# prabhudayala@gmail.com_NYC Final

December 1, 2019

## 1 Taxi demand prediction in New York City

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

     import pandas as pd#pandas to create small dataframes

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

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

     import time #Convert to unix time

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

     # matplotlib: used to plot graphs
     import matplotlib
     # matplotlib.use('nbagg') : matplotlib uses this protocall which makes plots␣
      ↪more user intractive like zoom in and zoom out
     matplotlib.use('nbagg')
     import matplotlib.pylab as plt
     import seaborn as sns#Plots
     from matplotlib import rcParams#Size of plots
```

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

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

# download migwin: https://mingw-w64.org/doku.php/download/mingw-builds
# install it in your system and keep the path, migw_path ='installed path'
mingw_path = 'C:\\Program Files\\mingw-w64\\x86_64-5.3.
 ↪0-posix-seh-rt_v4-rev0\\mingw64\\bin'
os.environ['PATH'] = mingw_path + ';' + os.environ['PATH']

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

# to install sklearn: pip install -U scikit-learn
from sklearn.ensemble import RandomForestRegressor
from sklearn.metrics import mean_squared_error
from sklearn.metrics import mean_absolute_error

from sklearn.model_selection import RandomizedSearchCV
from sklearn.model_selection import GridSearchCV

import warnings
warnings.filterwarnings("ignore")
```

## 2  Data Information

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

### 2.1  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 consid-

ered 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

# 3 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

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

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

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

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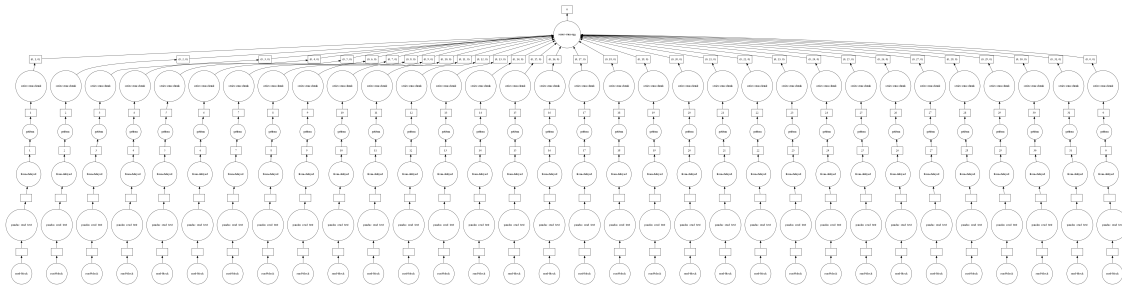
```
month.visualize()
```

[3]:



[4]:
```
month.fare_amount.sum().visualize()
```

[4]:



## 3.1 Features in the dataset:

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

Newark

Nassau or Westchester

Negotiated fare

Group ride

Store_and_fwd_flag

This flag indicates whether the trip record was held in vehicle memory before sending to the vendor,<br> aka "store and forward," because the vehicle did not have a connection to the server. <br>Y= store and forward trip <br>N= not a store and forward trip

Dropoff_longitude

Longitude where the meter was disengaged.

Dropoff_ latitude

Latitude where the meter was disengaged.

Payment_type

A numeric code signifying how the passenger paid for the trip.

Credit card

Cash

No charge

Dispute

Unknown

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.

# 4   ML Problem Formulation

Time-series forecasting and Regression

- To find number of pickups, given location cordinates(latitude and longitude) and time, in the query reigion 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.

# 5 Performance metrics

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

## 5.1 Data Cleaning

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

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

[5]:
```
   VendorID tpep_pickup_datetime tpep_dropoff_datetime  passenger_count  \
0         2  2015-01-15 19:05:39   2015-01-15 19:23:42                1
1         1  2015-01-10 20:33:38   2015-01-10 20:53:28                1
2         1  2015-01-10 20:33:38   2015-01-10 20:43:41                1
3         1  2015-01-10 20:33:39   2015-01-10 20:35:31                1
4         1  2015-01-10 20:33:39   2015-01-10 20:52:58                1

   trip_distance  pickup_longitude  pickup_latitude  RateCodeID  \
0           1.59        -73.993896        40.750111           1
1           3.30        -74.001648        40.724243           1
2           1.80        -73.963341        40.802788           1
3           0.50        -74.009087        40.713818           1
4           3.00        -73.971176        40.762428           1

   store_and_fwd_flag  dropoff_longitude  dropoff_latitude  payment_type  \
0                   N         -73.974785         40.750618             1
1                   N         -73.994415         40.759109             1
2                   N         -73.951820         40.824413             2
3                   N         -74.004326         40.719986             2
4                   N         -74.004181         40.742653             2

   fare_amount  extra  mta_tax  tip_amount  tolls_amount  \
0         12.0    1.0      0.5        3.25           0.0
1         14.5    0.5      0.5        2.00           0.0
2          9.5    0.5      0.5        0.00           0.0
3          3.5    0.5      0.5        0.00           0.0
4         15.0    0.5      0.5        0.00           0.0

   improvement_surcharge  total_amount
0                    0.3         17.05
1                    0.3         17.80
2                    0.3         10.80
3                    0.3          4.80
4                    0.3         16.30
```

### 5.1.1 1. Pickup Latitude and Pickup Longitude

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

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

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

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

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

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

```
[6]: <folium.folium.Map at 0x2123d95f940>
```

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

### 5.1.2 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 cordinates(lat,long) - (40.5774, -74.15) & (40.9176,-73.7004) so hence any cordinates not within these cordinates are not considered by us as we are only concerned with dropoffs which are within New York.

```
[7]: # Plotting dropoff cordinates which are outside the bounding box of New-York
     # we will collect all the points outside the bounding box of newyork city to␣
     ↪outlier_locations
```

```
outlier_locations = month[((month.dropoff_longitude <= -74.15) | (month.
 ↪dropoff_latitude <= 40.5774)| \
                    (month.dropoff_longitude >= -73.7004) | (month.
 ↪dropoff_latitude >= 40.9176))]

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

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

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

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

[7]: `<folium.folium.Map at 0x2123d96c5f8>`

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

### 5.1.3  3. Trip Durations:

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

[8]:
```
#The timestamps are converted to unix so as to get duration(trip-time) & speed␣
 ↪also pickup-times in unix are used while binning

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




# we return a data frame which contains the columns
# 1.'passenger_count' : self explanatory
# 2.'trip_distance' : self explanatory
```

```python
# 3.'pickup_longitude' : self explanatory
# 4.'pickup_latitude' : self explanatory
# 5.'dropoff_longitude' : self explanatory
# 6.'dropoff_latitude' : self explanatory
# 7.'total_amount' : total fair that was paid
# 8.'trip_times' : duration of each trip
# 9.'pickup_times : pickup time converted into unix time
# 10.'Speed' : velocity of each trip
def return_with_trip_times(month):
    duration = month[['tpep_pickup_datetime','tpep_dropoff_datetime']].compute()
    print(type(duration))
    #pickups and dropoffs to unix time
    duration_pickup = [convert_to_unix(x) for x in␣
 ↪duration['tpep_pickup_datetime'].values]
    duration_drop = [convert_to_unix(x) for x in␣
 ↪duration['tpep_dropoff_datetime'].values]
    #calculate duration of trips
    durations = (np.array(duration_drop) - np.array(duration_pickup))/float(60)

    #append durations of trips and speed in miles/hr to a new dataframe
    new_frame =␣
 ↪month[['passenger_count','trip_distance','pickup_longitude','pickup_latitude','dropoff_long
 ↪compute()

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

    return new_frame

# print(frame_with_durations.head())
#  ␣
 ↪passenger_count        trip_distance        pickup_longitude        pickup_latitude
#   1                1.59                -73.993896                40.750111    ␣
 ↪       -73.974785                40.750618                17.05            ␣
 ↪18.050000        1.421329e+09        5.285319
#   1                 3.30                -74.001648                40.724243  ␣
 ↪       -73.994415                40.759109                17.80  ␣
 ↪        19.833333        1.420902e+09        9.983193
#   1                 1.80                -73.963341                40.802788 ␣
 ↪        -73.951820                40.824413                10.80  ␣
 ↪        10.050000        1.420902e+09        10.746269
#   1                 0.50                -74.009087                40.713818 ␣
 ↪        -74.004326                40.719986                4.80  ␣
 ↪        1.866667        1.420902e+09        16.071429
```

```
#    1                    3.00              -73.971176                    40.762428  ␣
↪        -74.004181              40.742653                    16.30  ␣
↪        19.316667        1.420902e+09        9.318378
frame_with_durations = return_with_trip_times(month)
```

<class 'pandas.core.frame.DataFrame'>

[9]:
```
# the skewed box plot shows us the presence of outliers
sns.boxplot(y="trip_times", data =frame_with_durations)
plt.show()
```

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

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

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

[11]:
```
#looking further from the 99th percecntile
for i in range(90,100):
    var =frame_with_durations["trip_times"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/
↪100))]))
print ("100 percentile value is ",var[-1])
```

```
90 percentile value is 23.45
91 percentile value is 24.35
92 percentile value is 25.383333333333
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 34.46666666666667
98 percentile value is 38.71666666666667
99 percentile value is 46.75
100 percentile value is  548555.6333333333
```

[12]:
```python
#removing data based on our analysis and TLC regulations
frame_with_durations_modified=frame_with_durations[(frame_with_durations.
 ↪trip_times>1) & (frame_with_durations.trip_times<720)]
```

[13]:
```python
#box-plot after removal of outliers
sns.boxplot(y="trip_times", data =frame_with_durations_modified)
plt.show()
```

```
<IPython.core.display.Javascript object>
```

```
<IPython.core.display.HTML object>
```

[14]:
```python
#pdf of trip-times after removing the outliers
sns.FacetGrid(frame_with_durations_modified,size=6) \
      .map(sns.kdeplot,"trip_times") \
      .add_legend();
plt.show()
```

```
<IPython.core.display.Javascript object>
```

```
<IPython.core.display.HTML object>
```

[15]:
```python
#converting the values to log-values to chec for log-normal
import math
frame_with_durations_modified['log_times']=[math.log(i) for i in␣
 ↪frame_with_durations_modified['trip_times'].values]
```

[16]:
```python
#pdf of log-values
sns.FacetGrid(frame_with_durations_modified,size=6) \
      .map(sns.kdeplot,"log_times") \
      .add_legend();
plt.show();
```

```
<IPython.core.display.Javascript object>
```

```
<IPython.core.display.HTML object>
```

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

```
<IPython.core.display.Javascript object>
```

```
<IPython.core.display.HTML object>
```

### 5.1.4   4. Speed

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

```
<IPython.core.display.Javascript object>
```

```
<IPython.core.display.HTML object>
```

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

```
0 percentile value is 0.0
10 percentile value is 6.409495548961425
20 percentile value is 7.80952380952381
30 percentile value is 8.929133858267717
40 percentile value is 9.98019801980198
```

```
50 percentile value is 11.06865671641791
60 percentile value is 12.286689419795222
70 percentile value is 13.796407185628745
80 percentile value is 15.963224893917962
90 percentile value is 20.186915887850468
100 percentile value is  192857142.85714284
```

[20]:
```python
#calculating speed values at each percntile 90,91,92,93,94,95,96,97,98,99,100
for i in range(90,100):
    var =frame_with_durations_modified["Speed"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/
  ↪100))]))
print("100 percentile value is ",var[-1])
```

```
90 percentile value is 20.186915887850468
91 percentile value is 20.91645569620253
92 percentile value is 21.752988047808763
93 percentile value is 22.721893491124263
94 percentile value is 23.844155844155843
95 percentile value is 25.182552504038775
96 percentile value is 26.80851063829787
97 percentile value is 28.84304932735426
98 percentile value is 31.591128254580514
99 percentile value is 35.7513566847558
100 percentile value is  192857142.85714284
```

[21]:
```python
#calculating speed values at each percntile 99.0,99.1,99.2,99.3,99.4,99.5,99.
  ↪6,99.7,99.8,99.9,100
for i in np.arange(0.0, 1.0, 0.1):
    var =frame_with_durations_modified["Speed"].values
    var = np.sort(var,axis = None)
    print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/
  ↪100))]))
print("100 percentile value is ",var[-1])
```

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

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

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

[23]: 12.450173996027528

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

### 5.1.5  4. Trip Distance

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

```
<IPython.core.display.Javascript object>
```

```
<IPython.core.display.HTML object>
```

```
[25]: #calculating trip distance values at each percntile␣
       ↪0,10,20,30,40,50,60,70,80,90,100
      for i in range(0,100,10):
          var =frame_with_durations_modified["trip_distance"].values
          var = np.sort(var,axis = None)
          print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/
       ↪100))]))
      print("100 percentile value is ",var[-1])
```

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

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

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

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

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

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

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

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

### 5.1.6   5. Total Fare

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

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

```
[31]: #calculating total fare amount values at each percntile␣
      ↪0,10,20,30,40,50,60,70,80,90,100
      for i in range(0,100,10):
          var = frame_with_durations_modified["total_amount"].values
          var = np.sort(var,axis = None)
          print("{} percentile value is {}".format(i,var[int(len(var)*(float(i)/
      ↪100))]))
      print("100 percentile value is ",var[-1])
```

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

18

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

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

```
[33]:  #calculating total fare amount values at each percntile 99.0,99.1,99.2,99.3,99.
       ↪4,99.5,99.6,99.7,99.8,99.9,100
       for i in np.arange(0.0, 1.0, 0.1):
           var = frame_with_durations_modified["total_amount"].values
           var = np.sort(var,axis = None)
           print("{} percentile value is {}".format(99+i,var[int(len(var)*(float(99+i)/
       ↪100))]))
       print("100 percentile value is ",var[-1])
```

```
99.0 percentile value is 66.13
99.1 percentile value is 68.13
99.2 percentile value is 69.6
99.3 percentile value is 69.6
99.4 percentile value is 69.73
99.5 percentile value is 69.75
99.6 percentile value is 69.76
99.7 percentile value is 72.58
99.8 percentile value is 75.35
99.9 percentile value is 88.28
100 percentile value is  3950611.6
```

Observation:- As even the 99.9th percentile value doesnt 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 analyis

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

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

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

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

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

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

## 5.2   Remove all outliers/erronous points.

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


           a = new_frame.shape[0]
           print ("Number of pickup records = ",a)
           temp_frame = new_frame[((new_frame.dropoff_longitude >= -74.15) &␣
       ↪(new_frame.dropoff_longitude <= -73.7004) &\
```

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


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


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

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

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


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

    new_frame = new_frame[(new_frame.trip_times > 0) & (new_frame.trip_times <␣
↪720)]
    new_frame = new_frame[(new_frame.trip_distance > 0) & (new_frame.
↪trip_distance < 23)]
    new_frame = new_frame[(new_frame.Speed < 45.31) & (new_frame.Speed > 0)]
```

```
    new_frame = new_frame[(new_frame.total_amount <1000) & (new_frame.
 ↪total_amount >0)]

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

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

```
Removing outliers in the month of Jan-2015
----
Number of pickup records =  12748986
Number of outlier coordinates lying outside NY boundaries: 293919
Number of outliers from trip times analysis: 23889
Number of outliers from trip distance analysis: 92597
Number of outliers from speed analysis: 24473
Number of outliers from fare analysis: 5275
Total outliers removed 377910
---
fraction of data points that remain after removing outliers 0.9703576425607495
```

## 6 Data-preperation

### 6.1 Clustering/Segmentation

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

      def find_min_distance(cluster_centers, cluster_len):
          nice_points = 0
          wrong_points = 0
          less2 = []
          more2 = []
          min_dist=1000
          for i in range(0, cluster_len):
              nice_points = 0
              wrong_points = 0
              for j in range(0, cluster_len):
                  if j!=i:
                      distance = gpxpy.geo.haversine_distance(cluster_centers[i][0],␣
        ↪cluster_centers[i][1],cluster_centers[j][0], cluster_centers[j][1])
```

```
                min_dist = min(min_dist,distance/(1.60934*1000))
                if (distance/(1.60934*1000)) <= 2:
                    nice_points +=1
                else:
                    wrong_points += 1
        less2.append(nice_points)
        more2.append(wrong_points)
    neighbours.append(less2)
    print ("On choosing a cluster size of ",cluster_len,"\nAvg. Number of␣
 ↪Clusters within the vicinity (i.e. intercluster-distance < 2):", np.
 ↪ceil(sum(less2)/len(less2)), "\nAvg. Number of Clusters outside the vicinity␣
 ↪(i.e. intercluster-distance > 2):", np.ceil(sum(more2)/len(more2)),"\nMin␣
 ↪inter-cluster distance = ",min_dist,"\n---")

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

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

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

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

### 6.1.1 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

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

      # Getting 40 clusters using the kmeans
      kmeans = MiniBatchKMeans(n_clusters=30, batch_size=10000,random_state=0).
      →fit(coords)
      frame_with_durations_outliers_removed['pickup_cluster'] = kmeans.
      →predict(frame_with_durations_outliers_removed[['pickup_latitude',
      →'pickup_longitude']])
```

### 6.1.2 Plotting the cluster centers:

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

```
[42]: <folium.folium.Map at 0x2130abb7208>
```

### 6.1.3 Plotting the clusters:

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

      plot_clusters(frame_with_durations_outliers_removed)
```

```
<IPython.core.display.Javascript object>


<IPython.core.display.HTML object>
```

## 6.2 Time-binning

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

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


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

          start_pickup_unix=unix_times[year-2015][month-1]
          # https://www.timeanddate.com/time/zones/est
          # (int((i-start_pickup_unix)/600)+33) : our unix time is in gmt to we are␣
       ↪converting it to est
          tenminutewise_binned_unix_pickup_times=[(int((i-start_pickup_unix)/600)+33)␣
       ↪for i in unix_pickup_times]
          frame['pickup_bins'] = np.array(tenminutewise_binned_unix_pickup_times)
          return frame
```

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

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

```
[46]:    passenger_count  trip_distance  pickup_longitude  pickup_latitude  \
      0                1           1.59        -73.993896        40.750111
      1                1           3.30        -74.001648        40.724243
      2                1           1.80        -73.963341        40.802788
      3                1           0.50        -74.009087        40.713818
      4                1           3.00        -73.971176        40.762428

         dropoff_longitude  dropoff_latitude  total_amount  trip_times  \
      0         -73.974785         40.750618         17.05   18.050000
      1         -73.994415         40.759109         17.80   19.833333
      2         -73.951820         40.824413         10.80   10.050000
      3         -74.004326         40.719986          4.80    1.866667
      4         -74.004181         40.742653         16.30   19.316667

           pickup_times       Speed  pickup_cluster  pickup_bins
      0   1.421329e+09    5.285319              14         2130
      1   1.420902e+09    9.983193              25         1419
      2   1.420902e+09   10.746269               8         1419
      3   1.420902e+09   16.071429              21         1419
      4   1.420902e+09    9.318378              28         1419
```

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

```
[47]:                               trip_distance
      pickup_cluster pickup_bins
      0              1                       138
                     2                       262
                     3                       311
                     4                       326
                     5                       381
```

```
[48]: # upto now we cleaned data and prepared data for the month 2015,

      # now do the same operations for months Jan, Feb, March of 2016
      # 1. get the dataframe which inlcudes only required colums
      # 2. adding trip times, speed, unix time stamp of pickup_time
      # 4. remove the outliers based on trip_times, speed, trip_duration,
      →total_amount
      # 5. add pickup_cluster to each data point
```

```
# 6. add pickup_bin (index of 10min intravel to which that trip belongs to)
# 7. group by data, based on 'pickup_cluster' and 'pickuo_bin'

# Data Preparation for the months of Jan,Feb and March 2016
def datapreparation(month,kmeans,month_no,year_no):

    print ("Return with trip times..")

    frame_with_durations = return_with_trip_times(month)

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

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

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

    return final_updated_frame,final_groupby_frame

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

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

```
Return with trip times..
<class 'pandas.core.frame.DataFrame'>
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..
<class 'pandas.core.frame.DataFrame'>
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..
<class 'pandas.core.frame.DataFrame'>
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..
```

[49]: `jan_2016_frame.head()`

[49]:

|   | passenger_count | trip_distance | pickup_longitude | pickup_latitude \ |
|---|---|---|---|---|
| 5 | 2 | 5.52 | -73.980118 | 40.743050 |
| 6 | 2 | 7.45 | -73.994057 | 40.719990 |
| 7 | 1 | 1.20 | -73.979424 | 40.744614 |
| 8 | 1 | 6.00 | -73.947151 | 40.791046 |
| 9 | 1 | 3.21 | -73.998344 | 40.723896 |

|   | dropoff_longitude | dropoff_latitude | total_amount | trip_times \ |
|---|---|---|---|---|
| 5 | -73.913490 | 40.763142 | 20.3 | 18.50 |
| 6 | -73.966362 | 40.789871 | 27.3 | 26.75 |
| 7 | -73.992035 | 40.753944 | 10.3 | 11.90 |
| 8 | -73.920769 | 40.865578 | 19.3 | 11.20 |
| 9 | -73.995850 | 40.688400 | 12.8 | 11.10 |

```
     pickup_times       Speed  pickup_cluster  pickup_bins
5    1.451587e+09    17.902703              20            0
6    1.451587e+09    16.710280              11            0
7    1.451587e+09     6.050420              20            1
8    1.451587e+09    32.142857              26            1
9    1.451587e+09    17.351351              25            1
```

## 6.3 Smoothing

```python
[50]: # Gets the unique bins where pickup values are present for each each reigion

      # for each cluster region we will collect all the indices of 10min intravels in␣
       ↪which the pickups are happened
      # we got an observation that there are some pickpbins that doesnt have any␣
       ↪pickups
      def return_unq_pickup_bins(frame):
          values = []
          for i in range(0,30):
              new = frame[frame['pickup_cluster'] == i]
              list_unq = list(set(new['pickup_bins']))
              list_unq.sort()
              values.append(list_unq)
          return values
```

```python
[51]: # for every month we get all indices of 10min intravels in which atleast one␣
       ↪pickup got happened

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

      #feb
      feb_2016_unique = return_unq_pickup_bins(feb_2016_frame)

      #march
      mar_2016_unique = return_unq_pickup_bins(mar_2016_frame)
```

```python
[52]: # for each cluster number of 10min intravels with 0 pickups
      for i in range(30):
          print("for the ",i,"th cluster number of 10min intavels with zero pickups:␣
       ↪",4464 - len(set(jan_2015_unique[i])))
          print('-'*60)
```

```
for the   0 th cluster number of 10min intavels with zero pickups:   26
------------------------------------------------------------
for the   1 th cluster number of 10min intavels with zero pickups:   30
------------------------------------------------------------
```

30

```
for the  2 th cluster number of 10min intavels with zero pickups:  150
--------------------------------------------------------------
for the  3 th cluster number of 10min intavels with zero pickups:  35
--------------------------------------------------------------
for the  4 th cluster number of 10min intavels with zero pickups:  170
--------------------------------------------------------------
for the  5 th cluster number of 10min intavels with zero pickups:  40
--------------------------------------------------------------
for the  6 th cluster number of 10min intavels with zero pickups:  320
--------------------------------------------------------------
for the  7 th cluster number of 10min intavels with zero pickups:  35
--------------------------------------------------------------
for the  8 th cluster number of 10min intavels with zero pickups:  39
--------------------------------------------------------------
for the  9 th cluster number of 10min intavels with zero pickups:  46
--------------------------------------------------------------
for the  10 th cluster number of 10min intavels with zero pickups:  98
--------------------------------------------------------------
for the  11 th cluster number of 10min intavels with zero pickups:  32
--------------------------------------------------------------
for the  12 th cluster number of 10min intavels with zero pickups:  37
--------------------------------------------------------------
for the  13 th cluster number of 10min intavels with zero pickups:  326
--------------------------------------------------------------
for the  14 th cluster number of 10min intavels with zero pickups:  35
--------------------------------------------------------------
for the  15 th cluster number of 10min intavels with zero pickups:  29
--------------------------------------------------------------
for the  16 th cluster number of 10min intavels with zero pickups:  25
--------------------------------------------------------------
for the  17 th cluster number of 10min intavels with zero pickups:  40
--------------------------------------------------------------
for the  18 th cluster number of 10min intavels with zero pickups:  30
--------------------------------------------------------------
for the  19 th cluster number of 10min intavels with zero pickups:  35
--------------------------------------------------------------
for the  20 th cluster number of 10min intavels with zero pickups:  40
--------------------------------------------------------------
for the  21 th cluster number of 10min intavels with zero pickups:  38
--------------------------------------------------------------
for the  22 th cluster number of 10min intavels with zero pickups:  34
--------------------------------------------------------------
for the  23 th cluster number of 10min intavels with zero pickups:  49
--------------------------------------------------------------
for the  24 th cluster number of 10min intavels with zero pickups:  49
--------------------------------------------------------------
for the  25 th cluster number of 10min intavels with zero pickups:  27
--------------------------------------------------------------
```

```
for the  26 th cluster number of 10min intavels with zero pickups:  26
-----------------------------------------------------------
for the  27 th cluster number of 10min intavels with zero pickups:  720
-----------------------------------------------------------
for the  28 th cluster number of 10min intavels with zero pickups:  34
-----------------------------------------------------------
for the  29 th cluster number of 10min intavels with zero pickups:  29
-----------------------------------------------------------
```

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)

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

      # for every 10min intravel(pickup_bin) we will check it is there in our unique
       ↪bin,
      # if it is there we will add the count_values[index] to smoothed data
      # if not we add 0 to the smoothed data
      # we finally return smoothed data
      def fill_missing(count_values,values):
          smoothed_regions=[]
          ind=0
          for r in range(0,30):
              smoothed_bins=[]
              for i in range(4464):
                  if i in values[r]:
                      smoothed_bins.append(count_values[ind])
                      ind+=1
                  else:
                      smoothed_bins.append(0)
              smoothed_regions.extend(smoothed_bins)
          return smoothed_regions
```

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

```python
# for every 10min intravel(pickup_bin) we will check it is there in our unique
 →bin,
# if it is there we will add the count_values[index] to smoothed data
# if not we add smoothed data (which is calculated based on the methods that
 →are discussed in the above markdown cell)
# we finally return smoothed data
def smoothing(count_values,values):
    smoothed_regions=[] # stores list of final smoothed values of each region
    ind=0
    repeat=0
    smoothed_value=0
    for r in range(0,30):
        smoothed_bins=[] #stores the final smoothed values
        repeat=0
        for i in range(4464):
            if repeat!=0: # prevents iteration for a value which is already
→visited/resolved
                repeat-=1
                continue
            if i in values[r]: #checks if the pickup-bin exists
                smoothed_bins.append(count_values[ind]) # appends the value of
→the pickup bin if it exists
            else:
                if i!=0:
                    right_hand_limit=0
                    for j in range(i,4464):
                        if  j not in values[r]: #searches for the left-limit or
→the pickup-bin value which has a pickup value
                            continue
                        else:
                            right_hand_limit=j
                            break
                    if right_hand_limit==0:
                    #Case 1: When we have the last/last few values are found to
→be missing,hence we have no right-limit here
                        smoothed_value=count_values[ind-1]*1.0/((4463-i)+2)*1.0
                        for j in range(i,4464):
                            smoothed_bins.append(math.ceil(smoothed_value))
                        smoothed_bins[i-1] = math.ceil(smoothed_value)
                        repeat=(4463-i)
                        ind-=1
                    else:
                    #Case 2: When we have the missing values between two known
→values
```

```
                    ⎵
→smoothed_value=(count_values[ind-1]+count_values[ind])*1.0/
→((right_hand_limit-i)+2)*1.0
                        for j in range(i,right_hand_limit+1):
                                smoothed_bins.append(math.ceil(smoothed_value))
                        smoothed_bins[i-1] = math.ceil(smoothed_value)
                        repeat=(right_hand_limit-i)
                    else:
                        #Case 3: When we have the first/first few values are found⎵
→to be missing,hence we have no left-limit here
                        right_hand_limit=0
                        for j in range(i,4464):
                            if  j not in values[r]:
                                continue
                            else:
                                right_hand_limit=j
                                break
                        smoothed_value=count_values[ind]*1.0/
→((right_hand_limit-i)+1)*1.0
                        for j in range(i,right_hand_limit+1):
                                smoothed_bins.append(math.ceil(smoothed_value))
                        repeat=(right_hand_limit-i)
                ind+=1
            smoothed_regions.extend(smoothed_bins)
        return smoothed_regions
```

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

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

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

```
number of 10min intravels among all the clusters  133920
```

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

<IPython.core.display.Javascript object>


<IPython.core.display.HTML object>

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

      # Ans: consider we have data of some month in 2015 jan 1st, 10 _ _ _ 20, i.e␣
       ↪there are 10 pickups that are happened in 1st
      # 10st 10min intravel, 0 pickups happened in 2nd 10mins intravel, 0 pickups␣
       ↪happened in 3rd 10min intravel
      # and 20 pickups happened in 4th 10min intravel.
      # in fill_missing method we replace these values like 10, 0, 0, 20
      # where as in smoothing method we replace these values as 6,6,6,6,6, if you can␣
       ↪check the number of pickups
      # that are happened in the first 30min are same in both cases, but if you can␣
       ↪observe that we looking at the future values
      # wheen you are using smoothing we are looking at the future number of pickups␣
       ↪which might cause a data leakage.

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

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

```
[60]: # Making list of all the values of pickup data in every bin for a period of 3␣
       ↪months and storing them region-wise
      regions_cum = []
```

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

# number of 10min indices for jan 2015= 24*31*60/10 = 4464
# number of 10min indices for jan 2016 = 24*31*60/10 = 4464
# number of 10min indices for feb 2016 = 24*29*60/10 = 4176
# number of 10min indices for march 2016 = 24*31*60/10 = 4464
# regions_cum: it will contain 30 lists, each list will contain 4464+4176+4464␣
 ↪values which represents the number of pickups
# that are happened for three months in 2016 data

for i in range(0,30):
    regions_cum.append(jan_2016_smooth[4464*i:
 ↪4464*(i+1)]+feb_2016_smooth[4176*i:4176*(i+1)]+mar_2016_smooth[4464*i:
 ↪4464*(i+1)])
    #
# print(len(regions_cum))
# 30
#print(len(regions_cum[0]))
# 13104
```

## 6.4   Time series and Fourier Transforms

```
[61]: def uniqueish_color():
          """There're better ways to generate unique colors, but this isn't awful."""
          return plt.cm.gist_ncar(np.random.random())
      first_x = list(range(0,4464))
      second_x = list(range(4464,8640))
      third_x = list(range(8640,13104))
      for i in range(30):
          plt.figure(figsize=(10,4))
          plt.plot(first_x,regions_cum[i][:4464], color=uniqueish_color(),␣
       ↪label='2016 Jan month data')
          plt.plot(second_x,regions_cum[i][4464:8640], color=uniqueish_color(),␣
       ↪label='2016 feb month data')
          plt.plot(third_x,regions_cum[i][8640:], color=uniqueish_color(),␣
       ↪label='2016 march month data')
          plt.legend()
          plt.show()
```

```
<IPython.core.display.Javascript object>
```

```
<IPython.core.display.HTML object>
```

```
<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>
```

```
<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>
```

```
<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>

<IPython.core.display.HTML object>

<IPython.core.display.Javascript object>
```

```
<IPython.core.display.HTML object>


<IPython.core.display.Javascript object>


<IPython.core.display.HTML object>


<IPython.core.display.Javascript object>


<IPython.core.display.HTML object>


<IPython.core.display.Javascript object>


<IPython.core.display.HTML object>
```

```
[62]: # getting peaks: https://blog.ytotech.com/2015/11/01/findpeaks-in-python/
      # read more about fft function : https://docs.scipy.org/doc/numpy/reference/
       ↪generated/numpy.fft.fft.html
      Y    = np.fft.fft(np.array(jan_2016_smooth)[0:4464])
      # read more about the fftfreq: https://docs.scipy.org/doc/numpy/reference/
       ↪generated/numpy.fft.fftfreq.html
      freq = np.fft.fftfreq(4464, 1)
      n = len(freq)
      plt.figure()
      plt.plot( freq[:int(n/2)], np.abs(Y)[:int(n/2)] )
      plt.xlabel("Frequency")
      plt.ylabel("Amplitude")
      plt.show()
```

```
<IPython.core.display.Javascript object>


<IPython.core.display.HTML object>
```

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

```
[64]: final_fft=[]
      #final_fft_frequency=[]
      freq = np.fft.fftfreq(4464, 1)
      for i in range(0,30):
          temp=[]
          #temp_frequency =[]
          for i,j in enumerate(regions_cum[i]):
              if(i<=4464):
                  list_is = [0]*4464
              else:
                  list_is = regions_cum[0][i-4464:i]
              Y = np.fft.fft(np.array(list_is))
              temp.append(np.abs(Y)[np.argsort(np.abs(Y))[-5:]])
          #temp_frequency.append(freq)
      #    print(len(temp))
      #    print(temp[0])
      #    print(temp[1])
          final_fft.append(temp)
          #final_fft_frequency.append(temp_frequency)
      print(len(final_fft))
```

```
30
```

## 6.5 Double Exponential Smoothing

```
[65]: #https://grisha.org/blog/2016/02/16/
      ↪triple-exponential-smoothing-forecasting-part-ii/
      def double_exponential_smoothing(series, alpha, beta):
          result = [series[0]]
          for n in range(1, len(series)+1):
              #print(series[0])
              if n == 1:
                  level, trend = series[0], series[1] - series[0]
              if n >= len(series): # we are forecasting
                  value = result[-1]
              else:
                  value = series[n]
              #print(level, alpha*value,(1-alpha),(level+trend))
              last_level, level = level, alpha*value + (1-alpha)*(level+trend)
              trend = beta*(level-last_level) + (1-beta)*trend
              result.append(level+trend)
          return result
```

```
[66]: double_exponential_smoothing_feat = []
      for i in range(0,30):
```

```
        double_exponential_smoothing_feat.
     ↪append(double_exponential_smoothing(regions_cum[i], alpha=0.9, beta=0.9)[:
     ↪-1])
```

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

```python
[68]: def initial_seasonal_components(series, slen):
          seasonals = {}
          season_averages = []
          n_seasons = int(len(series)/slen)
          # compute season averages
          for j in range(n_seasons):
              season_averages.append(sum(series[slen*j:slen*j+slen])/float(slen))
          # compute initial values
          for i in range(slen):
              sum_of_vals_over_avg = 0.0
              for j in range(n_seasons):
                  sum_of_vals_over_avg += series[slen*j+i]-season_averages[j]
              seasonals[i] = sum_of_vals_over_avg/n_seasons
          return seasonals
```

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

```python
[70]: triple_exponential_smoothing_fet = []
      for i in range(0,30):
```

```
    triple_exponential_smoothing_fet.
 ↪append(triple_exponential_smoothing(regions_cum[i], 12, 0.716, 0.029, 0.993,␣
 ↪0))
```

## 6.6 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 2. Using Previous known values of the 2016 data itself to predict the future values

### 6.6.1 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 -

```python
[71]: def MA_R_Predictions(ratios,month):
          predicted_ratio=(ratios['Ratios'].values)[0]
          error=[]
          predicted_values=[]
          window_size=3
          predicted_ratio_values=[]
          for i in range(0,4464*30):
              if i%4464==0:
                  predicted_ratio_values.append(0)
                  predicted_values.append(0)
                  error.append(0)
                  continue
              predicted_ratio_values.append(predicted_ratio)
              predicted_values.append(int(((ratios['Given'].
     ↪values)[i])*predicted_ratio))
              error.append(abs((math.pow(int(((ratios['Given'].
     ↪values)[i])*predicted_ratio)-(ratios['Prediction'].values)[i],1))))
              if i+1>=window_size:
                  predicted_ratio=sum((ratios['Ratios'].values)[(i+1)-window_size:
     ↪(i+1)])/window_size
              else:
                  predicted_ratio=sum((ratios['Ratios'].values)[0:(i+1)])/(i+1)


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

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

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

```python
[72]: def MA_P_Predictions(ratios,month):
          predicted_value=(ratios['Prediction'].values)[0]
          error=[]
          predicted_values=[]
          window_size=1
          predicted_ratio_values=[]
          for i in range(0,4464*30):
              predicted_values.append(predicted_value)
              error.append(abs((math.pow(predicted_value-(ratios['Prediction'].
      ↪values)[i],1))))
              if i+1>=window_size:
                  predicted_value=int(sum((ratios['Prediction'].
      ↪values)[(i+1)-window_size:(i+1)])/window_size)
              else:
                  predicted_value=int(sum((ratios['Prediction'].values)[0:(i+1)])/
      ↪(i+1))

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

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

### 6.6.2 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 -

```python
[73]: def WA_R_Predictions(ratios,month):
          predicted_ratio=(ratios['Ratios'].values)[0]
          alpha=0.5
          error=[]
          predicted_values=[]
          window_size=5
          predicted_ratio_values=[]
```

44

```python
    for i in range(0,4464*30):
        if i%4464==0:
            predicted_ratio_values.append(0)
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_ratio_values.append(predicted_ratio)
        predicted_values.append(int(((ratios['Given'].
↪values)[i])*predicted_ratio))
        error.append(abs((math.pow(int(((ratios['Given'].
↪values)[i])*predicted_ratio)-(ratios['Prediction'].values)[i],1))))
        if i+1>=window_size:
            sum_values=0
            sum_of_coeff=0
            for j in range(window_size,0,-1):
                sum_values += j*(ratios['Ratios'].values)[i-window_size+j]
                sum_of_coeff+=j
            predicted_ratio=sum_values/sum_of_coeff
        else:
            sum_values=0
            sum_of_coeff=0
            for j in range(i+1,0,-1):
                sum_values += j*(ratios['Ratios'].values)[j-1]
                sum_of_coeff+=j
            predicted_ratio=sum_values/sum_of_coeff

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

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

Weighted Moving Averages using Previous 2016 Values -

```python
[74]: def WA_P_Predictions(ratios,month):
    predicted_value=(ratios['Prediction'].values)[0]
    error=[]
    predicted_values=[]
    window_size=2
    for i in range(0,4464*30):
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].
↪values)[i],1))))
        if i+1>=window_size:
```

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

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

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

For the above the Hyperparameter is the window-size (n) which is tuned manually and it is found that the window-size of 2 is optimal for getting the best results using Weighted Moving Averages using previous 2016 values therefore we get

### 6.6.3   Exponential Weighted Moving Averages

https://en.wikipedia.org/wiki/Moving_average#Exponential_moving_average        Through weighted averaged we have satisfied the analogy of giving higher weights to the latest value and decreasing weights to the subsequent ones but we still do not know which is the correct weighting scheme as there are infinetly 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  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  then the number of days on which the value of the current iteration is based is~ i.e. we consider values 10 days prior before we predict the value for the current iteration. Also the weights are assigned using  ,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.

```
[75]: def EA_R1_Predictions(ratios,month):
          predicted_ratio=(ratios['Ratios'].values)[0]
          alpha=0.6
          error=[]
          predicted_values=[]
```

```python
    predicted_ratio_values=[]
    for i in range(0,4464*30):
        if i%4464==0:
            predicted_ratio_values.append(0)
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_ratio_values.append(predicted_ratio)
        predicted_values.append(int(((ratios['Given'].
↪values)[i])*predicted_ratio))
        error.append(abs((math.pow(int(((ratios['Given'].
↪values)[i])*predicted_ratio)-(ratios['Prediction'].values)[i],1))))
        predicted_ratio = (alpha*predicted_ratio) +␣
↪(1-alpha)*((ratios['Ratios'].values)[i])

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

```python
[76]: def EA_P1_Predictions(ratios,month):
    predicted_value= (ratios['Prediction'].values)[0]
    alpha=0.3
    error=[]
    predicted_values=[]
    for i in range(0,4464*30):
        if i%4464==0:
            predicted_values.append(0)
            error.append(0)
            continue
        predicted_values.append(predicted_value)
        error.append(abs((math.pow(predicted_value-(ratios['Prediction'].
↪values)[i],1))))
        predicted_value =int((alpha*predicted_value) +␣
↪(1-alpha)*((ratios['Prediction'].values)[i]))

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

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

### 6.7 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

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

```
Error Metric Matrix (Forecasting Methods) - MAPE & MSE
----------------------------------------------------------------------------
------------------------
Moving Averages (Ratios) -                            MAPE:  0.1628685972027342
MSE:   561.0487977897252
Moving Averages (2016 Values) -                       MAPE:
0.12651735674574424          MSE:   241.14901433691756
----------------------------------------------------------------------------
------------------------
Weighted Moving Averages (Ratios) -                   MAPE:  0.1599661761243253
```

48

```
MSE:   548.5285170250896
Weighted Moving Averages (2016 Values) -                    MAPE:
0.12121086157001072        MSE:   229.33734318996414
--------------------------------------------------------------------------------
------------------------
Exponential Moving Averages (Ratios) -                    MAPE:   0.1593403910652334
MSE:   546.5861260454002
Exponential Moving Averages (2016 Values) -                    MAPE:   0.1209639974233378
MSE:   226.0377688172043
```

Plese 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:- i.e Exponential Moving Averages using 2016 Values

## 6.8 Regression Models

### 6.8.1 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

```python
[79]: # Preparing data to be split into train and test, The below prepares data in␣
      ↪cumulative form which will be later split into test and train
      # number of 10min indices for jan 2015= 24*31*60/10 = 4464
      # number of 10min indices for jan 2016 = 24*31*60/10 = 4464
      # number of 10min indices for feb 2016 = 24*29*60/10 = 4176
      # number of 10min indices for march 2016 = 24*31*60/10 = 4464
      # regions_cum: it will contain 30 lists, each list will contain 4464+4176+4464␣
      ↪values which represents the number of pickups
      # that are happened for three months in 2016 data

      print(len(regions_cum))
      # 30
      print(len(regions_cum[0]))
      # 12960

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

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


      # tsne_lat will contain 13104-5=13099 times lattitude of cluster center for␣
      ↪every cluster
```

```python
# Ex: [[cent_lat 13099times],[cent_lat 13099times], [cent_lat 13099times]....␣
 ↪30 lists]
# it is list of lists
tsne_lat = []



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


# we will code each day
# sunday = 0, monday=1, tue = 2, wed=3, thur=4, fri=5,sat=6
# for every cluster we will be adding 13099 values, each value represent to␣
 ↪which day of the week that pickup bin belongs to
# it is list of lists
tsne_weekday = []


# its an numbpy array, of shape (523960, 5)
# each row corresponds to an entry in out data
# for the first row we will have [f0,f1,f2,f3,f4] fi=number of pickups happened␣
 ↪in i+1th 10min intraval(bin)
# the second row will have [f1,f2,f3,f4,f5]
# the third row will have [f2,f3,f4,f5,f6]
# and so on...
tsne_feature = []



tsne_feature = [0]*number_of_time_stamps
for i in range(0,30):
    tsne_lat.append([kmeans.cluster_centers_[i][0]]*13099)
    tsne_lon.append([kmeans.cluster_centers_[i][1]]*13099)
    # jan 1st 2016 is thursday, so we start our day from 4: "(int(k/144))%7+4"
    # our prediction start from 5th 10min intraval since we need to have number␣
 ↪of pickups that are happened in last 5 pickup bins
    tsne_weekday.append([int(((int(k/144))%7+4)%7) for k in␣
 ↪range(5,4464+4176+4464)])
    # regions_cum is a list of lists [[x1,x2,x3..x13104], [x1,x2,x3..x13104],␣
 ↪[x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3..x13104], .. 30 lsits]
    tsne_feature = np.vstack((tsne_feature, [regions_cum[i][r:
 ↪r+number_of_time_stamps] for r in␣
 ↪range(0,len(regions_cum[i])-number_of_time_stamps)]))
    output.append(regions_cum[i][5:])
tsne_feature = tsne_feature[1:]
```

```
30
13104
```

[80]: 
```
len(tsne_lat[0])*len(tsne_lat) == tsne_feature.shape[0] ==␣
 ↪len(tsne_weekday)*len(tsne_weekday[0]) == 30*13099 ==␣
 ↪len(output)*len(output[0])
```

[80]: True

[81]: 
```python
# Getting the predictions of exponential moving averages to be used as a␣
 ↪feature in cumulative form

# upto now we computed 8 features for every data point that starts from 50th␣
 ↪min of the day
# 1. cluster center lattitude
# 2. cluster center longitude
# 3. day of the week
# 4. f_t_1: number of pickups that are happened previous t-1th 10min intravel
# 5. f_t_2: number of pickups that are happened previous t-2th 10min intravel
# 6. f_t_3: number of pickups that are happened previous t-3th 10min intravel
# 7. f_t_4: number of pickups that are happened previous t-4th 10min intravel
# 8. f_t_5: number of pickups that are happened previous t-5th 10min intravel

# from the baseline models we said the exponential weighted moving avarage␣
 ↪gives us the best error
# we will try to add the same exponential weighted moving avarage at t as a␣
 ↪feature to our data
# exponential weighted moving avarage => p'(t) = alpha*p'(t-1) +␣
 ↪(1-alpha)*P(t-1)
alpha=0.3

# it is a temporary array that store exponential weighted moving avarage for␣
 ↪each 10min intravel,
# for each cluster it will get reset
# for every cluster it contains 13104 values
predicted_values=[]

# it is similar like tsne_lat
# it is list of lists
# predict_list is a list of lists [[x5,x6,x7..x13104], [x5,x6,x7..x13104],␣
 ↪[x5,x6,x7..x13104], [x5,x6,x7..x13104], [x5,x6,x7..x13104], .. 30 lsits]
predict_list = []
tsne_flat_exp_avg = []
for r in range(0,30):
    for i in range(0,13104):
        if i==0:
            predicted_value= regions_cum[r][0]
            predicted_values.append(0)
```

```
                continue
            predicted_values.append(predicted_value)
            predicted_value =int((alpha*predicted_value) +
    ↪(1-alpha)*(regions_cum[r][i]))
        predict_list.append(predicted_values[5:])
        predicted_values=[]
```

```
[82]: # train, test split : 70% 30% split
      # Before we start predictions using the tree based regression models we take 3
          ↪months of 2016 pickup data
      # and split it such that for every region we have 70% data in train and 30% in
          ↪test,
      # ordered date-wise for every region
      print("size of train data :", int(13099*0.7))
      print("size of test data :", int(13099*0.3))
```

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

```
[83]: # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps)
          ↪for our training data
      train_features =  [tsne_feature[i*13099:(13099*i+9169)] for i in range(0,30)]
      # temp = [0]*(12955 - 9068)
      test_features = [tsne_feature[(13099*(i))+9169:13099*(i+1)] for i in
          ↪range(0,30)]
```

```
[84]: # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps)
          ↪for our training data
      train_features =  [tsne_feature[i*13099:(13099*i+9169)] for i in range(0,30)]
      # temp = [0]*(12955 - 9068)
      test_features = [tsne_feature[(13099*(i))+9169:13099*(i+1)] for i in
          ↪range(0,30)]
```

```
[85]: print("Number of data clusters",len(train_features), "Number of data points in
          ↪trian data", len(train_features[0]), "Each data point contains",
          ↪len(train_features[0][0]),"features")
      print("Number of data clusters",len(train_features), "Number of data points in
          ↪test data", len(test_features[0]), "Each data point contains",
          ↪len(test_features[0][0]),"features")
```

```
Number of data clusters 30 Number of data points in trian data 9169 Each data
point contains 5 features
Number of data clusters 30 Number of data points in test data 3930 Each data
point contains 5 features
```

```
[86]: # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps)
          ↪for our training data
      tsne_train_flat_lat = [i[:9169] for i in tsne_lat]
```

```
tsne_train_flat_lon = [i[:9169] for i in tsne_lon]
tsne_train_flat_weekday = [i[:9169] for i in tsne_weekday]
tsne_train_flat_output = [i[:9169] for i in output]
tsne_train_flat_exp_avg = [i[:9169] for i in predict_list]
tsne_train_flat_fft = [i[:9169] for i in final_fft]
tsne_train_double_exponential_smoothing_feat = [i[:9169] for i in␣
 ↪double_exponential_smoothing_feat]
tsne_train_triple_exponential_smoothing_fet = [i[:9169] for i in␣
 ↪triple_exponential_smoothing_fet]
```

[87]:
```
# extracting the rest of the timestamp values i.e 30% of 12956 (total␣
 ↪timestamps) for our test data
tsne_test_flat_lat = [i[9169:] for i in tsne_lat]
tsne_test_flat_lon = [i[9169:] for i in tsne_lon]
tsne_test_flat_weekday = [i[9169:] for i in tsne_weekday]
tsne_test_flat_output = [i[9169:] for i in output]
tsne_test_flat_exp_avg = [i[9169:] for i in predict_list]
tsne_test_flat_fft = [i[9169:] for i in final_fft]
tsne_test_double_exponential_smoothing_feat = [i[9169:] for i in␣
 ↪double_exponential_smoothing_feat]
tsne_test_triple_exponential_smoothing_fet = [i[9169:] for i in␣
 ↪triple_exponential_smoothing_fet]
```

[88]:
```
# the above contains values in the form of list of lists (i.e. list of values␣
 ↪of each region), here we make all of them in one list
train_new_features = []
for i in range(0,30):
    train_new_features.extend(train_features[i])
test_new_features = []
for i in range(0,30):
    test_new_features.extend(test_features[i])
```

[89]:
```
# converting lists of lists into sinle list i.e flatten
# a  = [[1,2,3,4],[4,6,7,8]]
# print(sum(a,[]))
# [1, 2, 3, 4, 4, 6, 7, 8]

tsne_train_lat = sum(tsne_train_flat_lat, [])
tsne_train_lon = sum(tsne_train_flat_lon, [])
tsne_train_weekday = sum(tsne_train_flat_weekday, [])
tsne_train_output = sum(tsne_train_flat_output, [])
tsne_train_exp_avg = sum(tsne_train_flat_exp_avg,[])
tsne_train_fft = sum(tsne_train_flat_fft, [])
tsne_train_double_exponential =␣
 ↪sum(tsne_train_double_exponential_smoothing_feat, [])
tsne_train_triple_exponential =␣
 ↪sum(tsne_train_triple_exponential_smoothing_fet, [])
```

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

       tsne_test_lat = sum(tsne_test_flat_lat, [])
       tsne_test_lon = sum(tsne_test_flat_lon, [])
       tsne_test_weekday = sum(tsne_test_flat_weekday, [])
       tsne_test_output = sum(tsne_test_flat_output, [])
       tsne_test_exp_avg = sum(tsne_test_flat_exp_avg,[])
       tsne_test_fft = sum(tsne_test_flat_fft, [])
       tsne_test_double_exponential = sum(tsne_test_double_exponential_smoothing_feat,␣
        ↪[])
       tsne_test_triple_exponential = sum(tsne_test_triple_exponential_smoothing_fet,␣
        ↪[])
```

```
[91]:  # Preparing the data frame for our train data
       columns = ['ft_5','ft_4','ft_3','ft_2','ft_1']
       df_train = pd.DataFrame(data=train_new_features, columns=columns)
       df_train['lat'] = tsne_train_lat
       df_train['lon'] = tsne_train_lon
       df_train['weekday'] = tsne_train_weekday
       df_train['exp_avg'] = tsne_train_exp_avg
       df_train[['f1','f2','f3','f4','f5']] = pd.
        ↪DataFrame(tsne_train_fft,columns=['f1','f2','f3','f4','f5'])
       df_train['double_exponent'] = pd.
        ↪DataFrame(tsne_train_double_exponential,columns=['double_exponent'])
       df_train['triple_exponent'] = pd.
        ↪DataFrame(tsne_train_triple_exponential,columns=['triple_exponent'])
       print(df_train.shape)
```

```
(275070, 16)
```

```
[92]:  # Preparing the data frame for our train data
       df_test = pd.DataFrame(data=test_new_features, columns=columns)
       df_test['lat'] = tsne_test_lat
       df_test['lon'] = tsne_test_lon
       df_test['weekday'] = tsne_test_weekday
       df_test['exp_avg'] = tsne_test_exp_avg
       df_test[['f1','f2','f3','f4','f5']] = pd.
        ↪DataFrame(tsne_test_fft,columns=['f1','f2','f3','f4','f5'])
       df_test['double_exponent'] = pd.
        ↪DataFrame(tsne_test_double_exponential,columns=['double_exponent'])
       df_test['triple_exponent'] = pd.
        ↪DataFrame(tsne_test_triple_exponential,columns=['triple_exponent'])
       print(df_test.shape)
```

```
(117900, 16)
```

[93]:
```python
df_test.head()
```

[93]:
```
   ft_5  ft_4  ft_3  ft_2  ft_1        lat        lon  weekday  exp_avg  \
0   240   213   243   222   234  40.777809 -73.954054        4      231
1   213   243   222   234   291  40.777809 -73.954054        4      273
2   243   222   234   291   256  40.777809 -73.954054        4      261
3   222   234   291   256   266  40.777809 -73.954054        4      264
4   234   291   256   266   268  40.777809 -73.954054        4      266

             f1             f2             f3             f4        f5  \
0  111346.847420  111346.847420  277959.882127  277959.882127  798408.0
1  111335.281941  111335.281941  277976.590716  277976.590716  798427.0
2  111337.430741  111337.430741  277973.159507  277973.159507  798423.0
3  111338.354070  111338.354070  277971.490380  277971.490380  798421.0
4  111341.032479  111341.032479  277965.822248  277965.822248  798414.0

   double_exponent  triple_exponent
0       230.230988       240.007733
1       189.491581       212.368925
2       255.759460       242.561288
3       216.141085       221.485811
4       237.444969       233.618560
```

### 6.8.2   Using Linear Regression

[94]:
```python
# find more about LinearRegression function here http://scikit-learn.org/stable/
 ↪modules/generated/sklearn.linear_model.LinearRegression.html
# --------------------------
# default paramters
# sklearn.linear_model.LinearRegression(fit_intercept=True, normalize=False,␣
 ↪copy_X=True, n_jobs=1)

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

from sklearn.linear_model import LinearRegression
```

```
lr_reg=LinearRegression().fit(df_train, tsne_train_output)


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

### 6.8.3   Using Random Forest Regressor

```
[96]: n_estimators = [int(x) for x in np.linspace(start = 100, stop = 1200, num = 12)]
      # Number of features to consider at every split
      max_features = ['auto', 'sqrt']
      # Maximum number of levels in tree
      max_depth = [int(x) for x in np.linspace(5, 30, num = 6)]
      # max_depth.append(None)
      # Minimum number of samples required to split a node
      min_samples_split = [2, 5, 10, 15, 100]
      # Minimum number of samples required at each leaf node
      min_samples_leaf = [1, 2, 5, 10]
      # Method of selecting samples for training each tree
      # bootstrap = [True, False]

      # Create the random grid
      random_grid = {'n_estimators': n_estimators,
                     'max_features': max_features,
                     'max_depth': max_depth,
                     'min_samples_split': min_samples_split,
                     'min_samples_leaf': min_samples_leaf}

      regr1 = RandomForestRegressor()
      rf_random = RandomizedSearchCV(estimator = regr1, param_distributions =␣
       ↪random_grid, n_iter = 10, cv = 3, verbose=1, random_state=42, n_jobs = -1)
      rf_random.fit(df_train, tsne_train_output)
```

```
Fitting 3 folds for each of 10 candidates, totalling 30 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 6 concurrent workers.
[Parallel(n_jobs=-1)]: Done  30 out of  30 | elapsed: 46.2min finished
```

```
[96]: RandomizedSearchCV(cv=3, error_score='raise-deprecating',
                         estimator=RandomForestRegressor(bootstrap=True,
                                                         criterion='mse',
                                                         max_depth=None,
                                                         max_features='auto',
                                                         max_leaf_nodes=None,
                                                         min_impurity_decrease=0.0,
                                                         min_impurity_split=None,
```

```
                                    min_samples_leaf=1,
                                    min_samples_split=2,
                                    min_weight_fraction_leaf=0.0,
                                    n_estimators='warn',
                                    n_jobs=None, oob_score=False,
                                    random_sta...
                                    warm_start=False),
                    iid='warn', n_iter=10, n_jobs=-1,
                    param_distributions={'max_depth': [5, 10, 15, 20, 25, 30],
                                    'max_features': ['auto', 'sqrt'],
                                    'min_samples_leaf': [1, 2, 5, 10],
                                    'min_samples_split': [2, 5, 10, 15,
                                                        100],
                                    'n_estimators': [100, 200, 300, 400,
                                                    500, 600, 700, 800,
                                                    900, 1000, 1100,
                                                    1200]},
                    pre_dispatch='2*n_jobs', random_state=42, refit=True,
                    return_train_score=False, scoring=None, verbose=1)
```

[97]: 
```python
rf_random.best_params_
```

[97]: 
```
{'n_estimators': 1100,
 'min_samples_split': 10,
 'min_samples_leaf': 2,
 'max_features': 'sqrt',
 'max_depth': 15}
```

[98]: 
```python
regr1 = RandomForestRegressor(n_estimators= 1100, min_samples_split= 10,
  →min_samples_leaf= 2, max_features= 'sqrt', max_depth= 15)
regr1.fit(df_train, tsne_train_output)
y_pred = regr1.predict(df_test)
rndf_test_predictions = [round(value) for value in y_pred]
y_pred = regr1.predict(df_train)
rndf_train_predictions = [round(value) for value in y_pred]
```

[99]: 
```python
#feature importances based on analysis using random forest
print (df_train.columns)
print (regr1.feature_importances_)
```

```
Index(['ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1', 'lat', 'lon', 'weekday',
       'exp_avg', 'f1', 'f2', 'f3', 'f4', 'f5', 'double_exponent',
       'triple_exponent'],
      dtype='object')
[0.04003611 0.08834721 0.12377305 0.16359002 0.21713722 0.00191753
 0.00272756 0.00033876 0.27453693 0.00035212 0.00034736 0.00035473
 0.00036804 0.00037524 0.02274572 0.0630524 ]
```

### 6.8.4 Using XgBoost Regressor

```
[100]: x_cfl=xgb.XGBRegressor()
       prams={
           'learning_rate':[0.01,0.03,0.05,0.1,0.15,0.2],
            'n_estimators':[100,200,500,1000,2000],
            'max_depth':[3,5,10],
            'min_child_weight':[2,3,4,5],
           'colsample_bytree':[0.1,0.3,0.5,1],
           'subsample':[0.1,0.3,0.5,1],
            'gamma': [1,2,3],
            'eta': [0.001, 0.01, 0.02],

       }
       random_cfl=RandomizedSearchCV(x_cfl,param_distributions=prams,verbose=1,n_iter=10,␣
        ↪n_jobs=-1,)
       random_cfl.fit(df_train, tsne_train_output)
       print(random_cfl.best_params_)
```

```
Fitting 3 folds for each of 10 candidates, totalling 30 fits

[Parallel(n_jobs=-1)]: Using backend LokyBackend with 6 concurrent workers.
[Parallel(n_jobs=-1)]: Done  30 out of  30 | elapsed:  5.6min finished

[14:32:39] WARNING: src/objective/regression_obj.cu:152: reg:linear is now
deprecated in favor of reg:squarederror.
{'subsample': 1, 'n_estimators': 100, 'min_child_weight': 2, 'max_depth': 3,
'learning_rate': 0.15, 'gamma': 1, 'eta': 0.01, 'colsample_bytree': 1}
```

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

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

       # some of methods of RandomForestRegressor()
       # fit(X, y, sample_weight=None, eval_set=None, eval_metric=None,␣
        ↪early_stopping_rounds=None, verbose=True, xgb_model=None)
       # get_params([deep])          Get parameters for this estimator.
```

```python
# predict(data, output_margin=False, ntree_limit=0) : Predict with data. NOTE:␣
 ↪This function is not thread safe.
# get_score(importance_type='weight') -> get the feature importance
# ----------------------
# video link1: https://www.appliedaicourse.com/course/applied-ai-course-online/
 ↪lessons/regression-using-decision-trees-2/
# video link2: https://www.appliedaicourse.com/course/applied-ai-course-online/
 ↪lessons/what-are-ensembles/
# ----------------------

x_model = xgb.XGBRegressor(subsample= 1, n_estimators= 100, min_child_weight=␣
 ↪2, max_depth= 3, learning_rate= 0.15, gamma= 1, eta= 0.01, colsample_bytree=␣
 ↪1)
x_model.fit(df_train, tsne_train_output)
```

```
[14:43:38] WARNING: src/objective/regression_obj.cu:152: reg:linear is now
deprecated in favor of reg:squarederror.
```

[101]: XGBRegressor(base_score=0.5, booster='gbtree', colsample_bylevel=1,
                colsample_bynode=1, colsample_bytree=1, eta=0.01, gamma=1,
                importance_type='gain', learning_rate=0.15, max_delta_step=0,
                max_depth=3, min_child_weight=2, missing=None, n_estimators=100,
                n_jobs=1, nthread=None, objective='reg:linear', random_state=0,
                reg_alpha=0, reg_lambda=1, scale_pos_weight=1, seed=None,
                silent=None, subsample=1, verbosity=1)

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

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

```python
[103]: #feature importances
x_model.get_booster().get_score(importance_type='weight')
```

[103]: {'exp_avg': 136,
  'ft_1': 125,
  'ft_3': 41,
  'ft_2': 79,
  'lat': 21,
  'ft_4': 28,
  'triple_exponent': 35,
  'ft_5': 77,
  'lon': 35,
  'double_exponent': 62,

```
 'f5': 18,
 'weekday': 6,
 'f1': 19,
 'f3': 17}
```

### 6.8.5  Calculating the error metric values for various models

```
[104]: train_mape=[]
       test_mape=[]

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

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

### 6.8.6  Error Metric Matrix

```
[105]: print ("Error Metric Matrix (Tree Based Regression Methods) -  MAPE")
       print␣
        ↪("-------------------------------------------------------------------------------------
       print ("Baseline Model -                       Train: ",train_mape[0]," ␣
        ↪   Test: ",test_mape[0])
       print ("Exponential Averages Forecasting -     Train: ",train_mape[1]," ␣
        ↪   Test: ",test_mape[1])
       print ("Linear Regression -                    Train: ",train_mape[4]," ␣
        ↪   Test: ",test_mape[4])
```

```
print ("Random Forest Regression -                    Train: ",train_mape[2]," ⊔
 ↪  Test: ",test_mape[2])
print ("XgBoost Regression -                    Train: ",train_mape[3]," ⊔
 ↪   Test: ",test_mape[3])
print⊔
 ↪("---------------------------------------------------------------------------
```

```
Error Metric Matrix (Tree Based Regression Methods) -  MAPE
-------------------------------------------------------------------------------
-----------------------
Baseline Model -                           Train:  0.12477882091940766
Test:  0.12137217161272074
Exponential Averages Forecasting -         Train:  0.11976904266333344
Test:  0.11613179453264473
Linear Regression -                        Train:  0.11920005631976266
Test:  0.2783937780258345
Random Forest Regression -                 Train:  0.10128958808870862
Test:  0.11722052532789988
XgBoost Regression -                       Train:  0.1180291770034226
Test:  0.12265070396690994
-------------------------------------------------------------------------------
-----------------------
```

# 7   Assignments

```
Random Forest Regression model is the best model with test MAPE as 11.72%
```