201
2 feb_2016_frame,feb_2016_groupby = datapreparation(month_feb_2016,kmeans,2,201
3 mar_2016_frame,mar_2016_groupby = datapreparation(month_mar_2016,kmeans,3,201
```

```
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..
```

```
In [0]: 1 jan_2016_frame.to_csv('jan_2016_frame.csv',index = False)
2 jan_2016_groupby.to_csv('jan_2016_groupby.csv',index = False)
3 feb_2016_frame.to_csv('feb_2016_frame.csv',index = False)
4 feb_2016_groupby.to_csv('feb_2016_groupby.csv',index = False)
5 mar_2016_frame.to_csv('march_2016_frame.csv',index = False)
6 mar_2016_groupby.to_csv('march_2016_groupby.csv',index = False)
```

Type *Markdown* and LaTeX:  $\alpha^2$

```
In [0]: 1 del month_jan_2016
2 del month_feb_2016
3 del month_mar_2016
```

## Smoothing

```
In [0]: 1 jan_2015_frame = pd.read_csv('drive/My Drive/NYTD/jan_2015_frame.csv')
2 feb_2016_frame = pd.read_csv('drive/My Drive/NYTD/feb_2016_frame.csv')
3 march_2016_frame = pd.read_csv('drive/My Drive/NYTD/march_2016_frame.csv')
4 jan_2016_frame = pd.read_csv('drive/My Drive/NYTD/jan_2016_frame.csv')
5 jan_2015_groupby = pd.read_csv('drive/My Drive/NYTD/jan_2015_groupby.csv')
6 jan_2016_groupby = pd.read_csv('drive/My Drive/NYTD/jan_2016_groupby.csv')
7 feb_2016_groupby = pd.read_csv('drive/My Drive/NYTD/feb_2016_groupby.csv')
8 march_2016_groupby = pd.read_csv('drive/My Drive/NYTD/march_2016_groupby.csv')
```

```
In [0]: 1 # Gets the unique bins where pickup values are present for each each reigion
2
3 # for each cluster region we will collect all the indices of 10min intravels
4 # we got an observation that there are some pickpbins that doesnt have any pi
5 def return_unq_pickup_bins(frame):
6     values = []
7     for i in range(0,40):
8         new = frame[frame['pickup_cluster'] == i]
9         list_unq = list(set(new['pickup_bins']))
10        list_unq.sort()
11        values.append(list_unq)
12    return values
```

```
In [0]: 1 # for every month we get all indices of 10min intravels in which atleast one
2
3 #jan
4 jan_2015_unique = return_unq_pickup_bins(jan_2015_frame)
5 jan_2016_unique = return_unq_pickup_bins(jan_2016_frame)
6
7 #feb
8 feb_2016_unique = return_unq_pickup_bins(feb_2016_frame)
```

```
In [0]: 1 #march
2 mar_2016_unique = return_unq_pickup_bins(march_2016_frame)
```

```
In [11]: 1 # for each cluster number of 10min intravels with 0 pickups
2 for i in range(40):
3     print("for the ",i,"th cluster number of 10min intavels with zero pickups")
4     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 intervals with zero pickups: 31
-----
for the 27 th cluster number of 10min intervals with zero pickups: 214
-----
for the 28 th cluster number of 10min intervals with zero pickups: 36
-----
for the 29 th cluster number of 10min intervals with zero pickups: 41
-----
for the 30 th cluster number of 10min intervals with zero pickups: 1180
-----
for the 31 th cluster number of 10min intervals with zero pickups: 42
-----
for the 32 th cluster number of 10min intervals with zero pickups: 44
-----
for the 33 th cluster number of 10min intervals with zero pickups: 43
-----
for the 34 th cluster number of 10min intervals with zero pickups: 39
-----
for the 35 th cluster number of 10min intervals with zero pickups: 42
-----
for the 36 th cluster number of 10min intervals with zero pickups: 36
-----
for the 37 th cluster number of 10min intervals with zero pickups: 321
-----
for the 38 th cluster number of 10min intervals with zero pickups: 36
-----
for the 39 th cluster number of 10min intervals 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 => ceil(x/4), ceil(x/4), ceil(x/4), ceil(x/4)`
    - Ex2: `__ x => ceil(x/3), ceil(x/3), ceil(x/3)`
  - Case 2:(values missing in middle)
    - Ex1: `x __ y => ceil((x+y)/4), ceil((x+y)/4), ceil((x+y)/4), ceil((x+y)/4)`
    - Ex2: `x __ y => ceil((x+y)/5), ceil((x+y)/5), ceil((x+y)/5), ceil((x+y)/5), ceil((x+y)/5)`
  - Case 3:(values missing at the end)
    - Ex1: `x __ => ceil(x/4), ceil(x/4), ceil(x/4), ceil(x/4)`
    - Ex2: `x _ => ceil(x/2), ceil(x/2)`



```
In [0]: 1 # Fills a value of zero for every bin where no pickup data is present
2 # the count_values: number pickups that are happened in each region for each 10min
3 # there wont be any value if there are no pickups.
4 # values: number of unique bins
5
6 # for every 10min intravel(pickup_bin) we will check it is there in our unique bins
7 # if it is there we will add the count_values[index] to smoothed data
8 # if not we add 0 to the smoothed data
9 # we finally return smoothed data
10 def fill_missing(count_values,values):
11     smoothed_regions=[]
12     ind=0
13     for r in range(0,40):
14         smoothed_bins=[]
15         for i in range(4464):
16             if i in values[r]:
17                 smoothed_bins.append(count_values[ind])
18                 ind+=1
19             else:
20                 smoothed_bins.append(0)
21         smoothed_regions.extend(smoothed_bins)
22     return smoothed_regions
```

```

In [0]: 1 # Fills a value of zero for every bin where no pickup data is present
2 # the count_values: number pickups that are happened in each region for each 10min
3 # there wont be any value if there are no pickups.
4 # values: number of unique bins
5
6 # for every 10min intravel(pickup_bin) we will check it is there in our unique bins
7 # if it is there we will add the count_values[index] to smoothed data
8 # if not we add smoothed data (which is calculated based on the methods that we have)
9 # we finally return smoothed data
10 def smoothing(count_values,values):
11     smoothed_regions=[] # stores list of final smoothed values of each region
12     ind=0
13     repeat=0
14     smoothed_value=0
15     for r in range(0,40):
16         smoothed_bins=[] #stores the final smoothed values
17         repeat=0
18         for i in range(4464):
19             if repeat!=0: # prevents iteration for a value which is already visited
20                 repeat-=1
21                 continue
22             if i in values[r]: #checks if the pickup-bin exists
23                 smoothed_bins.append(count_values[ind]) # appends the value of count_values
24             else:
25                 if i!=0:
26                     right_hand_limit=0
27                     for j in range(i,4464):
28                         if j not in values[r]: #searches for the left-limit
29                             continue
30                         else:
31                             right_hand_limit=j
32                             break
33                     if right_hand_limit==0:
34                         #Case 1: When we have the last/last few values are found
35                         smoothed_value=count_values[ind-1]*1.0/((4463-i)+2)*1
36                         for j in range(i,4464):
37                             smoothed_bins.append(math.ceil(smoothed_value))
38                         smoothed_bins[i-1] = math.ceil(smoothed_value)
39                         repeat=(4463-i)
40                         ind-=1
41                     else:
42                         #Case 2: When we have the missing values between two known values
43                         smoothed_value=(count_values[ind-1]+count_values[ind])/2
44                         for j in range(i,right_hand_limit+1):
45                             smoothed_bins.append(math.ceil(smoothed_value))
46                         smoothed_bins[i-1] = math.ceil(smoothed_value)
47                         repeat=(right_hand_limit-i)
48                 else:
49                     #Case 3: When we have the first/first few values are found
50                     right_hand_limit=0
51                     for j in range(i,4464):
52                         if j not in values[r]:
53                             continue
54                         else:
55                             right_hand_limit=j
56                             break

```

```

57         smoothed_value=count_values[ind]*1.0/((right_hand_limit-i
58         for j in range(i,right_hand_limit+1):
59             smoothed_bins.append(math.ceil(smoothed_value))
60         repeat=(right_hand_limit-i)
61         ind+=1
62         smoothed_regions.extend(smoothed_bins)
63     return smoothed_regions
64

```

```
In [0]: 1 jan_2015_groupby = pd.read_csv('drive/My Drive/NYTD/jan_2015_groupby.csv')
```

```
In [0]: 1 #Filling Missing values of Jan-2015 with 0
2 # here in jan_2015_groupby dataframe the trip_distance represents the number
3 jan_2015_fill = fill_missing(jan_2015_groupby['trip_distance'].values,jan_201
4
5 #Smoothing Missing values of Jan-2015
6 jan_2015_smooth = smoothing(jan_2015_groupby['trip_distance'].values,jan_2015

```

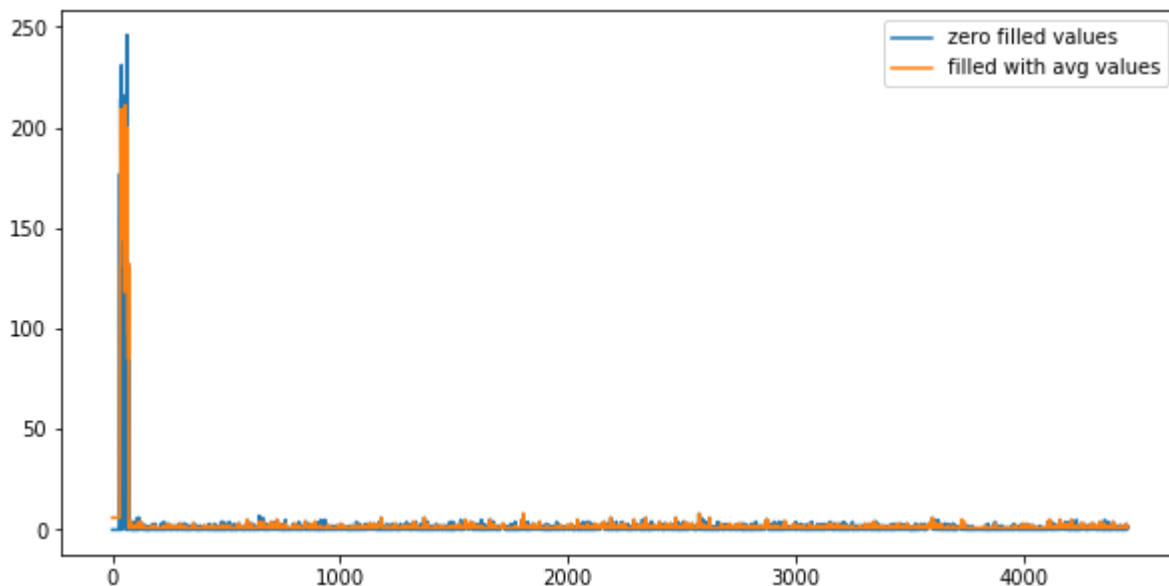
Type Markdown and LaTeX:  $\alpha^2$

```
In [16]: 1 # number of 10min indices for jan 2015= 24*31*60/10 = 4464
2 # number of 10min indices for jan 2016 = 24*31*60/10 = 4464
3 # number of 10min indices for feb 2016 = 24*29*60/10 = 4176
4 # number of 10min indices for march 2016 = 24*30*60/10 = 4320
5 # for each cluster we will have 4464 values, therefore 40*4464 = 178560 (Leng
6 print("number of 10min intravels among all the clusters ",len(jan_2015_fill))

```

number of 10min intravels among all the clusters 178560

```
In [17]: 1 # Smoothing vs Filling
2 # sample plot that shows two variations of filling missing values
3 # we have taken the number of pickups for cluster region 2
4 plt.figure(figsize=(10,5))
5 plt.plot(jan_2015_fill[4464:8920], label="zero filled values")
6 plt.plot(jan_2015_smooth[4464:8920], label="filled with avg values")
7 plt.legend()
8 plt.show()
```



```
In [0]: 1 # why we choose, these methods and which method is used for which data?
2
3 # Ans: consider we have data of some month in 2015 jan 1st, 10 __ 20, i.e
4 # 10st 10min intravel, 0 pickups happened in 2nd 10mins intravel, 0 pickups h
5 # and 20 pickups happened in 4th 10min intravel.
6 # in fill_missing method we replace these values like 10, 0, 0, 20
7 # where as in smoothing method we replace these values as 6,6,6,6,6, if you c
8 # that are happened in the first 40min are same in both cases, but if you can
9 # when you are using smoothing we are looking at the future number of pickup
10
11 # so we use smoothing for jan 2015th data since it acts as our training data
12 # and we use simple fill_missing method for 2016th data.
```

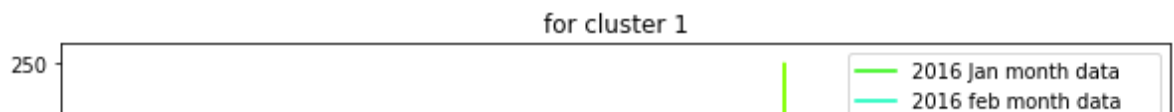
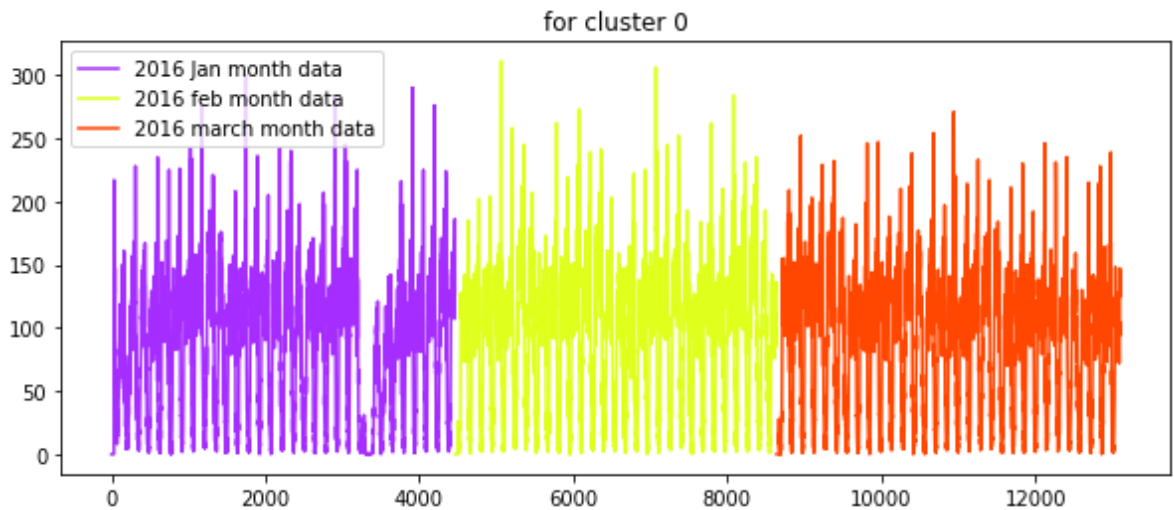
```

In [0]: 1 # Jan-2015 data is smoothed, Jan, Feb & March 2016 data missing values are fil
2 jan_2015_smooth = smoothing(jan_2015_groupby['trip_distance'].values, jan_2015
3 jan_2016_smooth = fill_missing(jan_2016_groupby['trip_distance'].values, jan_2
4 feb_2016_smooth = fill_missing(feb_2016_groupby['trip_distance'].values, feb_2
5 mar_2016_smooth = fill_missing(march_2016_groupby['trip_distance'].values, mar
6
7 # Making list of all the values of pickup data in every bin for a period of 3
8 regions_cum = []
9
10 # a = [1, 2, 3]
11 # b = [2, 3, 4]
12 # a+b = [1, 2, 3, 2, 3, 4]
13
14 # number of 10min indices for jan 2015 = 24*31*60/10 = 4464
15 # number of 10min indices for jan 2016 = 24*31*60/10 = 4464
16 # number of 10min indices for feb 2016 = 24*29*60/10 = 4176
17 # number of 10min indices for march 2016 = 24*31*60/10 = 4464
18 # regions_cum: it will contain 40 lists, each list will contain 4464+4176+446
19 # that are happened for three months in 2016 data
20
21 for i in range(0, 40):
22     regions_cum.append(jan_2016_smooth[4464*i:4464*(i+1)] + feb_2016_smooth[417
23
24 # print(len(regions_cum))
25 # 40
26 # print(len(regions_cum[0]))
27 # 13104

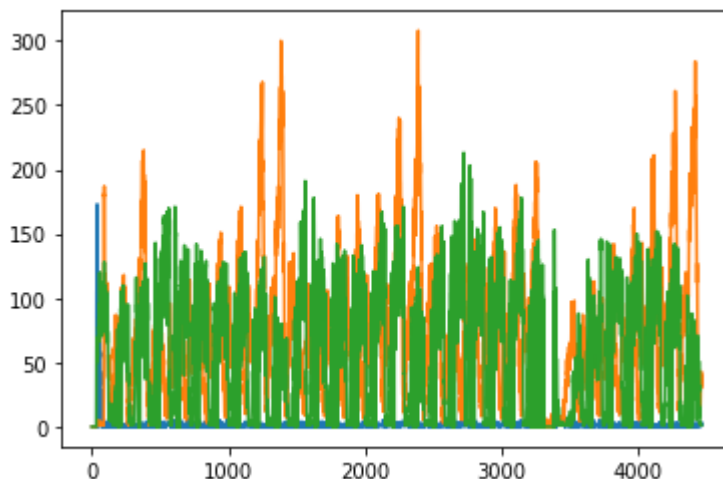
```

## Time series and Fourier Transforms

```
In [0]: 1 def uniqueish_color():
2         """There're better ways to generate unique colors, but this isn't awful."""
3         return plt.cm.gist_ncar(np.random.random())
4 first_x = list(range(0,4464))
5 second_x = list(range(4464,8640))
6 third_x = list(range(8640,13104))
7 for i in range(40):
8     plt.figure(figsize=(10,4))
9     plt.plot(first_x,regions_cum[i][:4464], color=uniqueish_color(), label='2016 jan month data')
10    plt.plot(second_x,regions_cum[i][4464:8640], color=uniqueish_color(), label='2016 feb month data')
11    plt.plot(third_x,regions_cum[i][8640:], color=uniqueish_color(), label='2016 march month data')
12    plt.title('for cluster ' + str(i))
13    plt.legend()
14    plt.show()
```



```
In [0]: 1 plt.plot(regions_cum[1][0:4464])
2         plt.plot(regions_cum[2][0:4464])
3         plt.plot(regions_cum[3][0:4464])
4
5         plt.show()
```

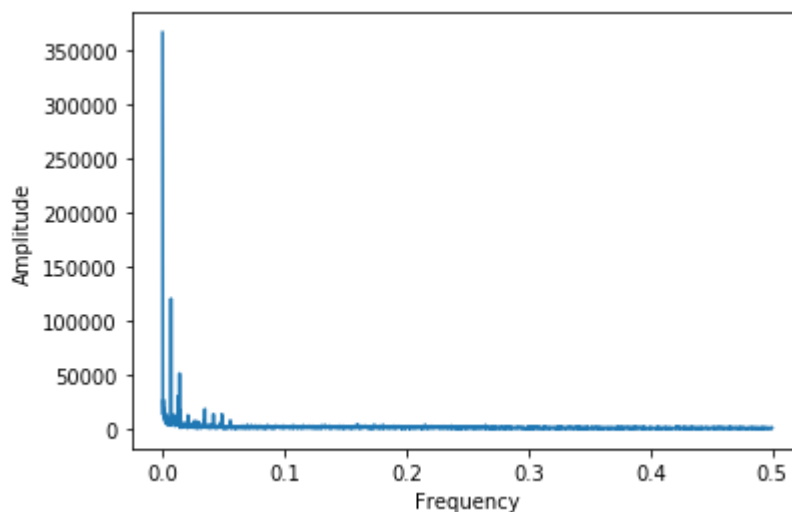


```
In [0]: 1 len(regions_cum)
```

```
Out[50]: 40
```

```
In [0]: 1 %matplotlib inline
```

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



```
In [0]: 1 #Preparing the Dataframe only with x(i) values as jan-2015 data and y(i) va
2 ratios_jan = pd.DataFrame()
3 ratios_jan['Given']=jan_2015_smooth
4 ratios_jan['Prediction']=jan_2016_smooth
5 ratios_jan['Ratios']=ratios_jan['Prediction']*1.0/ratios_jan['Given']*1.0
```

## Modelling: Baseline Models

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

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

## Simple Moving Averages

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

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

```
In [0]: 1 def MA_R_Predictions(ratios,month):
2         predicted_ratio=(ratios['Ratios'].values)[0]
3         error=[]
4         predicted_values=[]
5         window_size=3
6         predicted_ratio_values=[]
7         for i in range(0,4464*40):
8             if i%4464==0:
9                 predicted_ratio_values.append(0)
10                predicted_values.append(0)
11                error.append(0)
12                continue
13                predicted_ratio_values.append(predicted_ratio)
14                predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
15                error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio),2))-((ratios['Given'].values)[i])**2))
16                if i+1>=window_size:
17                    predicted_ratio=sum((ratios['Ratios'].values)[(i+1)-window_size:(i+1)])/(i+1)
18                else:
19                    predicted_ratio=sum((ratios['Ratios'].values)[0:(i+1)])/(i+1)
20
21
22                ratios['MA_R_Predicted'] = predicted_values
23                ratios['MA_R_Error'] = error
24                mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
25                mse_err = sum([e**2 for e in error])/len(error)
26                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 [0]: 1 def MA_P_Predictions(ratios,month):
2         predicted_value=(ratios['Prediction'].values)[0]
3         error=[]
4         predicted_values=[]
5         window_size=1
6         predicted_ratio_values=[]
7         for i in range(0,4464*40):
8             predicted_values.append(predicted_value)
9             error.append(abs((math.pow(predicted_value-(ratios['Prediction'].valu
10             if i+1>=window_size:
11                 predicted_value=int(sum((ratios['Prediction'].values)[(i+1)-windo
12             else:
13                 predicted_value=int(sum((ratios['Prediction'].values)[0:(i+1)))/(
14
15         ratios['MA_P_Predicted'] = predicted_values
16         ratios['MA_P_Error'] = error
17         mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(
18         mse_err = sum([e**2 for e in error])/len(error)
19         return ratios,mape_err,mse_err

```

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

## Weighted Moving Averages

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

Weighted Moving Averages using Ratio Values -

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

```

In [0]: 1 def WA_R_Predictions(ratios,month):
2         predicted_ratio=(ratios['Ratios'].values)[0]
3         alpha=0.5
4         error=[]
5         predicted_values=[]
6         window_size=5
7         predicted_ratio_values=[]
8         for i in range(0,4464*40):
9             if i%4464==0:
10                predicted_ratio_values.append(0)
11                predicted_values.append(0)
12                error.append(0)
13                continue
14                predicted_ratio_values.append(predicted_ratio)
15                predicted_values.append(int(((ratios['Given'].values)[i])*predicted_r
16                error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicte
17                if i+1>=window_size:
18                    sum_values=0
19                    sum_of_coeff=0
20                    for j in range(window_size,0,-1):
21                        sum_values += j*(ratios['Ratios'].values)[i-window_size+j]
22                        sum_of_coeff+=j
23                    predicted_ratio=sum_values/sum_of_coeff
24                else:
25                    sum_values=0
26                    sum_of_coeff=0
27                    for j in range(i+1,0,-1):
28                        sum_values += j*(ratios['Ratios'].values)[j-1]
29                        sum_of_coeff+=j
30                    predicted_ratio=sum_values/sum_of_coeff
31
32                ratios['WA_R_Predicted'] = predicted_values
33                ratios['WA_R_Error'] = error
34                mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(
35                mse_err = sum([e**2 for e in error])/len(error)
36                return ratios,mape_err,mse_err

```

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 5 is optimal for getting the best results using Weighted Moving Averages using previous Ratio values therefore we get

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

Weighted Moving Averages using Previous 2016 Values -

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

```

In [0]: 1 def WA_P_Predictions(ratios,month):
2         predicted_value=(ratios['Prediction'].values)[0]
3         error=[]
4         predicted_values=[]
5         window_size=2
6         for i in range(0,4464*40):
7             predicted_values.append(predicted_value)
8             error.append(abs((math.pow(predicted_value-(ratios['Prediction'].valu
9             if i+1>=window_size:
10                sum_values=0
11                sum_of_coeff=0
12                for j in range(window_size,0,-1):
13                    sum_values += j*(ratios['Prediction'].values)[i-window_size+j
14                    sum_of_coeff+=j
15                predicted_value=int(sum_values/sum_of_coeff)
16
17            else:
18                sum_values=0
19                sum_of_coeff=0
20                for j in range(i+1,0,-1):
21                    sum_values += j*(ratios['Prediction'].values)[j-1]
22                    sum_of_coeff+=j
23                predicted_value=int(sum_values/sum_of_coeff)
24
25        ratios['WA_P_Predicted'] = predicted_values
26        ratios['WA_P_Error'] = error
27        mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(
28        mse_err = sum([e**2 for e in error])/len(error)
29        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)

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

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

For eg. If  $\alpha = 0.9$  then the number of days on which the value of the current iteration is based is  $\sim 1/(1 - \alpha) = 10$  i.e. we consider values 10 days prior before we predict the value for the current

iteration. Also the weights are assigned using  $2/(N + 1) = 0.18$ , where  $N$  = number of prior values being considered, hence from this it is implied that the first or latest value is assigned a weight of 0.18 which keeps exponentially decreasing for the subsequent values.

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

```
In [0]: 1 def EA_R1_Predictions(ratios,month):
2     predicted_ratio=(ratios['Ratios'].values)[0]
3     alpha=0.6
4     error=[]
5     predicted_values=[]
6     predicted_ratio_values=[]
7     for i in range(0,4464*40):
8         if i%4464==0:
9             predicted_ratio_values.append(0)
10            predicted_values.append(0)
11            error.append(0)
12            continue
13            predicted_ratio_values.append(predicted_ratio)
14            predicted_values.append(int(((ratios['Given'].values)[i])*predicted_r
15            error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicte
16            predicted_ratio = (alpha*predicted_ratio) + (1-alpha)*((ratios['Ratio
17
18    ratios['EA_R1_Predicted'] = predicted_values
19    ratios['EA_R1_Error'] = error
20    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(
21    mse_err = sum([e**2 for e in error])/len(error)
22    return ratios,mape_err,mse_err
```

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

```
In [0]: 1 def EA_P1_Predictions(ratios,month):
2     predicted_value= (ratios['Prediction'].values)[0]
3     alpha=0.3
4     error=[]
5     predicted_values=[]
6     for i in range(0,4464*40):
7         if i%4464==0:
8             predicted_values.append(0)
9             error.append(0)
10            continue
11            predicted_values.append(predicted_value)
12            error.append(abs((math.pow(predicted_value-(ratios['Prediction'].valu
13            predicted_value =int((alpha*predicted_value) + (1-alpha)*((ratios['Pr
14
15    ratios['EA_P1_Predicted'] = predicted_values
16    ratios['EA_P1_Error'] = error
17    mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(
18    mse_err = sum([e**2 for e in error])/len(error)
19    return ratios,mape_err,mse_err
```

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

Type *Markdown* and LaTeX:  $\alpha^2$

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

Error Metric Matrix (Forecasting Methods) - MAPE & MSE

```
-----
Moving Averages (Ratios) - MAPE: 0.182115517339
MSE: 400.0625504032258
Moving Averages (2016 Values) - MAPE: 0.14292849687
MSE: 174.84901993727598
-----
Weighted Moving Averages (Ratios) - MAPE: 0.178486925438
MSE: 384.01578741039424
Weighted Moving Averages (2016 Values) - MAPE: 0.135510884362
MSE: 162.46707549283155
-----
Exponential Moving Averages (Ratios) - MAPE: 0.177835501949
MSE: 378.34610215053766
Exponential Moving Averages (2016 Values) - MAPE: 0.135091526367
MSE: 159.73614471326164
```

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

From the above matrix it is inferred that the best forecasting model for our prediction would be:-

$$P'_t = \alpha * P_{t-1} + (1 - \alpha) * P'_{t-1} \text{ i.e Exponential Moving Averages using 2016 Values}$$

## Regression Models

### Top fourier features(frequencies and amplitudes)

#### Train-Test Split

Before we start predictions using the tree based regression models we take 3 months of 2016 pickup data and split it such that for every region we have 70% data in train and 30% in test, ordered date-wise for every region

```

In [0]: 1 #now we will calculate the top frequencies and amplitudes
2
3 #all the fourier features
4 final_df = pd.DataFrame(columns = ["Amplitude1", "Amplitude2", "Amplitude3", "Am
5 for i in range(0,40):
6
7
8     #we will finally append the values in dataframe
9     jan_df = pd.DataFrame()
10    feb_df = pd.DataFrame()
11    march_df = pd.DataFrame()
12
13    #we will calculate fft for each month
14    #as length of regions_cum is where it holds data for each of the clusters
15
16    #for month of january
17
18    january_data = regions_cum[i][0:4464]
19    jan_df['Amplitude'] = sorted(np.fft.fft(january_data), reverse = True)[0:5]
20    jan_df['Frequency'] = sorted(np.fft.fftfreq(4464,1), reverse = True)[0:5]
21    jan_df = jan_df.reset_index(drop = True).T
22    jan_list = []
23    #jan_final = fourier(jan_df, 4464)
24
25
26    #feb data
27    february_data = regions_cum[i][4464:4464 + 4176]
28    feb_df['Amplitude'] = sorted(np.fft.fft(february_data), reverse = True)[0:5]
29    feb_df['Frequency'] = sorted(np.fft.fftfreq(4464,1), reverse = True)[0:5]
30    feb_df = feb_df.reset_index(drop = True).T
31    feb_list = []
32    #feb_final = fourier(feb_df, 4176)
33
34    #march data
35    march_data = regions_cum[i][4464+4176:4464+4176+4464]
36    march_df['Amplitude'] = sorted(np.fft.fft(march_data), reverse = True)[0:5]
37    march_df['Frequency'] = sorted(np.fft.fftfreq(4464,1), reverse = True)[0:5]
38    march_df = march_df.reset_index(drop = True).T
39    march_list = []
40
41
42    jan_list = []
43    feb_list = []
44    march_list = []
45    for i in range(0,5):
46
47        jan_list.append(float(fr_am_jan_sorted[i]['Frequency']))
48        jan_list.append(float(fr_am_jan_sorted[i]['Amplitude']))
49
50        feb_list.append(float(fr_am_feb_sorted[i]['Frequency']))
51        feb_list.append(float(fr_am_feb_sorted[i]['Amplitude']))
52
53        march_list.append(float(fr_am_mar_sorted[i]['Frequency']))
54        march_list.append(float(fr_am_mar_sorted[i]['Amplitude']))
55
56    new_jan = pd.DataFrame([fr_am_list_jan]*4464)

```

```
57 new_feb = pd.DataFrame([fr_am_list_feb]*4176)
58 new_mar = pd.DataFrame([fr_am_list_mar]*4464)
59
60 new_jan.columns = ["Amplitude1", "Amplitude2", "Amplitude3", "Amplitude4", "A
61 new_feb.columns = ["Amplitude1", "Amplitude2", "Amplitude3", "Amplitude4", "A
62 new_mar.columns = ["Amplitude1", "Amplitude2", "Amplitude3", "Amplitude4", "A
63
64
65 final_df = final_df.append(fr_am_new_jan, ignore_index=True)
66 final_df = final_df.append(fr_am_new_feb, ignore_index=True)
67 final_df = final_df.append(fr_am_new_mar, ignore_index=True)
```

```
In [0]: 1 final_df = fourier_features_df
```

```
In [22]: 1 final_df.head()
```

Out[22]:

	Amplitude1	Amplitude2	Amplitude3	Amplitude4	Amplitude5	Freq1	Freq2	Fre
0	367173.0	94490.188858	94490.188858	14349.849101	14349.849101	0.499776	0.499552	0.4993
1	367173.0	94490.188858	94490.188858	14349.849101	14349.849101	0.499776	0.499552	0.4993
2	367173.0	94490.188858	94490.188858	14349.849101	14349.849101	0.499776	0.499552	0.4993
3	367173.0	94490.188858	94490.188858	14349.849101	14349.849101	0.499776	0.499552	0.4993
4	367173.0	94490.188858	94490.188858	14349.849101	14349.849101	0.499776	0.499552	0.4993



In [23]:

```

1  ## Feature engineering part for the data
2  # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps)
3
4  f_train = pd.DataFrame(columns=['Amplitude1','Amplitude2','Amplitude3','Ampli
5  f_test = pd.DataFrame(columns=['Amplitude1','Amplitude2','Amplitude3','Ampli
6  #Train data
7  for i in range(40):
8      f_train = f_train.append(final_df[i*13099 : 13099*i + 9169])
9
10 #Test data
11 for i in range(40):
12     f_test = f_test.append(final_df[i*13099 + 9169 : 13099*(i+1)])
13
14 #Reset all the indexes in train and test data
15 f_train.reset_index(inplace=True)
16 f_test.reset_index(inplace = True)
17
18 fourier_train_df=f_train.drop(labels="index", axis=1)
19 fourier_test_df=f_test.drop(labels="index", axis=1)
20
21 print("Shape of the fourier transformed train dataframe: ",fourier_train_df.s
22 print("Shape of the fourier transformed test dataframe: ",fourier_test_df.sha
23 fourier_test_df.head()
24

```

Shape of the fourier transformed train dataframe: (366760, 10)

Shape of the fourier transformed test dataframe: (157200, 10)

Out[23]:

	Amplitude1	Amplitude2	Amplitude3	Amplitude4	Amplitude5	Freq1	Freq2	Freq
0	387761.0	91160.781939	91160.781939	17509.351171	17509.351171	0.499776	0.499552	0.4993:
1	387761.0	91160.781939	91160.781939	17509.351171	17509.351171	0.499776	0.499552	0.4993:
2	387761.0	91160.781939	91160.781939	17509.351171	17509.351171	0.499776	0.499552	0.4993:
3	387761.0	91160.781939	91160.781939	17509.351171	17509.351171	0.499776	0.499552	0.4993:
4	387761.0	91160.781939	91160.781939	17509.351171	17509.351171	0.499776	0.499552	0.4993:

## Time Series data feature engineering

In this part I have considered two different parts ,first one is High band pass filter and Wavelet Denposing while the second part is Holt's Trend Model adn smoothing

### Part 1:High Band pass filer + Wavelet Denoising

for understanding about the wavelet denosing part I have referred to following blogs and kernels whcih have explained the concepts and code fluently and in an excellent manner

1.<https://www.kaggle.com/jackvial/dwt-signal-denoising> (<https://www.kaggle.com/jackvial/dwt-signal-denoising>)

2.<https://machinelearningmastery.com/time-series-forecasting-methods-in-python-cheat-sheet/>  
(<https://machinelearningmastery.com/time-series-forecasting-methods-in-python-cheat-sheet/>)

3.<https://www.kaggle.com/theoviel/fast-fourier-transform-denoising>  
(<https://www.kaggle.com/theoviel/fast-fourier-transform-denoising>)

The high amplitude, unstable parts of the signal (with high peaks and valleys) are the actual signal we are looking for. The medium amplitude patches of the signal (between the high amplitude regions) represent the unnecessary noise and artificial impulse, and the wavelet method seems to do better at removing these patches. This denoising illustrated in the image at the beginning of the kernel.

```
In [0]: 1 import os
        2 import gc
        3 from numpy.fft import *
        4 import pywt
        5 from statsmodels.robust import mad
        6 import scipy
        7 from scipy import signal
        8 from scipy.signal import butter, deconvolve
        9 import warnings
       10 warnings.filterwarnings('ignore')
```

```
In [0]: 1 SIGNAL_LEN = 4464 + 4176 + 4464 #signal length corresponds to the length size
        2 SAMPLE_RATE = 40000 #this is the sampling rate we are considering here
```


```
In [0]: 1 signals = []
        2 for i in range(40):
        3     signals.append(np.array(regions_cum[i][0:4464+4176+4464]))
        4
        5
```

```

In [0]: 1 def maddest(d, axis=None):
2         """The mean absolute deviation
3         This calculates the mean of the absolute values of the deviations of the
4         of the time series. It is a measure of entropy or disorder in the time se
5         The greater the MAD value, the more disorderly and unpredictable the time
6
7         return np.mean(np.absolute(d - np.mean(d, axis)), axis)
8
9         #=====
10 def high_pass_filter(x, low_cutoff=1000, SAMPLE_RATE=SAMPLE_RATE):
11     """
12     From @randxie https://github.com/randxie/Kaggle-VSB-Baseline/blob/master/
13     Modified to work with scipy version 1.1.0 which does not have the fs para
14     """
15
16     # nyquist frequency is half the sample rate https://en.wikipedia.org/wiki
17     nyquist = 0.5 * SAMPLE_RATE
18     norm_low_cutoff = low_cutoff / nyquist
19
20     # Fault pattern usually exists in high frequency band. According to liter
21     sos = butter(10, Wn=[norm_low_cutoff], btype='highpass', output='sos')
22     filtered_sig = signal.sosfilt(sos, x)
23
24     return filtered_sig
25
26     #=====
27     #after passing from the filter we will finally denoise the signal
28 def denoise_signal(x, wavelet='db4', level=1):
29     """
30     1. Adapted from waveletSmooth function found here:
31     http://connor-johnson.com/2016/01/24/using-pywavelets-to-remove-high-freq
32     2. Threshold equation and using hard mode in threshold as mentioned
33     in section '3.2 denoising based on optimized singular values' from paper
34     http://dspace.vsb.cz/bitstream/handle/10084/133114/VAN431_FEI_P1807_1801V
35     """
36
37     # Decompose to get the wavelet coefficients
38     coeff = pywt.wavedec(x, wavelet, mode="per")
39
40     # Calculate sigma for threshold as defined in http://dspace.vsb.cz/bitstr
41     # As noted by @harshit92 MAD referred to in the paper is Mean Absolute De
42     sigma = (1/0.6745) * maddest(coeff[-level])
43
44     # Calculate the universal threshold
45     uthresh = sigma * np.sqrt(2*np.log(len(x)))
46     coeff[1:] = (pywt.threshold(i, value=uthresh, mode='hard') for i in coeff
47
48     # Reconstruct the signal using the thresholded coefficients
49     return pywt.waverec(coeff, wavelet, mode='per')

```

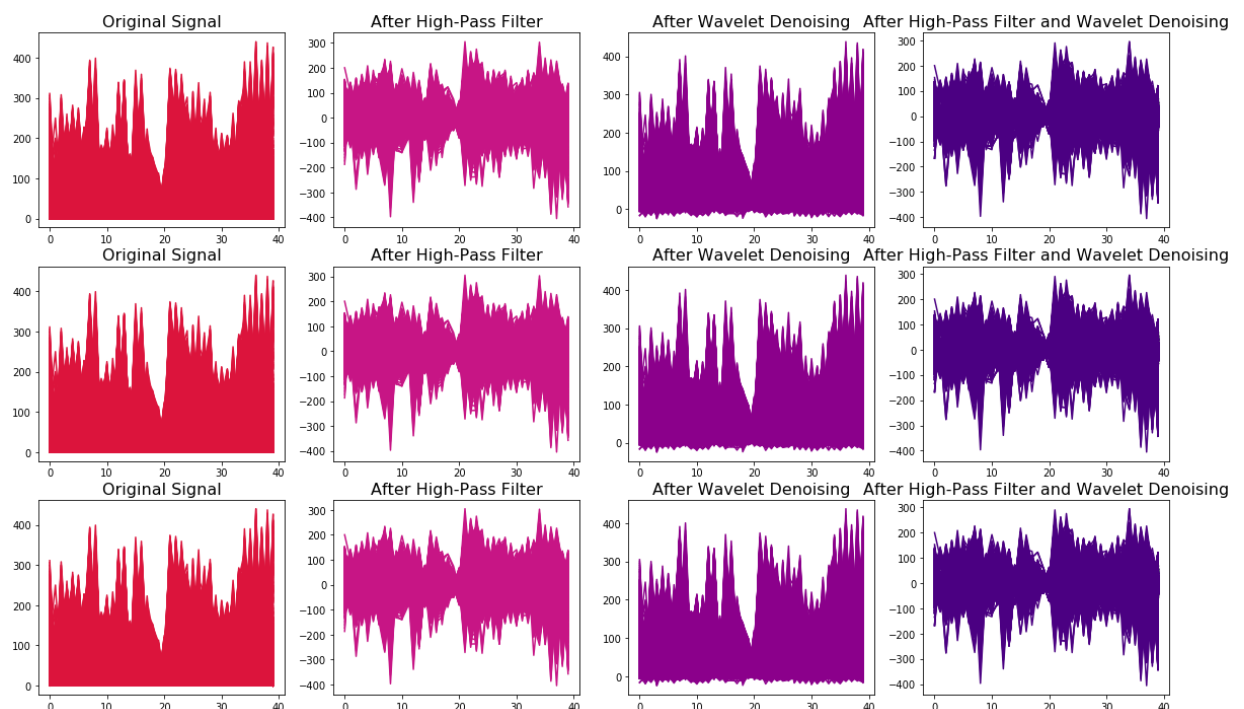
```
In [0]: 1 #storing the denosied signals for all the regions
2 denoised_signals = []
3 for i in range(40):
4     recon_signal = denoise_signal(high_pass_filter(signals[i], low_cutoff=100
5     denoised_signals.append(recon_signal)
```



```

In [0]: 1 #plotting for filter,denosing and then both of them combined
2 fig, ax = plt.subplots(nrows=3, ncols=4, figsize=(20, 12))
3
4 ax[0, 0].plot(signals[0], 'crimson')
5 ax[0, 0].set_title('Original Signal', fontsize=16)
6 ax[0, 1].plot(high_pass_filter(signals[0], low_cutoff=10000, SAMPLE_RATE=4000)
7 ax[0, 1].set_title('After High-Pass Filter', fontsize=16)
8 ax[0, 2].plot(denoise_signal(signals[0]), 'darkmagenta')
9 ax[0, 2].set_title('After Wavelet Denoising', fontsize=16)
10 ax[0, 3].plot(denoise_signal(high_pass_filter(signals[0], low_cutoff=10000, S
11 ax[0, 3].set_title('After High-Pass Filter and Wavelet Denoising', fontsize=1
12
13 ax[1, 0].plot(signals[1], 'crimson')
14 ax[1, 0].set_title('Original Signal', fontsize=16)
15 ax[1, 1].plot(high_pass_filter(signals[1], low_cutoff=10000, SAMPLE_RATE=4000
16 ax[1, 1].set_title('After High-Pass Filter', fontsize=16)
17 ax[1, 2].plot(denoise_signal(signals[1]), 'darkmagenta')
18 ax[1, 2].set_title('After Wavelet Denoising', fontsize=16)
19 ax[1, 3].plot(denoise_signal(high_pass_filter(signals[1], low_cutoff=10000, S
20 ax[1, 3].set_title('After High-Pass Filter and Wavelet Denoising', fontsize=1
21
22 ax[2, 0].plot(signals[2], 'crimson')
23 ax[2, 0].set_title('Original Signal', fontsize=16)
24 ax[2, 1].plot(high_pass_filter(signals[2], low_cutoff=10000, SAMPLE_RATE=4000
25 ax[2, 1].set_title('After High-Pass Filter', fontsize=16)
26 ax[2, 2].plot(denoise_signal(signals[2]), 'darkmagenta')
27 ax[2, 2].set_title('After Wavelet Denoising', fontsize=16)
28 ax[2, 3].plot(denoise_signal(high_pass_filter(signals[2], low_cutoff=10000, S
29 ax[2, 3].set_title('After High-Pass Filter and Wavelet Denoising', fontsize=1
30
31 plt.show()

```



## Part 2: Holt's Linear Trend Model

It is an extension of simple exponential smoothing to allow forecasting of data with a trend. This method takes into account the trend of the dataset. The forecast function in this method is a function of level and trend. First, let's visualize the trend, seasonality and error in the series and then decompose the time series in four parts.

Observed, which is the original time series. Trend, which shows the trend in the time series, i.e., increasing or decreasing behaviour of the time series. Seasonal, which tells us about the seasonality in the time series. Residual, which is obtained by removing any trend or seasonality in the time series

For understanding about it and the code snippet I have used is:

1. <https://grisha.org/blog/2016/01/29/triple-exponential-smoothing-forecasting/>  
(<https://grisha.org/blog/2016/01/29/triple-exponential-smoothing-forecasting/>)
2. <https://grisha.org/blog/2016/02/16/triple-exponential-smoothing-forecasting-part-ii/>  
(<https://grisha.org/blog/2016/02/16/triple-exponential-smoothing-forecasting-part-ii/>)
3. <https://www.kaggle.com/lampubhutia/timeseries-modelling-predicting-traffic-growth>  
(<https://www.kaggle.com/lampubhutia/timeseries-modelling-predicting-traffic-growth>)

In [24]:

```

1 def initial_trend(series, slen):
2     #here series is the signal we want
3     sum = 0.0
4     for i in range(slen):
5         sum += float(series[i+slen] - series[i]) / slen
6     return sum / slen
7     #=====
8 def initial_seasonal_components(series, slen):
9     seasonals = {}
10    season_averages = []
11    n_seasons = int(len(series)/slen)
12    # compute season averages
13    for j in range(n_seasons):
14        season_averages.append(sum(series[slen*j:slen*j+slen])/float(slen))
15    # compute initial values
16    for i in range(slen):
17        sum_of_vals_over_avg = 0.0
18        for j in range(n_seasons):
19            sum_of_vals_over_avg += series[slen*j+i]-season_averages[j]
20        seasonals[i] = sum_of_vals_over_avg/n_seasons
21    return seasonals
22    #=====
23 def triple_exponential_smoothing(series, slen, alpha, beta, gamma, n_preds):
24     result = []
25     seasonals = initial_seasonal_components(series, slen)
26     for i in range(len(series)+n_preds):
27         if i == 0: # initial values
28             smooth = series[0]
29             trend = initial_trend(series, slen)
30             result.append(series[0])
31             continue
32         if i >= len(series): # we are forecasting
33             m = i - len(series) + 1
34             result.append((smooth + m*trend) + seasonals[i%slen])
35         else:
36             val = series[i]
37             last_smooth, smooth = smooth, alpha*(val-seasonals[i%slen]) + (1-
38             trend = beta * (smooth-last_smooth) + (1-beta)*trend
39             seasonals[i%slen] = gamma*(val-smooth) + (1-gamma)*seasonals[i%sl
40             result.append(smooth+trend+seasonals[i%slen])
41     return result
42
43 #Holt Winters initialization of variables: # https://robjhyndman.com/hyndsigh
44 alpha = 0.3
45 beta = 0.15
46 gamma = 0.2
47 season_len = 24
48
49 #Prepare the features for all points for all clusters
50 predict_values_three = []
51 predict_final = []
52 for r in range(0,40):
53     predict_values_three = triple_exponential_smoothing(regions_cum[r][0:1310
54     predict_final.append(predict_values_three[5:])
55

```

```
Out[24]: [0,
-0.15001182844932864,
0.04973103441697209,
0.1047856693471457,
0.2207239432930307,
-0.13916905851025202,
0.0798387775307404,
-0.26116276794373344,
0.09459909749222546,
0.17075822779303143]
```

## Using the final features

```
In [0]: 1
2 # we take number of pickups that are happened in last 5 10min intravels
3 number_of_time_stamps = 5
4
5 # output variable
6 # it is list of lists
7 # it will contain number of pickups 13099 for each cluster
8 output = []
9
10
11 # tsne_lat will contain 13104-5=13099 times latitude of cluster center for e
12 # Ex: [[cent_lat 13099times],[cent_lat 13099times], [cent_lat 13099times]....
13 # it is list of lists
14 tsne_lat = []
15 tsne_lon = []#similarly for Longitudes
16
17 # we will code each day
18 # sunday = 0, monday=1, tue = 2, wed=3, thur=4, fri=5,sat=6
19 # for every cluster we will be adding 13099 values, each value represent to w
20 # it is list of lists
21 tsne_weekday = []
22
23 # its an numpy array, of shape (523960, 5)
24 # each row corresponds to an entry in out data
25 # for the first row we will have [f0,f1,f2,f3,f4] fi=number of pickups happen
26 # the second row will have [f1,f2,f3,f4,f5]
27 # the third row will have [f2,f3,f4,f5,f6]
28 # and so on...
29 tsne_feature = []
30 tsne_feature = [0]*number_of_time_stamps
31 for i in range(0,40):
32     tsne_lat.append([kmeans.cluster_centers_[i][0]]*13099)
33     tsne_lon.append([kmeans.cluster_centers_[i][1]]*13099)
34     # jan 1st 2016 is thursday, so we start our day from 4: "(int(k/144))%7+4
35     # our prediction start from 5th 10min intravel since we need to have numb
36     tsne_weekday.append([int(((int(k/144))%7+4)%7) for k in range(5,4464+4176
37     # regions_cum is a list of lists [[x1,x2,x3..x13104], [x1,x2,x3..x13104],
38     tsne_feature = np.vstack((tsne_feature, [regions_cum[i][r:r+number_of_tim
39     output.append(regions_cum[i][5:]))
40 tsne_feature = tsne_feature[1:]
```



```
In [26]: 1 # train, test split : 70% 30% split
2 # Before we start predictions using the tree based regression models we take
3 # and split it such that for every region we have 70% data in train and 30% i
4 # ordered date-wise for every region
5 print("size of train data :", int(13099*0.7))
6 print("size of test data :", int(13099*0.3))
```

size of train data : 9169

size of test data : 3929

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

```
In [0]: 1 # Getting the predictions of exponential moving averages to be used as a feat
2
3 # upto now we computed 8 features for every data point that starts from 50th
4 # 1. cluster center latitude
5 # 2. cluster center longitude
6 # 3. day of the week
7 # 4. f_t_1: number of pickups that are happened previous t-1th 10min intravel
8 # 5. f_t_2: number of pickups that are happened previous t-2th 10min intravel
9 # 6. f_t_3: number of pickups that are happened previous t-3th 10min intravel
10 # 7. f_t_4: number of pickups that are happened previous t-4th 10min intravel
11 # 8. f_t_5: number of pickups that are happened previous t-5th 10min intravel
12
13 # from the baseline models we said the exponential weighted moving avarage gi
14 # we will try to add the same exponential weighted moving avarage at t as a f
15 # exponential weighted moving avarage => p'(t) = alpha*p'(t-1) + (1-alpha)*P(
16 alpha=0.3
17
18 # it is a temporary array that store exponential weighted moving avarage for
19 # for each cluster it will get reset
20 # for every cluster it contains 13104 values
21 predicted_values=[]
22
23 # it is similar like tsne_lat
24 # it is list of lists
25 # predict_list is a list of lists [[x5,x6,x7..x13104], [x5,x6,x7..x13104], [x
26 predict_list = []
27 tsne_flat_exp_avg = []
28 for r in range(0,40):
29     for i in range(0,13104):
30         if i==0:
31             predicted_value= regions_cum[r][0]
32             predicted_values.append(0)
33             continue
34             predicted_values.append(predicted_value)
35             predicted_value =int((alpha*predicted_value) + (1-alpha)*(regions_cum
36 predict_list.append(predicted_values[5:])
37 predicted_values=[]
```

In [41]:

```

1 print("Number of data clusters",len(train_features), "Number of data points i
2 print("Number of data clusters",len(train_features), "Number of data points i

```

Number of data clusters 40 Number of data points in trian data 9169 Each data p  
oint contains 5 features

Number of data clusters 40 Number of data points in test data 3930 Each data po  
int contains 5 features

In [0]:

```

1 # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps)
2 #Extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) f
3 tsne_train_flat_lat = [i[:9169] for i in tsne_lat]
4 tsne_train_flat_lon = [i[:9169] for i in tsne_lon]
5 tsne_train_flat_weekday = [i[:9169] for i in tsne_weekday]
6 tsne_train_flat_output = [i[:9169] for i in output]
7 tsne_train_flat_exp_avg = [i[:9169] for i in predict_list]
8 tsne_train_flat_triple_exp = [i[:9169] for i in predict_list_three]
9
10
11
12 # extracting the rest of the timestamp values i.e 30% of 12956 (total timesta
13 #Extracting the rest of the timestamp values i.e 30% of 12956 (total timestam
14 tsne_test_flat_lat = [i[9169:] for i in tsne_lat]
15 tsne_test_flat_lon = [i[9169:] for i in tsne_lon]
16 tsne_test_flat_weekday = [i[9169:] for i in tsne_weekday]
17 tsne_test_flat_output = [i[9169:] for i in output]
18 tsne_test_flat_exp_avg = [i[9169:] for i in predict_list]
19 #tsne_test_flat_double_exp = [i[9169:] for i in predict_list_two]
20 tsne_test_flat_triple_exp = [i[9169:] for i in predict_list_three]
21
22
23 # the above contains values in the form of list of lists (i.e. list of values
24 train_new_features = []
25 for i in range(0,40):
26     train_new_features.extend(train_features[i])
27 test_new_features = []
28 for i in range(0,40):
29     test_new_features.extend(test_features[i])
30
31
32
33 # converting lists of lists into sinle list i.e flatten
34 # a = [[1,2,3,4],[4,6,7,8]]
35 # print(sum(a,[]))
36 # [1, 2, 3, 4, 4, 6, 7, 8]
37 tsne_train_lat = sum(tsne_train_flat_lat, [])
38 tsne_train_lon = sum(tsne_train_flat_lon, [])
39 tsne_train_weekday = sum(tsne_train_flat_weekday, [])
40 tsne_train_output = sum(tsne_train_flat_output, [])
41 tsne_train_exp_avg = sum(tsne_train_flat_exp_avg,[])
42 tsne_train_flat_triple_exp = sum(tsne_train_flat_triple_exp,[])
43

```

```
In [48]: 1 # Preparing the data frame for our train data
2 # Preparing the data frame for our train data
3 columns = ['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1']
4 df_train = pd.DataFrame(data=train_new_features, columns=columns)
5 df_train['lat'] = tsne_train_lat
6 df_train['lon'] = tsne_train_lon
7 df_train['weekday'] = tsne_train_weekday
8 df_train['exp_avg'] = tsne_train_exp_avg
9 #df_train['double_exp'] = tsne_train_flat_double_exp
10 df_train['triple_exp'] = tsne_train_flat_triple_exp
11
12 # Preparing the data frame for our train data
13 df_test = pd.DataFrame(data=test_new_features, columns=columns)
14 df_test['lat'] = tsne_test_lat
15 df_test['lon'] = tsne_test_lon
16 df_test['weekday'] = tsne_test_weekday
17 df_test['exp_avg'] = tsne_test_exp_avg
18 #df_test['double_exp'] = tsne_test_flat_double_exp
19 df_test['triple_exp'] = tsne_test_flat_triple_exp
20 print(df_test.shape)
```

(366760, 10)

```
In [0]: 1 #Mering the train and test df with fourier features train and test df
2 df_train = pd.concat([df_train, fourier_train_df], axis = 1)
3 df_test = pd.concat([df_test, fourier_test_df], axis = 1)
```

```
In [0]: 1 #Save the final dataframes along with output
2 df_test.to_csv("df_test.csv", index=None)
3 df_train.to_csv("df_train.csv", index=None)
4
5 import pickle
6 with open('tsne_train_output.pkl', 'wb') as file:
7     pickle.dump(tsne_train_output, file)
8 with open('tsne_test_output.pkl', 'wb') as file:
9     pickle.dump(tsne_test_output, file)
```

## Machine Learning Models

We will use Linear Regression, Random Forest regression and XGBoost regression and tune the parameters using GridSearch Cross validation

### Model1: Using Linear Regression and Tuning

```

In [0]: 1 from sklearn.linear_model import LinearRegression
2 from sklearn.linear_model import SGDRegressor
3 from sklearn.model_selection import GridSearchCV
4
5 #declaring the parameters to be tuned
6 #Declaring parameters
7 params = {'alpha':[0.0001,0.001,0.01,0.1],
8           'penalty':['l2','l1'],
9           'fit_intercept':[True,False]}
10
11 #Using SGRRegressor with squared loss, in order to fine tune the hyperparameters
12 start =datetime.datetime.now()
13 print('Hyperparameter tuning: \n')
14 model= SGDRegressor(loss='squared_loss', l1_ratio=0.15, max_iter=1000, tol=0.
15                     verbose=0, epsilon=0.1, random_state=None, learning_rate=
16                     validation_fraction=0.1, n_iter_no_change=5, warm_start=False)
17 g_search =GridSearchCV(model, params, scoring='neg_mean_absolute_error',cv=3,
18 g_search.fit(df_train, tsne_train_output)
19 print('Time taken to perform Hyperparameter tuning :',datetime.datetime.now()
20
21 #Getting the best hyperparameter tuned model
22 best_model=g_search.best_estimator_
23 print("Best estimator: ",best_model)
24
25 #Fitting the best model to our training data
26 best_model.fit(df_train, tsne_train_output)

```

Hyperparameter tuning:

Time taken to perform Hyperparameter tuning : 0:34:04.541462

Best estimator: SGDRegressor(alpha=0.001, average=False, early\_stopping=False, epsilon=0.1,

eta=0.01, fit\_intercept=False, l1\_ratio=0.15, learning\_rate='invscaling', loss='squared\_loss', max\_iter=1000, n\_iter\_no\_change=5, penalty='l2', power\_t=0.25, random\_state=None, shuffle=True, tol=0.001, validation\_fraction=0.1, verbose=0, warm\_start=False)

Out[71]: SGDRegressor(alpha=0.001, average=False, early\_stopping=False, epsilon=0.1, eta=0.01, fit\_intercept=False, l1\_ratio=0.15, learning\_rate='invscaling', loss='squared\_loss', max\_iter=1000, n\_iter\_no\_change=5, penalty='l2', power\_t=0.25, random\_state=None, shuffle=True, tol=0.001, validation\_fraction=0.1, verbose=0, warm\_start=False)

```
In [65]: 1 from sklearn.linear_model import SGDRegressor
2 best_est = SGDRegressor(alpha=0.001, average=False, early_stopping=False, eps
3     eta0=0.01, fit_intercept=False, l1_ratio=0.15,
4     learning_rate='invscaling', loss='squared_loss', max_iter=1000,
5     n_iter_no_change=5, penalty='l2', power_t=0.25, random_state=None,
6     shuffle=True, tol=0.001, validation_fraction=0.1, verbose=0,
7     warm_start=False)
8 best_est.fit(df_train,tsne_train_output)
```

```
Out[65]: SGDRegressor(alpha=0.001, average=False, early_stopping=False, epsilon=0.1,
    eta0=0.01, fit_intercept=False, l1_ratio=0.15,
    learning_rate='invscaling', loss='squared_loss', max_iter=1000,
    n_iter_no_change=5, penalty='l2', power_t=0.25, random_state=None,
    shuffle=True, tol=0.001, validation_fraction=0.1, verbose=0,
    warm_start=False)
```

```
In [0]: 1 #Using the best model to make predictions
2 pred_test = (best_est.predict(df_test))
3 pred_train = (best_est.predict(df_train))
4 pred_train = [round(i) for i in pred_train]#rounding up for integer value
5 pred_test = [round(i) for i in pred_test]
6
```

## Using Random Forest Regressor

```
In [8]: 1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.model_selection import RandomizedSearchCV
3 from scipy.stats import randint as sp_randint #for uniform discrete random va
4 from scipy.stats import uniform
5
6 param_dist = {"n_estimators":sp_randint[500,600,800,1000,1200,1500,2000],#num
7     "max_depth": sp_randint(2,9),#the depth of the the trees
8     "min_samples_leaf": sp_randint(20,65)} #number of leafs
9 start = datetime.datetime.now()
10 clf = RandomForestRegressor(random_state=25,n_jobs=-1)
11
12 rf_random = RandomizedSearchCV(clf, param_distributions=param_dist,
13     n_iter=5,cv=10,scoring='neg_mean_absolute_
14 rf_random.fit(df_train,tsne_train_output)
15 print('total time taken for computation is:',datetime.datetime.now() - start)
```

time taken to perform the hyperparameter tuning: 1:12:29.091710

```
In [52]: 1 model = RandomForestRegressor(bootstrap=True, criterion='mse', max_depth=4,
2         max_features='auto', max_leaf_nodes=100,
3         min_impurity_decrease=0.0, min_impurity_split=None,
4         min_samples_leaf=25, min_samples_split=7,
5         min_weight_fraction_leaf=0.0, n_estimators=1500, n_jobs=-1,
6         oob_score=False, random_state=None, verbose=0, warm_start=False)
7
8 model.fit(df_train,tsne_train_output)
```

```
Out[52]: RandomForestRegressor(bootstrap=True, ccp_alpha=0.0, criterion='mse',
                                max_depth=4, max_features='auto', max_leaf_nodes=100,
                                max_samples=None, min_impurity_decrease=0.0,
                                min_impurity_split=None, min_samples_leaf=25,
                                min_samples_split=7, min_weight_fraction_leaf=0.0,
                                n_estimators=1500, n_jobs=-1, oob_score=False,
                                random_state=None, verbose=0, warm_start=False)
```

```
In [0]: 1 #Using the best model to make predictions
2 y_pred = model.predict(df_test)
3 rf_test_predictions = [round(value) for value in y_pred] #rounding the values
4 y_pred = model.predict(df_train)
5 rf_train_predictions = [round(value) for value in y_pred]
```

## Using XgBoost Regressor

```
In [9]: 1 from xgboost import XGBRegressor
2 #Declaring parameters
3 params = {'learning_rate':[0.1,0.01,0.001,0.0001],
4           'n_estimators':[250,500,750,1000,1500,2000,3000],
5           'subsample':[0.6,0.7,0.8,0.9],
6           'min_child_weight':[3,5,7,9],
7           'reg_lambda':[100,200,300,400],
8           'reg_alpha':[100,200,300, 400],
9           'max_depth': [3,4,5,6,7,9],
10          'colsample_bytree':[0.6,0.7,0.8],
11          'gamma':[0,0.5,1]}
12
13 #Tuning hyperparameters
14 start =datetime.now()
15 print('Hyperparameter tuning: \n')
16 model= XGBRegressor(random_state=0,n_jobs=-1)
17 rsearch = RandomizedSearchCV(model,params,n_iter=20,scoring='neg_mean_absolut
18 rsearch.fit(df_train, tsne_train_output)
19 print('Time taken to perform Hyperparameter tuning :',datetime.now()-start)
20
21 #Getting the best hyperparameter tuned model
22 best_model=rsearch.best_estimator_
23 print("Best estimator: ",best_model)
24
25 #Fitting the best model to our training data
26 #best_model.fit(df_train, tsne_train_output)
```

```
In [58]: 1 best_xgb = XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1
2         colsample_bytree=0.8, gamma=0.5, learning_rate=0.1,
3         max_delta_step=0, max_depth=9, min_child_weight=9, missing=None,
4         n_estimators=1500, n_jobs=-1, nthread=None, objective='reg:linear',
5         random_state=0, reg_alpha=100, reg_lambda=200, scale_pos_weight=1,
6         seed=None, silent=True, subsample=0.9)
7
8 best_xgb.fit(df_train,tsne_train_output)
```

```
Out[58]: XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
         colsample_bynode=1, colsample_bytree=0.8, gamma=0.5,
         importance_type='gain', learning_rate=0.1, max_delta_step=0,
         max_depth=9, min_child_weight=9, missing=None, n_estimators=1500,
         n_jobs=-1, nthread=None, objective='reg:linear', random_state=0,
         reg_alpha=100, reg_lambda=200, scale_pos_weight=1, seed=None,
         silent=True, subsample=0.9, verbosity=1)
```

```
In [0]: 1 #predicting with our trained Xg-Boost regressor
2 # the models x_model is already hyper parameter tuned
3 # the parameters that we got above are found using grid search
4
5 y_pred = best_xgb.predict(df_test)
6 xgb_test_predictions = [round(value) for value in y_pred]
7 y_pred = best_xgb.predict(df_train)
8 xgb_train_predictions = [round(value) for value in y_pred]
```

## Calculating the error metric values for various models

```
In [0]: 1 train_mape=[]
2 test_mape=[]
3
4 train_mape.append((mean_absolute_error(tsne_train_output,df_train['ft_1'].value
5 train_mape.append((mean_absolute_error(tsne_train_output,df_train['exp_avg'].value
6 train_mape.append((mean_absolute_error(tsne_train_output,rndf_train_predictions)
7 train_mape.append((mean_absolute_error(tsne_train_output, xgb_train_predictions)
8 train_mape.append((mean_absolute_error(tsne_train_output, pred_train))/(sum(tsne
9
10 test_mape.append((mean_absolute_error(tsne_test_output, df_test['ft_1'].value
11 test_mape.append((mean_absolute_error(tsne_test_output, df_test['exp_avg'].value
12 test_mape.append((mean_absolute_error(tsne_test_output, rndf_test_predictions)
13 test_mape.append((mean_absolute_error(tsne_test_output, xgb_test_predictions)
14 test_mape.append((mean_absolute_error(tsne_test_output,pred_test))/(sum(tsne
```

```
In [70]: 1 print ("Error Metric Matrix (Tree Based Regression Methods) - MAPE")
2 print ("-----")
3 print ("Baseline Model - Train: ",train_mape[0],")
4 print ("Exponential Averages Forecasting - Train: ",train_mape[1],")
5 print ("Linear Regression - Train: ",train_mape[4],")
6 print ("Random Forest Regression - Train: ",train_mape[2],")
```

Error Metric Matrix (Tree Based Regression Methods) - MAPE

```
-----
-----
Baseline Model - Train: 0.14870666996426116
Test: 0.14225522601041551
Exponential Averages Forecasting - Train: 0.14121603560900353
Test: 0.13490049942819257
Linear Regression - Train: 6.0207618934980696e+16
Test: 5.87620067570374e+16
Random Forest Regression - Train: 0.10861209273180993 T
est: 0.10484934300782964
```

## Error Metric Matrix

```
In [71]: 1 print ("Error Metric Matrix (Tree Based Regression Methods) - MAPE")
2 print ("-----")
3 print ("Baseline Model - Train: ",train_mape[0],")
4 print ("Exponential Averages Forecasting - Train: ",train_mape[1],")
5 print ("Linear Regression - Train: ",train_mape[4],")
6 print ("Random Forest Regression - Train: ",train_mape[2],")
7 print ("XgBoost Regression - Train: ",train_mape[3],")
8 print ("-----")
```

Error Metric Matrix (Tree Based Regression Methods) - MAPE

```
-----
-----
Baseline Model - Train: 0.14870666996426116
Test: 0.14225522601041551
Exponential Averages Forecasting - Train: 0.14121603560900353
Test: 0.13490049942819257
Linear Regression - Train: 6.0207618934980696e+16
Test: 5.87620067570374e+16
Random Forest Regression - Train: 0.10861209273180993 T
est: 0.10484934300782964
XgBoost Regression - Train: 0.07017864663012166
Test: 0.0875584710259383
-----
-----
```

In this case study the business problem we have undertaken was to predict the number of pickups in a region for a 10 min interval, the interpretability of the model is not important here from a perspective of end user, the error metric we have chosen here is mean absolute percentage error.

Then we clustered the data using segmentation i.e by using the latitude and longitude of pickups and then broke the regions using daily pickups in 10 min interval, the strategy we used here for clustering the data was that we wanted minimum inter cluster distance to be 0.5 and maximum



cluster distance to be 2 miles as that's how much a taxi travels with average speed in a 10 min interval, also we smoothed the data as for intervals where number of pickups tend to be zero we imputed the average over interval values

If we notice carefully some 10 minute intervals has 0 pick ups. 0 is not very relevant when we train the model. So instead of keeping those values as it is, we will smooth them. The process is as follows. We have time binned regions of 10 minute intervals. Let's take three continuous intervals and let's say each of the intervals have 50, 100 and 0 pickups respectively. Instead of taking 0 as the third value we will take the average of all the values and distribute it in the three regions as 50, 50, 50. In this way, each 30 minute interval will have the same number of pickups even though the pickups values changes in each of the three 10 minute interval.

**Time Series:** In the below graph, we can see repeated patterns like each day the number of cabs are highest during the office hours. It's least during midnight and increases as the day progresses. It decreases during noon time. Then again starts increasing during evening time. We will explore all such data dependencies using Fourier transformed features.

Fourier transform generally means decomposition of a wave into sum of multiple sine waves. Whenever there is repeated patterns in data, we can leverage the most out of them by using Fourier transform. Using Fourier transform we can decompose any given waveform (or a function) into its constituent frequencies. The graph obtained after Fourier transform will have unique frequencies and amplitudes corresponding to their most frequent occurrences in the original wave.

Each Fourier transformed feature is basically multiple sine waves. Each sine wave has a specific time period, frequency and amplitude. From the time period we can easily get the frequencies for each of the sine waves. We are basically converting the time series data from time domain to frequency domain. We will plot these frequencies and its corresponding amplitudes in a Fourier transformed graph. The X axis will represent the frequencies of the sine waves and the Y axis will represent the corresponding amplitudes.

Then we used High band pass filter + Wavenet denoising along with Holt's exponential smoothing as part of feature extracting process.

**I achieved the best MAPE with XGBRegressor - 0.0875**