

Import required libraries:

```
In [194]: 1 import warnings
          2 warnings.filterwarnings('ignore')
```

```
In [195]: 1 import pandas as pd
          2 import tsfresh as ts
          3 import numpy as np
          4 import pickle
          5 import seaborn as sns
          6 import matplotlib.pyplot as plt
          7 import math
          8 from sklearn.linear_model import SGDRegressor
          9 from sklearn.model_selection import GridSearchCV, RandomizedSearchCV, TimeSeriesSplit
         10 from sklearn.metrics import make_scorer, mean_absolute_error
         11 from sklearn.ensemble import RandomForestRegressor
         12 from sklearn.preprocessing import StandardScaler
         13 from prettytable import PrettyTable
         14 import xgboost as xgb
         15 from xgboost.sklearn import XGBRegressor
         16 from sklearn.linear_model import LinearRegression
         17 import pickle
         18 from collections import OrderedDict
         19 %matplotlib inline
```

Load and Store current state of object:

```
In [196]: 1 #Functions to save objects for later use and retrieve it
          2 def savetofile(obj,filename):
          3     pickle.dump(obj,open(filename+".pkl","wb"))
          4 def openfromfile(filename):
          5     temp = pickle.load(open(filename+".pkl","rb"))
          6     return temp
          7
```

Load Preprocessed Data:

```
In [197]: 1 jan_2015_smooth=openfromfile('jan_2015_smooth')
          2 jan_2016_smooth=openfromfile('jan_2016_smooth')
          3 kmeans=openfromfile('kmeans')
          4 regions_cum=openfromfile('regions_cum')
```

```
In [198]: 1 #Preparing the Dataframe only with x(i) values as jan-2015 data and y(i) values as jan-2016
          2 ratios_jan = pd.DataFrame()
          3 ratios_jan['Given']=jan_2015_smooth
          4 ratios_jan['Prediction']=jan_2016_smooth
          5 ratios_jan['Ratios']=ratios_jan['Prediction']*1.0/ratios_jan['Given']*1.0
```

```
In [199]: 1 ratios_jan.shape
```

```
Out[199]: (133920, 3)
```

```
In [200]: 1 ratios_jan.head(5)
```

```
Out[200]:
```

	Given	Prediction	Ratios
0	191	207	1.083770
1	381	315	0.826772
2	403	383	0.950372
3	374	364	0.973262
4	400	362	0.905000

Models:

1.Baseline Models

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

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

Function for BaseLine models:

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

```

In [201]: 1 def MA_R_Predictions(ratios,month):
2     predicted_ratio=(ratios['Ratios'].values)[0]
3     error=[]
4     predicted_values=[]
5     window_size=3
6     predicted_ratio_values=[]
7     for i in range(0,4464*30):
8         if i%4464==0:
9             # initially / start we dont have past data so ratio=0
10            predicted_ratio_values.append(0)
11            predicted_values.append(0)
12            error.append(0)
13            continue
14            #when we past data is available predicted_ratio is avg of all past ration in a recent window
15            predicted_ratio_values.append(predicted_ratio)# when i=1 we append ratio at i=0, after i=0 append avg ratio
16            #p(t)_2016 = p(t)_2015 * R(t)_2016
17            #where R(t)_2016 is calculated using past r(t-1),R(t-2),-----so on
18            predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
19            error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(ratios['Prediction'].values)[i])))
20            if i+1>=window_size:
21                predicted_ratio=sum((ratios['Ratios'].values)[(i+1)-window_size:(i+1)])/window_size
22            else:
23                # moving avg of past ratio, here in slicing i+1 is used bcz slice start from begin to end-1,
24                # divide by (i+1) bcz we started i from 0, and avg count from 1.
25                predicted_ratio=sum((ratios['Ratios'].values)[0:(i+1)])/(i+1)
26
27
28            ratios['MA_R_Predicted'] = predicted_values
29            ratios['MA_R_Error'] = error
30            mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
31            mse_err = sum([e**2 for e in error])/len(error)
32            return ratios,mape_err,mse_err

```

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

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

In [202]:

```

1  def MA_P_Predictions(ratios,month):
2      predicted_value=(ratios['Prediction'].values)[0]
3      error=[]
4      predicted_values=[]
5      window_size=1
6      predicted_ratio_values=[]
7      for i in range(0,4464*30):
8          predicted_values.append(predicted_value)
9          error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
10         if i+1>=window_size:
11             predicted_value=int(sum((ratios['Prediction'].values)[(i+1)-window_size:(i+1)])/window_size)
12         else:
13             predicted_value=int(sum((ratios['Prediction'].values)[0:(i+1)])/(i+1))
14
15     ratios['MA_P_Predicted'] = predicted_values
16     ratios['MA_P_Error'] = error
17     mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
18     mse_err = sum([e**2 for e in error])/len(error)
19     return ratios,mape_err,mse_err

```

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

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

In [203]:

```

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

```

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

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

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

In [204]:

```

1  def WA_P_Predictions(ratios,month):
2      predicted_value=(ratios['Prediction'].values)[0]
3      error=[]
4      predicted_values=[]
5      window_size=2
6      for i in range(0,4464*30):
7          predicted_values.append(predicted_value)
8          error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
9          if i+1>=window_size:
10             sum_values=0
11             sum_of_coeff=0
12             for j in range(window_size,0,-1):
13                 sum_values += j*(ratios['Prediction'].values)[i-window_size+j]
14                 sum_of_coeff+=j
15             predicted_value=int(sum_values/sum_of_coeff)
16
17         else:
18             sum_values=0
19             sum_of_coeff=0
20             for j in range(i+1,0,-1):
21                 sum_values += j*(ratios['Prediction'].values)[j-1]
22                 sum_of_coeff+=j
23             predicted_value=int(sum_values/sum_of_coeff)
24
25     ratios['WA_P_Predicted'] = predicted_values
26     ratios['WA_P_Error'] = error
27     mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
28     mse_err = sum([e**2 for e in error])/len(error)
29     return ratios,mape_err,mse_err

```

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

[1.3]Exponential Weighted Moving Averages

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

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

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

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

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

In [205]:

```

1 def EA_R1_Predictions(ratios,month):
2     predicted_ratio=(ratios['Ratios'].values)[0]
3     alpha=0.6
4     error=[]
5     predicted_values=[]
6     predicted_ratio_values=[]
7     for i in range(0,4464*30):
8         if i%4464==0:
9             predicted_ratio_values.append(0)
10            predicted_values.append(0)
11            error.append(0)
12            continue
13            predicted_ratio_values.append(predicted_ratio)
14            predicted_values.append(int(((ratios['Given'].values)[i])*predicted_ratio))
15            error.append(abs((math.pow(int(((ratios['Given'].values)[i])*predicted_ratio)-(ratios['Prediction'].values)[i]))
16            predicted_ratio = (alpha*predicted_ratio) + (1-alpha)*((ratios['Ratios'].values)[i])
17
18            ratios['EA_R1_Predicted'] = predicted_values
19            ratios['EA_R1_Error'] = error
20            mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
21            mse_err = sum([e**2 for e in error])/len(error)
22            return ratios,mape_err,mse_err

```

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

In [206]:

```

1 def EA_P1_Predictions(ratios,month):
2     predicted_value= (ratios['Prediction'].values)[0]
3     alpha=0.3
4     error=[]
5     predicted_values=[]
6     for i in range(0,4464*30):
7         if i%4464==0:
8             predicted_values.append(0)
9             error.append(0)
10            continue
11            predicted_values.append(predicted_value)
12            error.append(abs((math.pow(predicted_value-(ratios['Prediction'].values)[i],1))))
13            predicted_value =int((alpha*predicted_value) + (1-alpha)*((ratios['Prediction'].values)[i]))
14
15            ratios['EA_P1_Predicted'] = predicted_values
16            ratios['EA_P1_Error'] = error
17            mape_err = (sum(error)/len(error))/(sum(ratios['Prediction'].values)/len(ratios['Prediction'].values))
18            mse_err = sum([e**2 for e in error])/len(error)
19            return ratios,mape_err,mse_err

```

Applying Baseline Models:

In [207]:

```

1 ratios_jan, mape_sma_r, mse_sma_r = MA_R_Predictions(ratios_jan,'jan')
2 ratios_jan, mape_sma_p, mse_sma_p = MA_P_Predictions(ratios_jan,'jan')
3 ratios_jan, mape_wma_r, mse_wma_r = WA_R_Predictions(ratios_jan,'jan')
4 ratios_jan, mape_wma_p, mse_wma_p = WA_P_Predictions(ratios_jan,'jan')
5 ratios_jan, mape_ewma_r, mse_ewma_r = EA_R1_Predictions(ratios_jan,'jan')
6 ratios_jan, mape_ewma_p, mse_ewma_p = EA_P1_Predictions(ratios_jan,'jan')

```

Summary of baseline models:

We have chosen our error metric for comparison between models as **MAPE (Mean Absolute Percentage Error)** so that we can know that on an average how good is our model with predictions and **MSE (Mean Squared Error)** is also used so that we have a clearer understanding as to how well our forecasting model performs with outliers so that we make sure that there is not much of a error margin between our prediction and the actual value

In [208]:

```

1 print ("Error Metric Matrix (Forecasting Methods) - MAPE & MSE")
2 print ("-----")
3 #print(model_summary)
4 print ("Moving Averages (Ratios) - MAPE: %.3f"%(100*mape_sma_r), " MSE: %.4f"%mse_sma_r)
5 print ("Moving Averages (2016 Values) - MAPE: %.3f"%(100*mape_sma_p), " MSE: %.4f"%mse_sma_p)
6 print ("-----")
7 print ("Weighted Moving Averages (Ratios) - MAPE: %.3f"%(100*mape_wma_r), " MSE: %.4f"%mse_wma_r)
8 print ("Weighted Moving Averages (2016 Values) - MAPE: %.3f"%(100*mape_wma_p), " MSE: %.4f"%mse_wma_p)
9 print ("-----")
10 print ("Exponential Moving Averages (Ratios) - MAPE: %.3f"%(100*mape_ewma_r), " MSE: %.4f"%mse_ewma_r)
11 print ("Exponential Moving Averages (2016 Values) - MAPE: %.3f"%(100*mape_ewma_p), " MSE: %.4f"%mse_ewma_p)

```

Error Metric Matrix (Forecasting Methods) - MAPE & MSE

Moving Averages (Ratios) -	MAPE: 16.287	MSE: 496.6735
Moving Averages (2016 Values) -	MAPE: 12.761	MSE: 236.2616

Weighted Moving Averages (Ratios) -	MAPE: 16.005	MSE: 486.8578
Weighted Moving Averages (2016 Values) -	MAPE: 12.195	MSE: 222.2438

Exponential Moving Averages (Ratios) -	MAPE: 15.975	MSE: 490.0153
Exponential Moving Averages (2016 Values) -	MAPE: 12.176	MSE: 221.1123

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

From the above matrix it is inferred that the best forecasting model for our prediction would be:- $P'_t = \alpha * P_{t-1} + (1 - \alpha) * P'_{t-1}$ i.e Exponential Moving Averages using 2016 Values

2.Regression Models

[2.1]Preparing data for regression model:

1. cluster center latitude
2. cluster center longitude
3. day of the week
4. ft_1: number of pickups that are happened previous t-1th 10min intravel

5. ft_2: number of pickups that are happened previous t-2th 10min intravel
6. ft_3: number of pickups that are happened previous t-3th 10min intravel
7. ft_4: number of pickups that are happened previous t-4th 10min intravel
8. ft_5: number of pickups that are happened previous t-5th 10min intravel

In [209]:

```
1  # Preparing data to be split into train and test, The below prepares data in cumulative form which will be later spl
2  # number of 10min indices for jan 2015= 24*31*60/10 = 4464
3  # number of 10min indices for jan 2016 = 24*31*60/10 = 4464
4  # number of 10min indices for feb 2016 = 24*29*60/10 = 4176
5  # number of 10min indices for march 2016 = 24*31*60/10 = 4464
6  # regions_cum: it will contain 40 lists, each list will contain 4464+4176+4464 values which represents the number of
7  # that are happened for three months in 2016 data
8
9  # print(len(regions_cum))
10 # 40
11 # print(len(regions_cum[0]))
12 # 12960
13
14 # we take number of pickups that are happened in last 5 10min intravels
15 number_of_time_stamps = 5
16
17 # output variable
18 # it is list of lists
19 # it will contain number of pickups 4459 for each cluster
20 output = []
21
22
23 # tsne_lat will contain 13104-5=13099 times Latitude of cluster center for every cluster
24 # Ex: [[cent_lat 13099times],[cent_lat 13099times], [cent_lat 13099times].... 40 lists]
25 # it is list of lists
26 tsne_lat = []
27
28
29 # tsne_lon will contain 13104-5=13099 times Logitude of cluster center for every cluster
30 # Ex: [[cent_long 13099times],[cent_long 13099times], [cent_long 13099times].... 40 lists]
31 # it is list of lists
32 tsne_lon = []
33
34 # we will code each day
35 # sunday = 0, monday=1, tue = 2, wed=3, thur=4, fri=5,sat=6
36 # for every cluster we will be adding 13099 values, each value represent to which day of the week that pickup bin be
37 # it is list of lists
38 tsne_weekday = []
39
40 # its an numpy array, of shape (523960, 5)
41 # each row corresponds to an entry in out data
```

```

42 # for the first row we will have [f0,f1,f2,f3,f4] fi=number of pickups happened in i+1th 10min intravel(bin)
43 # the second row will have [f1,f2,f3,f4,f5]
44 # the third row will have [f2,f3,f4,f5,f6]
45 # and so on...
46 features = []
47 #tsne_feature=features
48
49 features = [0]*number_of_time_stamps
50 for i in range(0,30):
51     tsne_lat.append([kmeans.cluster_centers_[i][0]]*4459)# bcz first 5 are used for given
52     tsne_lon.append([kmeans.cluster_centers_[i][1]]*4459)
53     # jan 1st 2016 is friday, so we start our day from 4: "(int(k/144))%7+5"
54     # our prediction start from 5th 10min intravel since we need to have number of pickups that are happened in last
55     tsne_weekday.append([int(((int(k/144))%7+5)%7) for k in range(5,4464)])# 0,1,2,3,4 used as
56     # regions_cum is a list of lists [[x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3..x13104], [x1,x2,x3..x13104]
57     features = np.vstack((features, [regions_cum[i][r:r+number_of_time_stamps] for r in range(0,len(regions_cum[i))-
58     output.append(regions_cum[i][5:])
59 features = features[1:]

```

```

In [210]: 1 len(tsne_lat[0])*len(tsne_lat) == features.shape[0] == len(tsne_weekday)*len(tsne_weekday[0]) == 30*4459 == len(outp

```

Out[210]: True

[2.2] Feature Extraction:

[2.2.1]Add exponential moving average as new feature to data for regression model:

In [211]:

```

1  # Getting the predictions of exponential moving averages to be used as a feature in cumulative form
2
3  # upto now we computed 8 features for every data point that starts from 50th min of the day
4  # 1. cluster center Latitude
5  # 2. cluster center Longitude
6  # 3. day of the week
7  # 4. f_t_1: number of pickups that are happened previous t-1th 10min intravel
8  # 5. f_t_2: number of pickups that are happened previous t-2th 10min intravel
9  # 6. f_t_3: number of pickups that are happened previous t-3th 10min intravel
10 # 7. f_t_4: number of pickups that are happened previous t-4th 10min intravel
11 # 8. f_t_5: number of pickups that are happened previous t-5th 10min intravel
12
13 # from the baseline models we said the exponential weighted moving average gives us the best error
14 # we will try to add the same exponential weighted moving average at t as a feature to our data
15 # exponential weighted moving average =>  $p'(t) = \alpha * p'(t-1) + (1-\alpha) * P(t-1)$ 
16 alpha=0.3
17
18 # it is a temporary array that store exponential weighted moving average for each 10min intravel,
19 # for each cluster it will get reset
20 # for every cluster it contains 13104 values
21 predicted_values=[]
22
23 # it is similar like tsne_lat
24 # it is list of lists
25 # predict_list is a list of lists [[x5,x6,x7..x13104], [x5,x6,x7..x13104], [x5,x6,x7..x13104], [x5,x6,x7..x13104], [
26 predict_list = []
27 tsne_flat_exp_avg = []
28 for r in range(0,30):
29     for i in range(0,4464):
30         if i==0:
31             predicted_value= regions_cum[r][0]
32             predicted_values.append(0)
33             continue
34             predicted_values.append(predicted_value)
35             #EWMA
36             predicted_value =int((alpha*predicted_value) + (1-alpha)*(regions_cum[r][i]))
37             predict_list.append(predicted_values[5:])#discarding 1st 5 values bcz 1st 5 values we used as given values to pr
38             predicted_values=[]

```

[2.2.2]Time series :

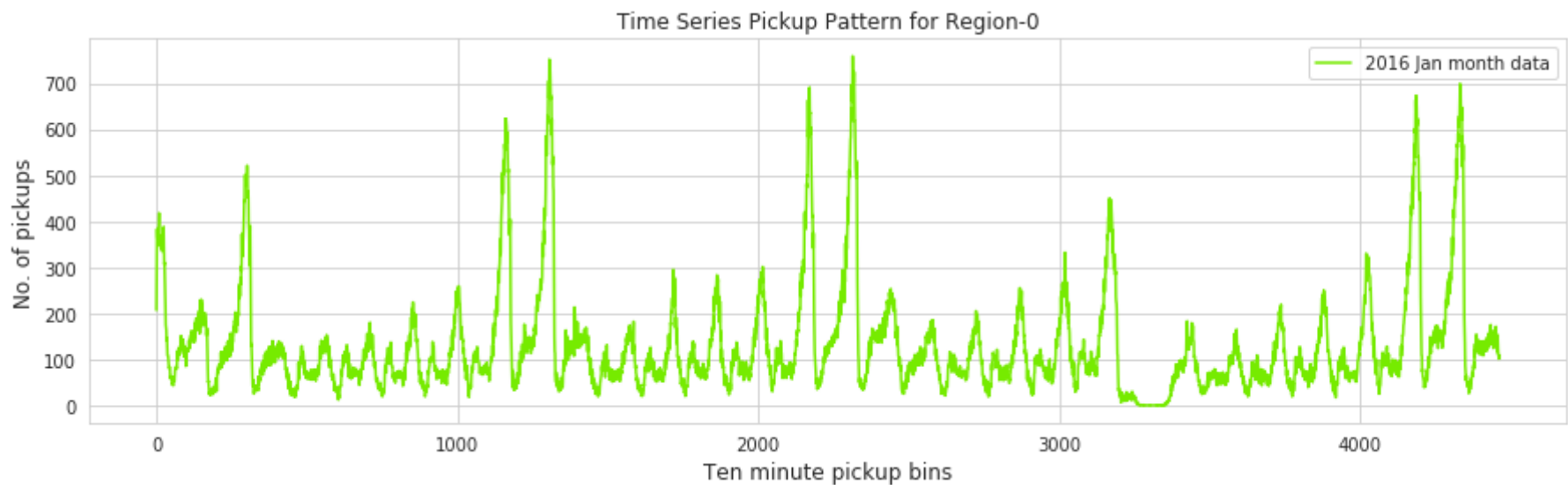
We plot our time series data to check whether it is having repetative pattern , If our time series have repetative pattern then we can able to use "Fourier Transform" to create new features .

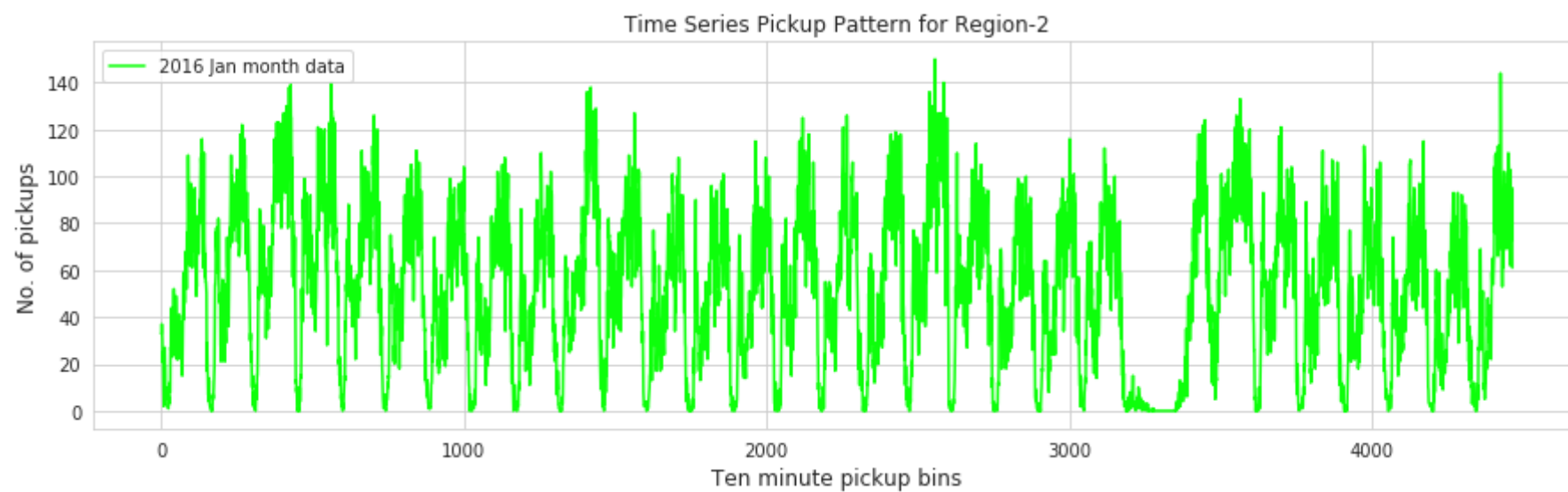
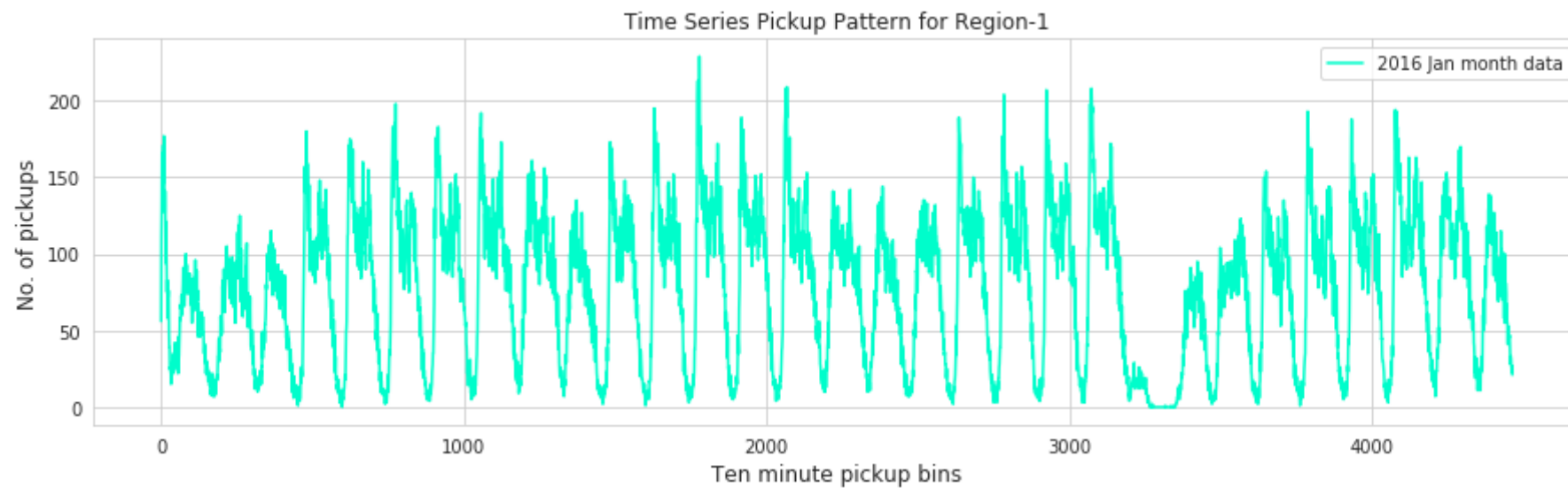
In [212]:

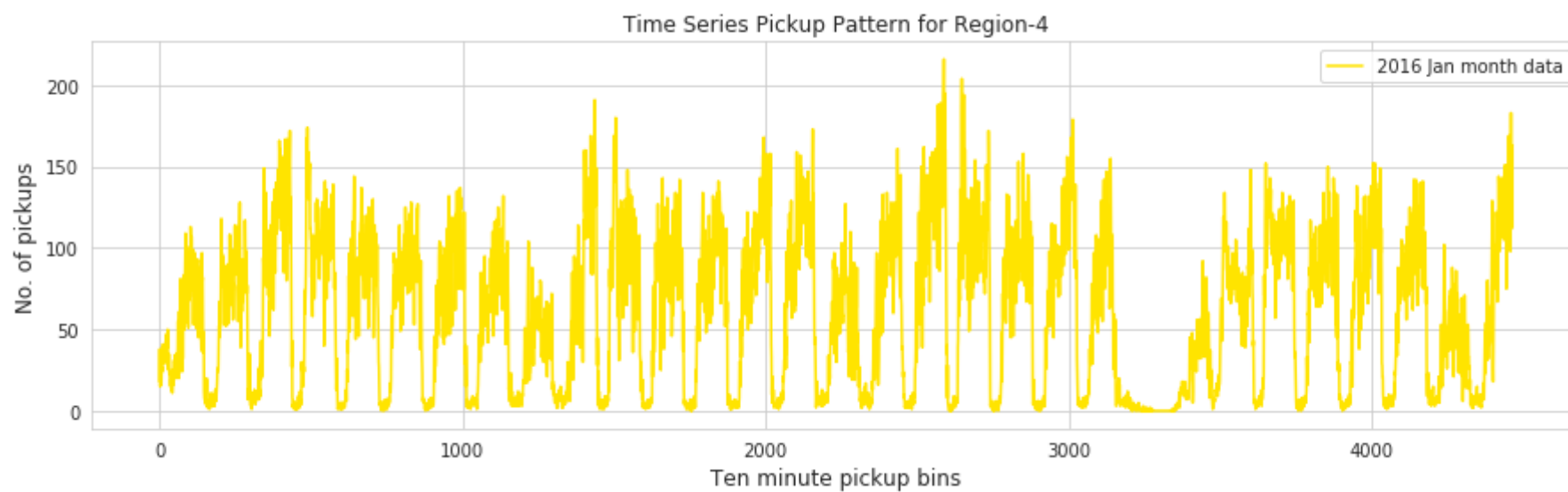
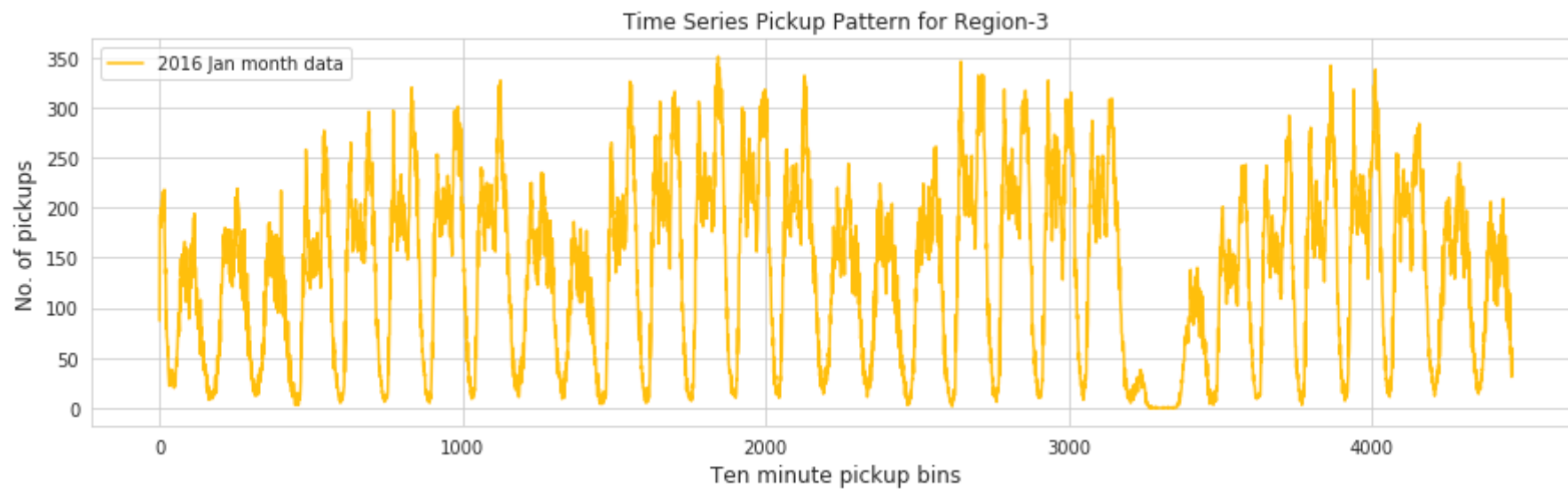
```

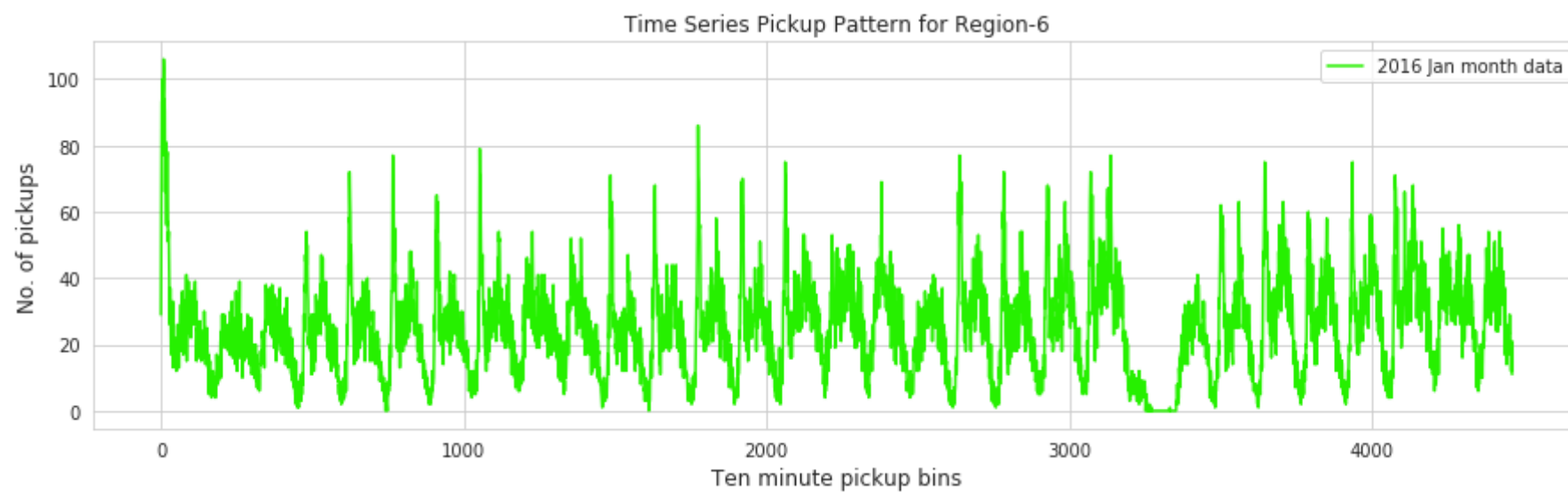
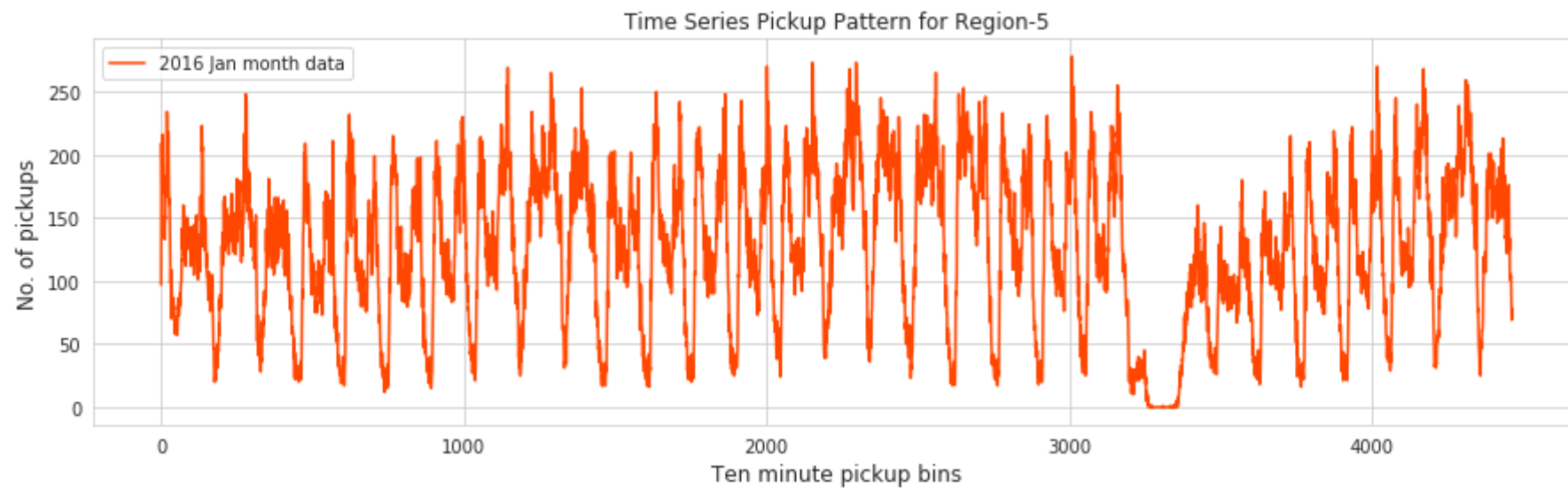
1 def uniqueish_color():
2     """There're better ways to generate unique colors, but this isn't awful."""
3     return plt.cm.gist_ncar(np.random.random())
4 first_x = list(range(0,4464))
5 #second_x = list(range(4464,8640))
6 #third_x = list(range(8640,13104))
7 for i in range(30):
8     plt.figure(figsize=(15,4))
9     sns.set_style('whitegrid')
10    plt.title('Time Series Pickup Pattern for Region-%d%i, size=12)
11    plt.xlabel('Ten minute pickup bins', fontsize=12)
12    plt.ylabel('No. of pickups', fontsize=12)
13    plt.plot(first_x,regions_cum[i][:4464], color=uniqueish_color(), label='2016 Jan month data')
14    #plt.plot(second_x,regions_cum[i][4464:8640], color=uniqueish_color(), label='2016 feb month data')
15    #plt.plot(third_x,regions_cum[i][8640:], color=uniqueish_color(), label='2016 march month data')
16    plt.legend()
17    plt.show()

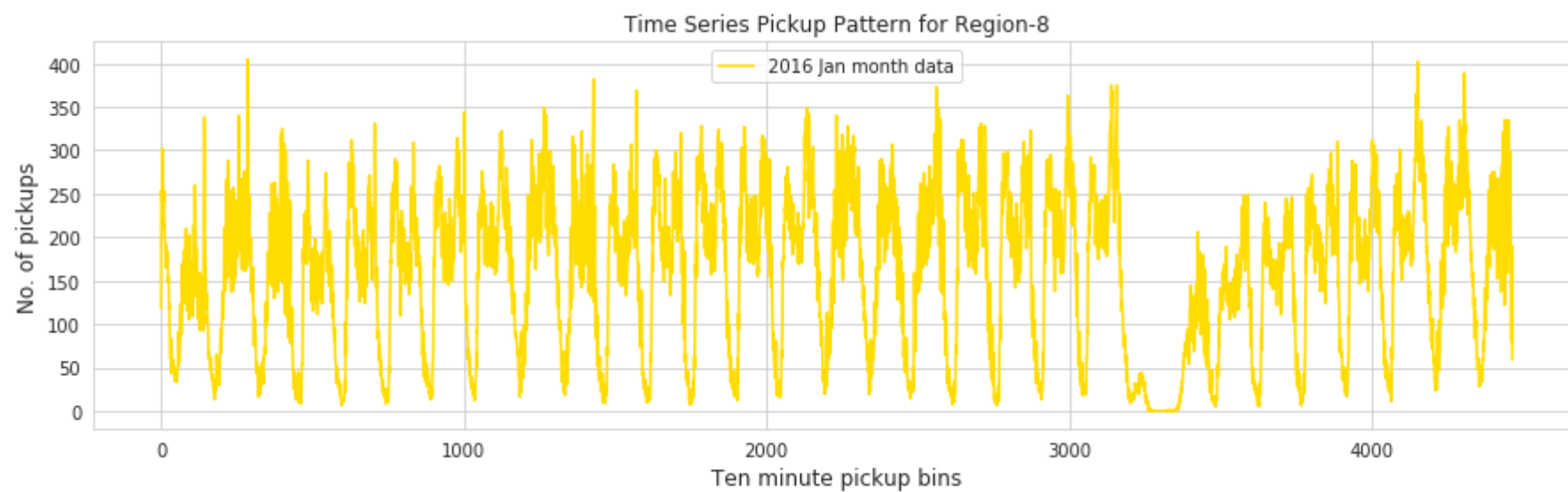
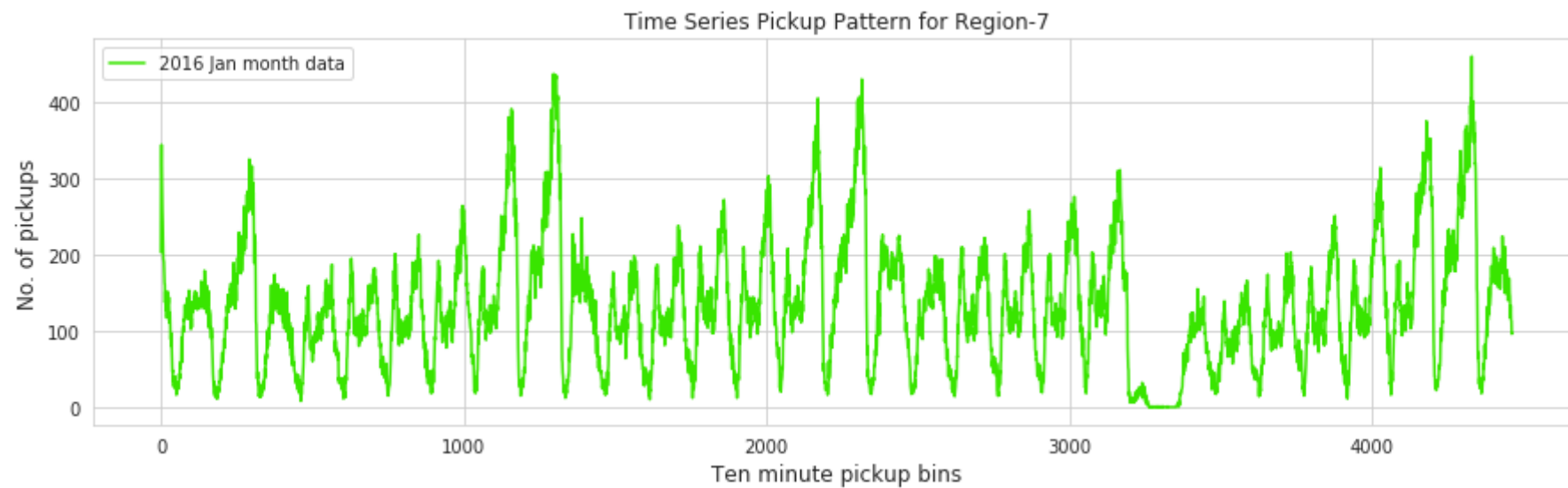
```

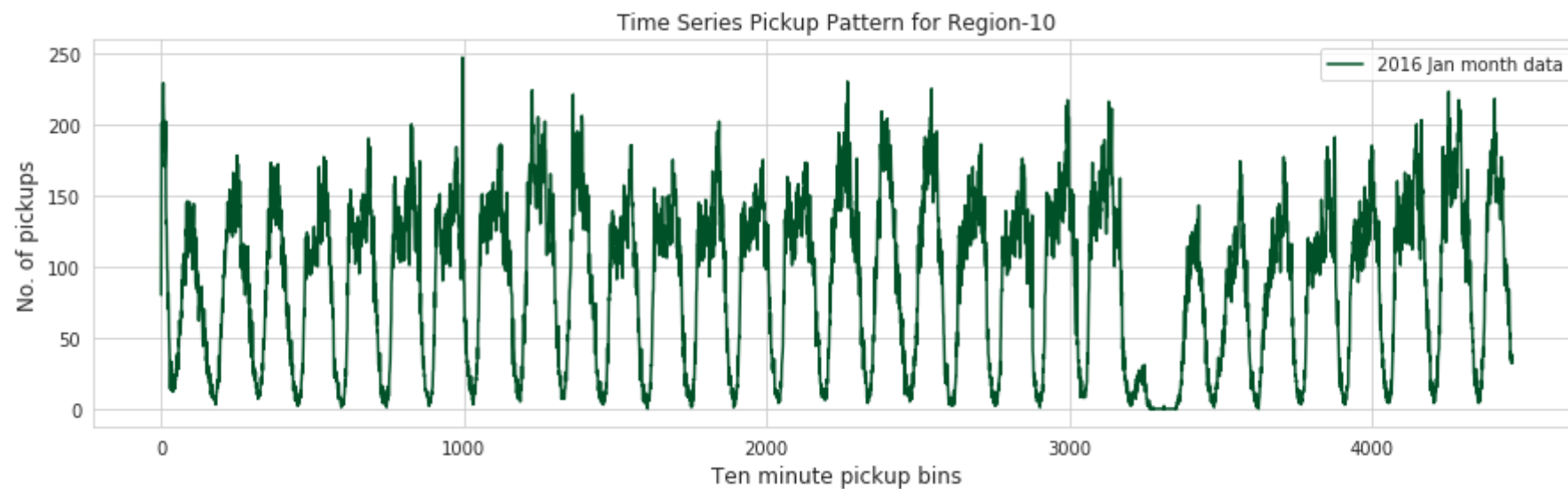
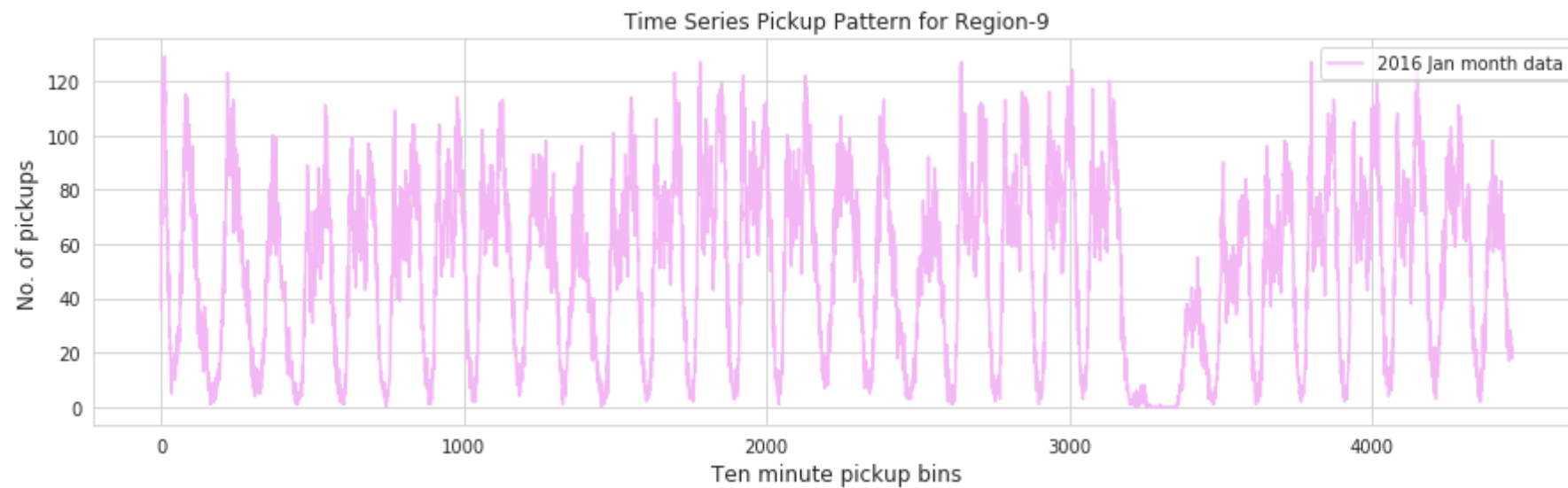


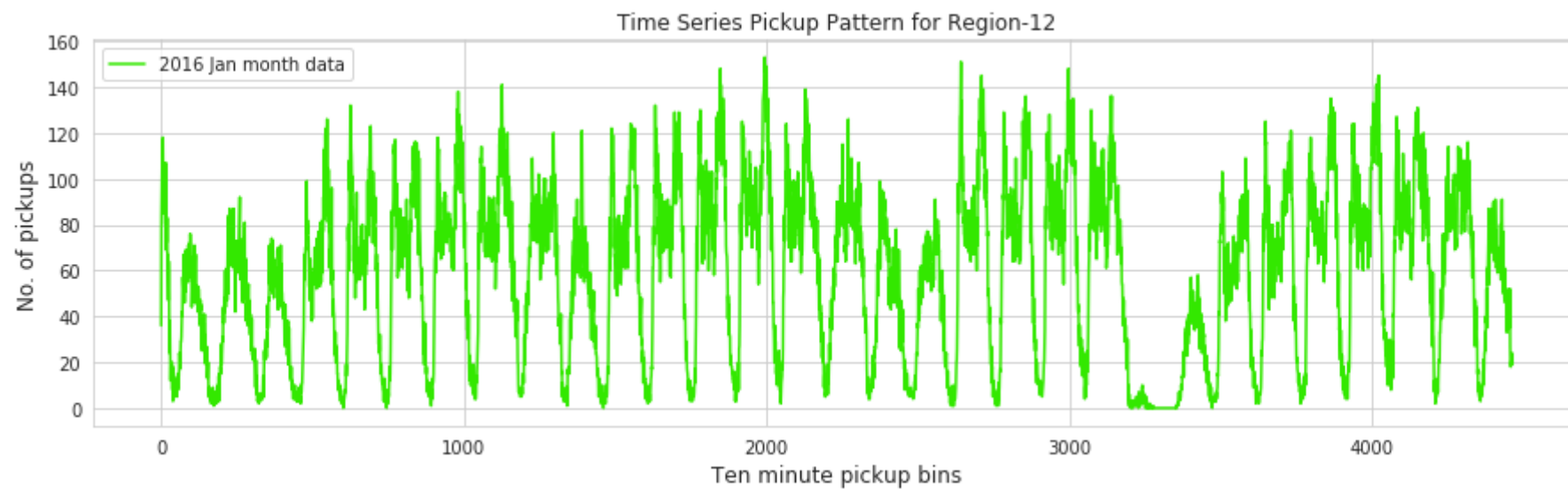
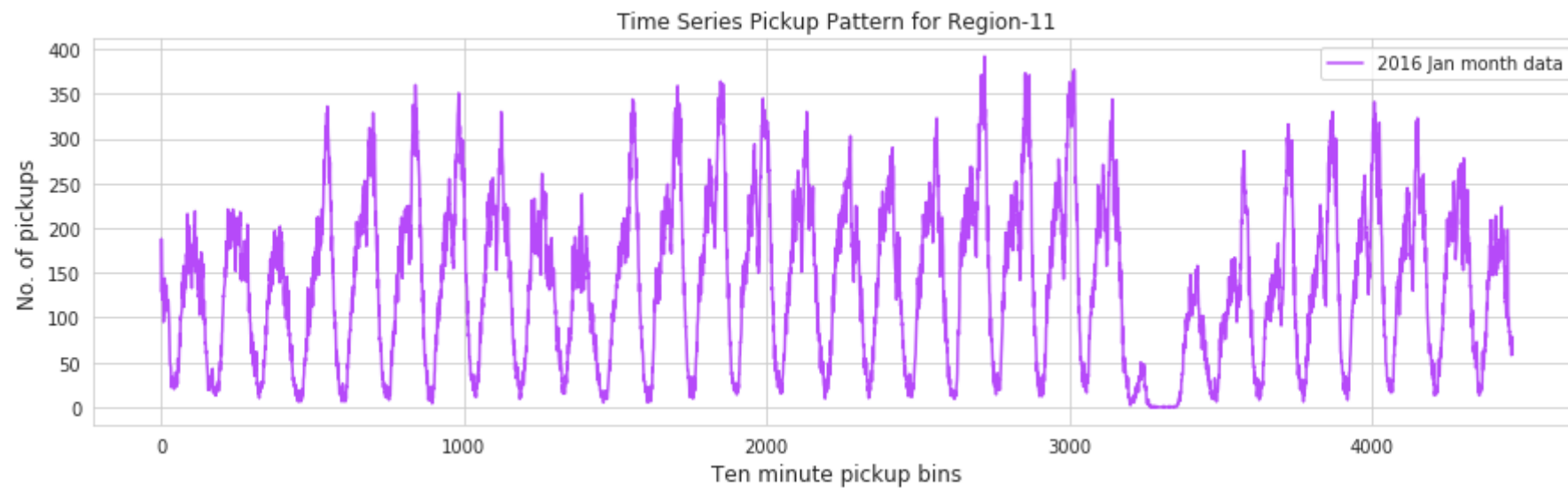


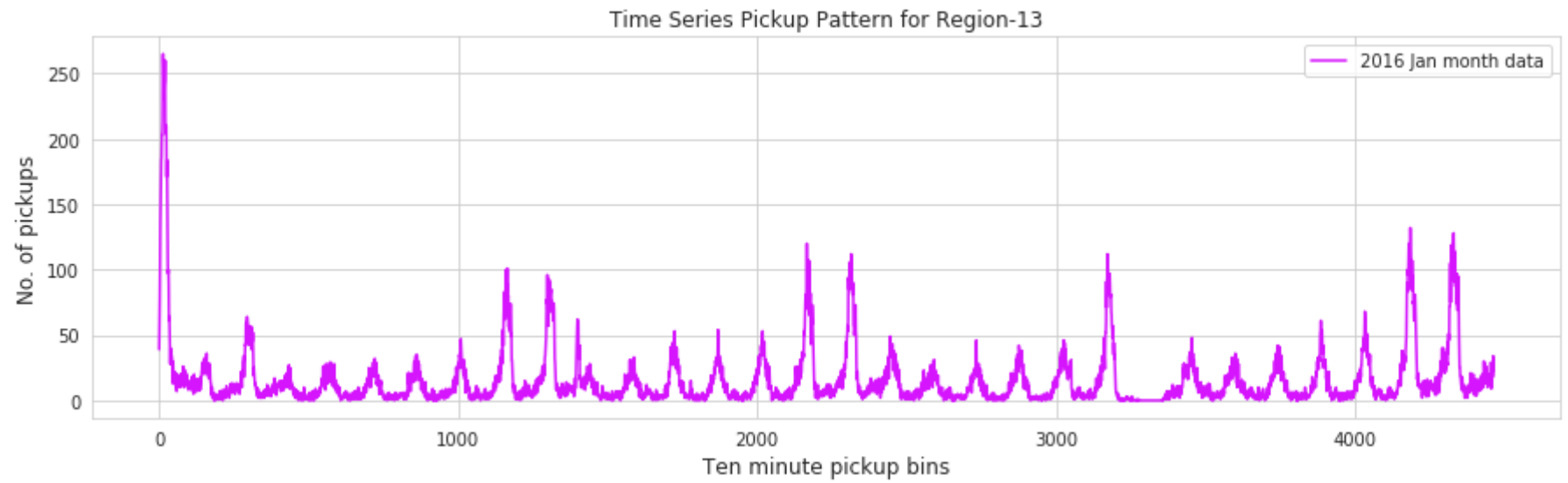


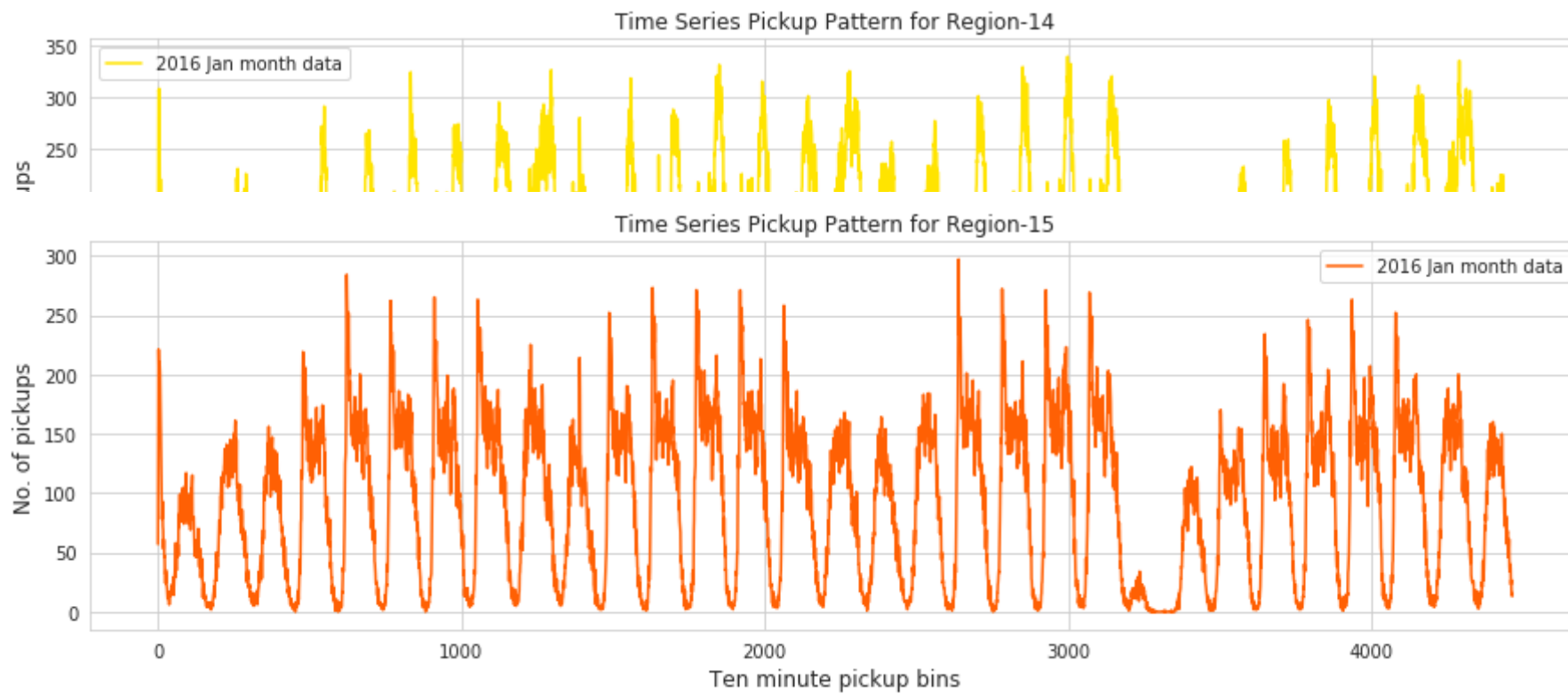


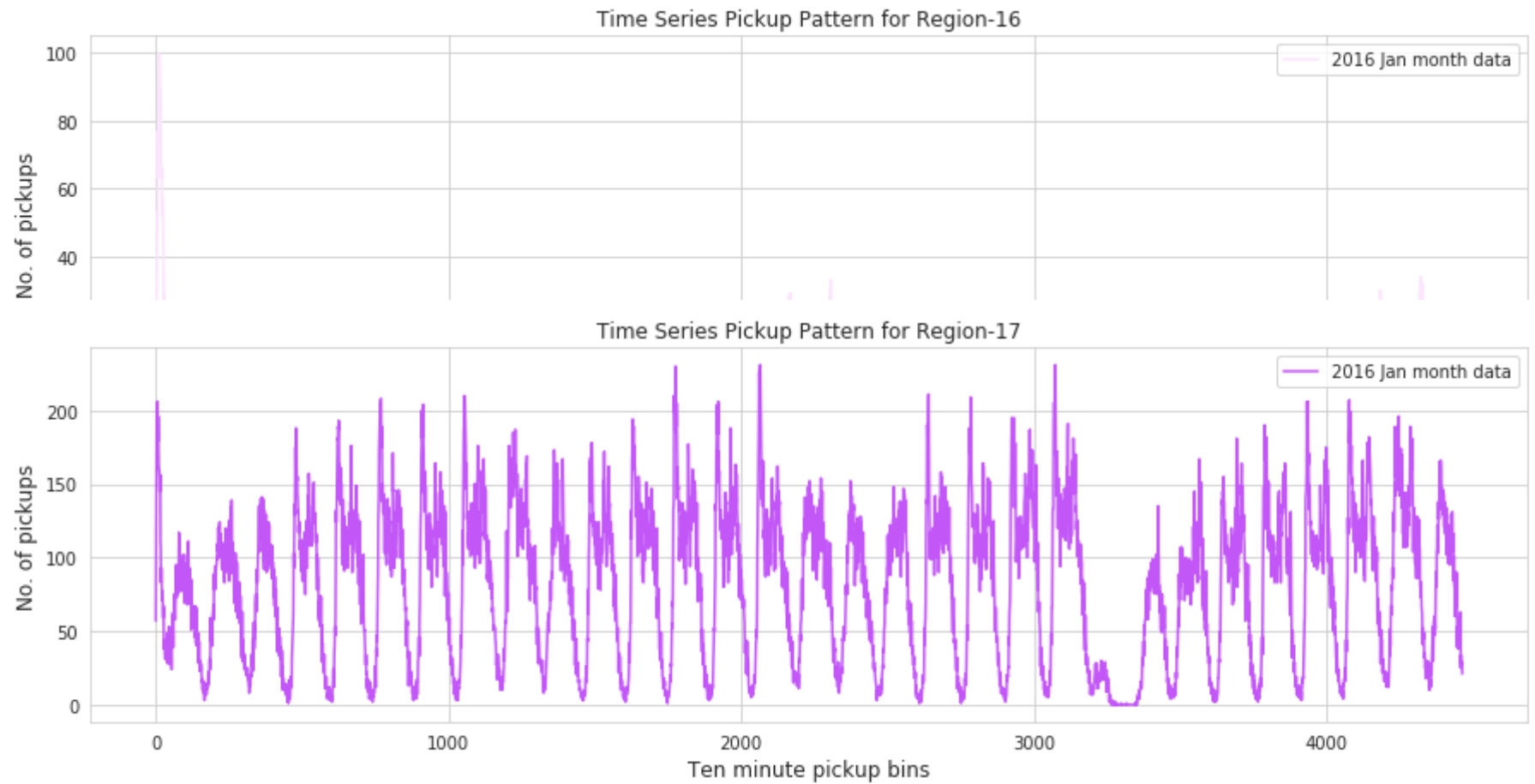


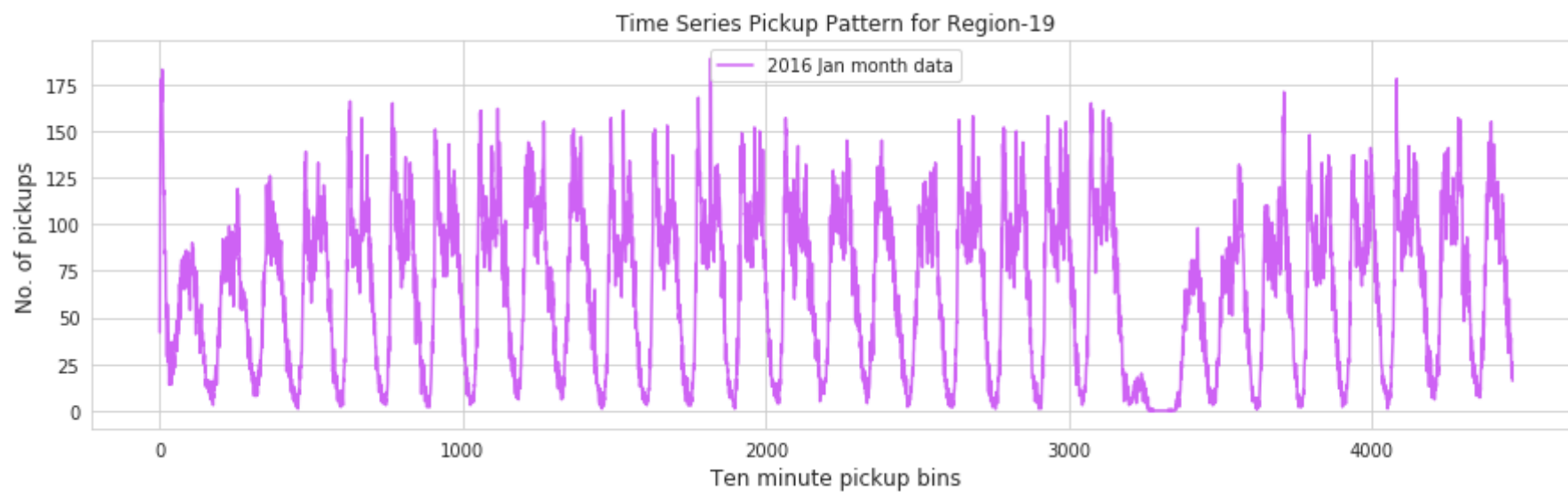
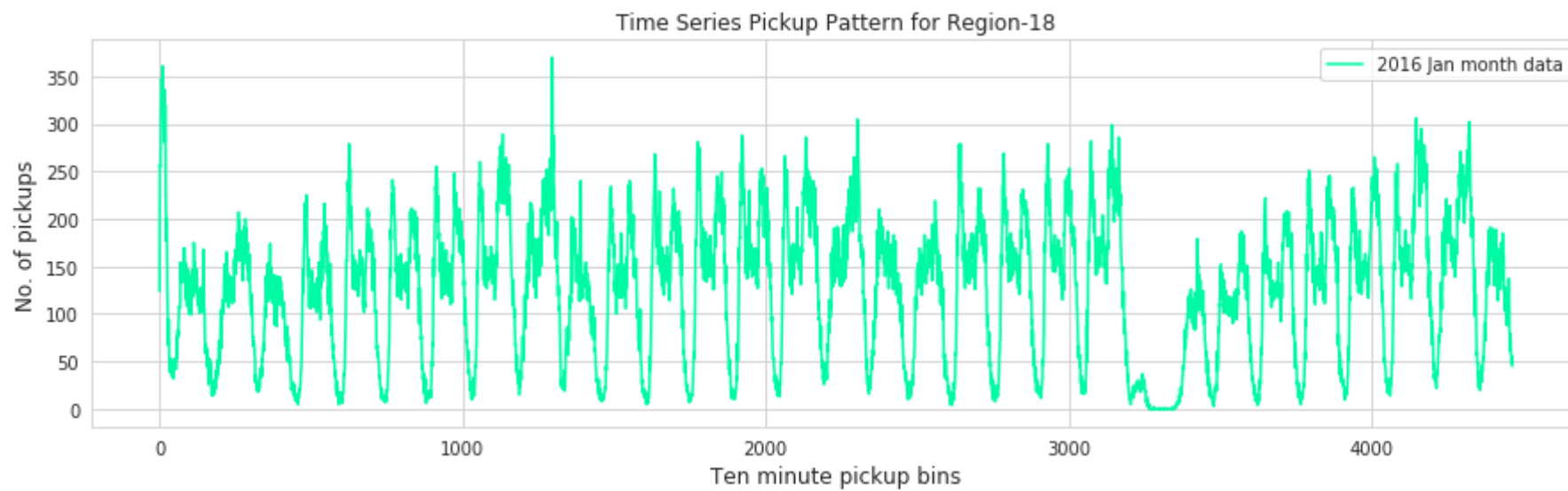


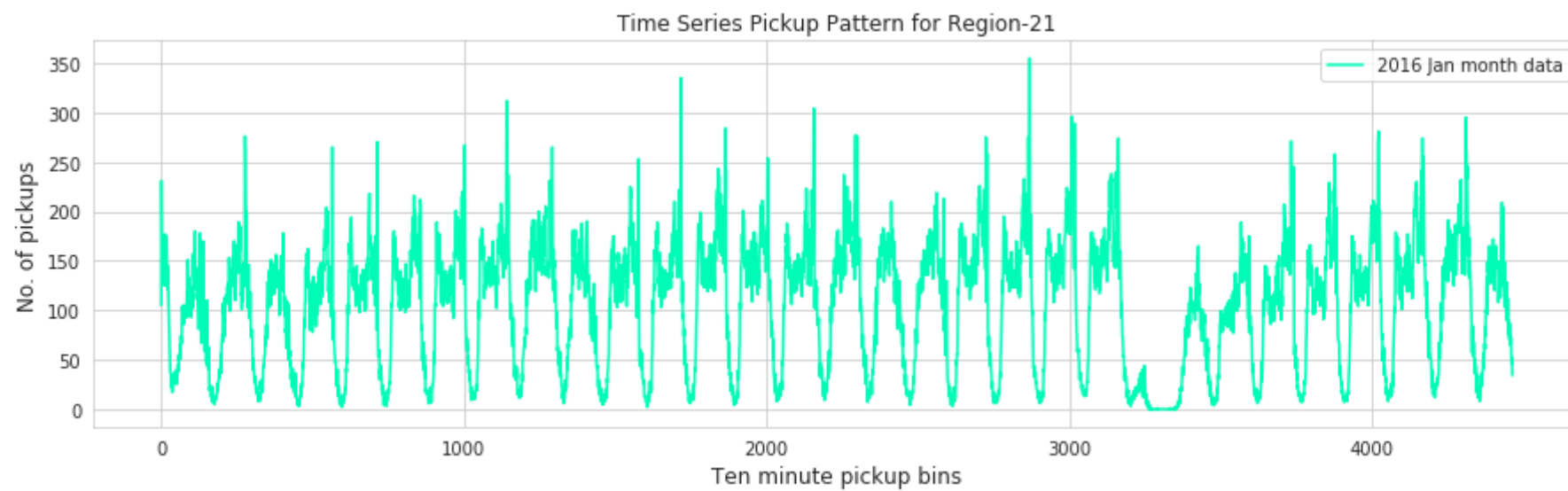
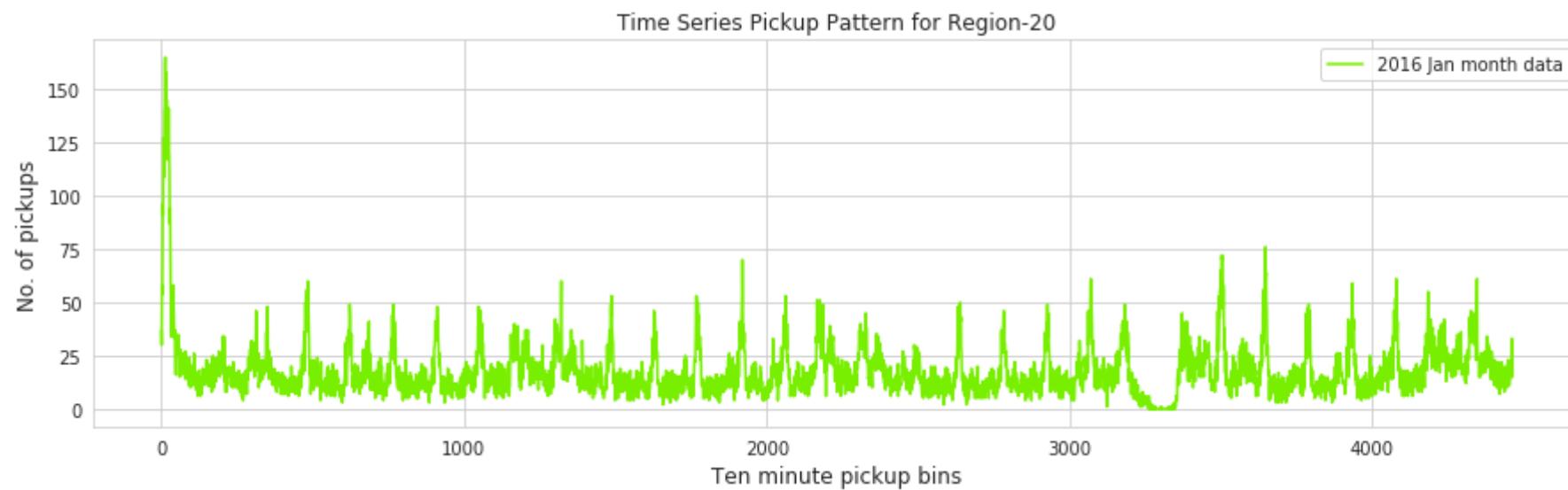


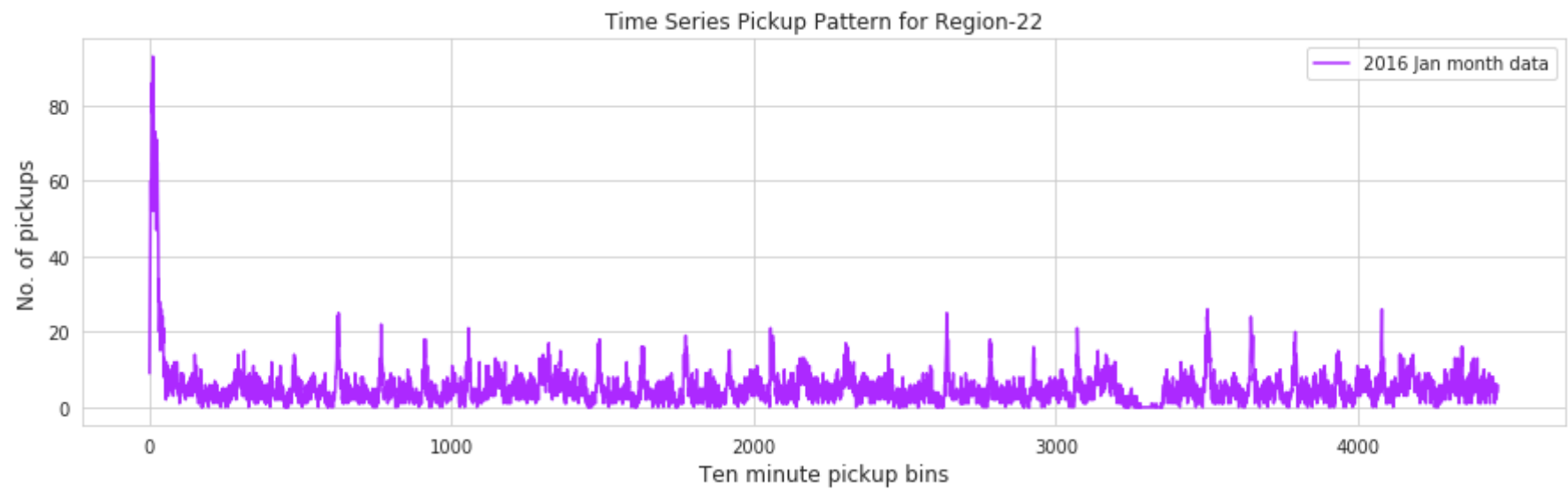


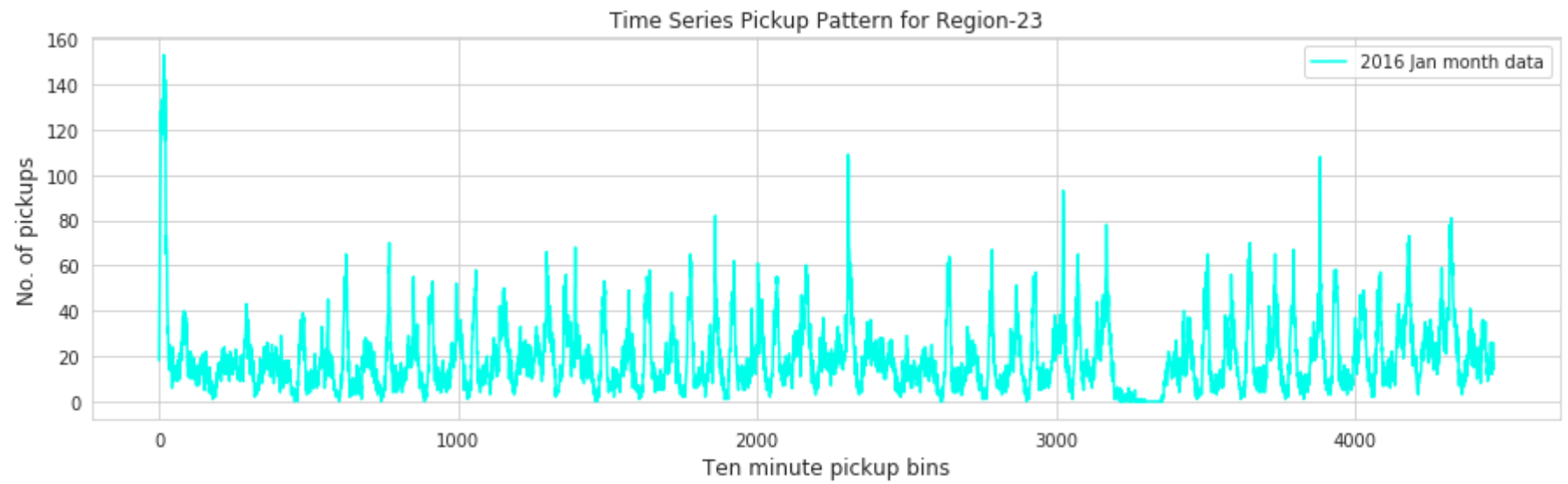


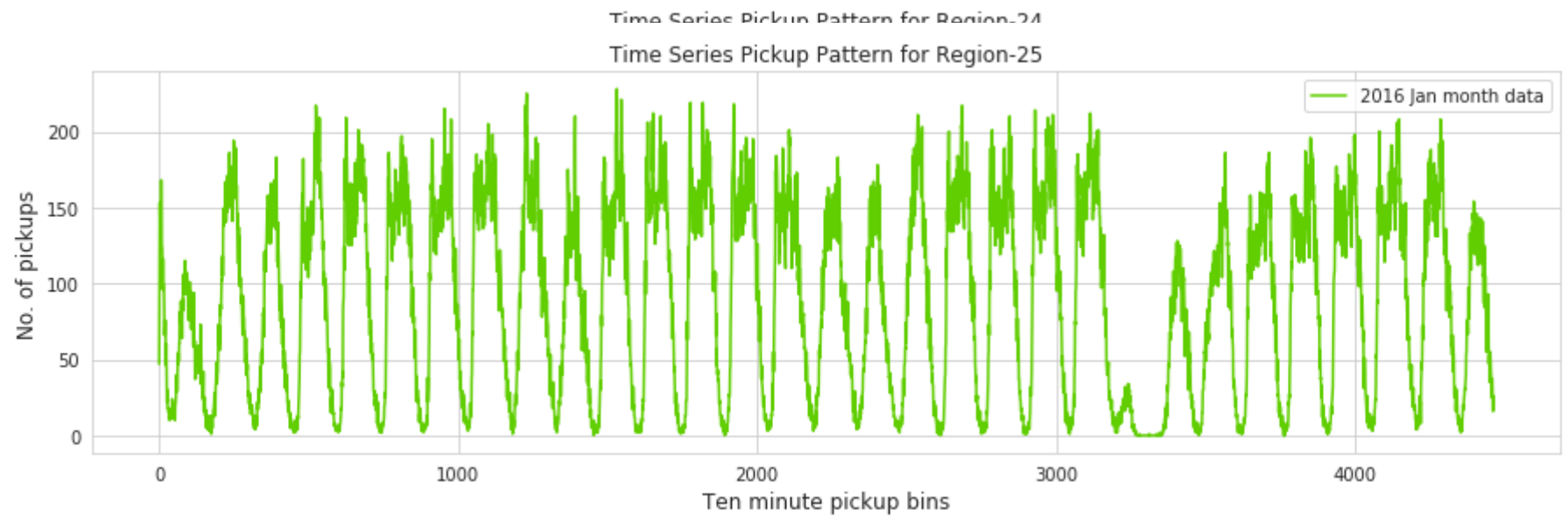


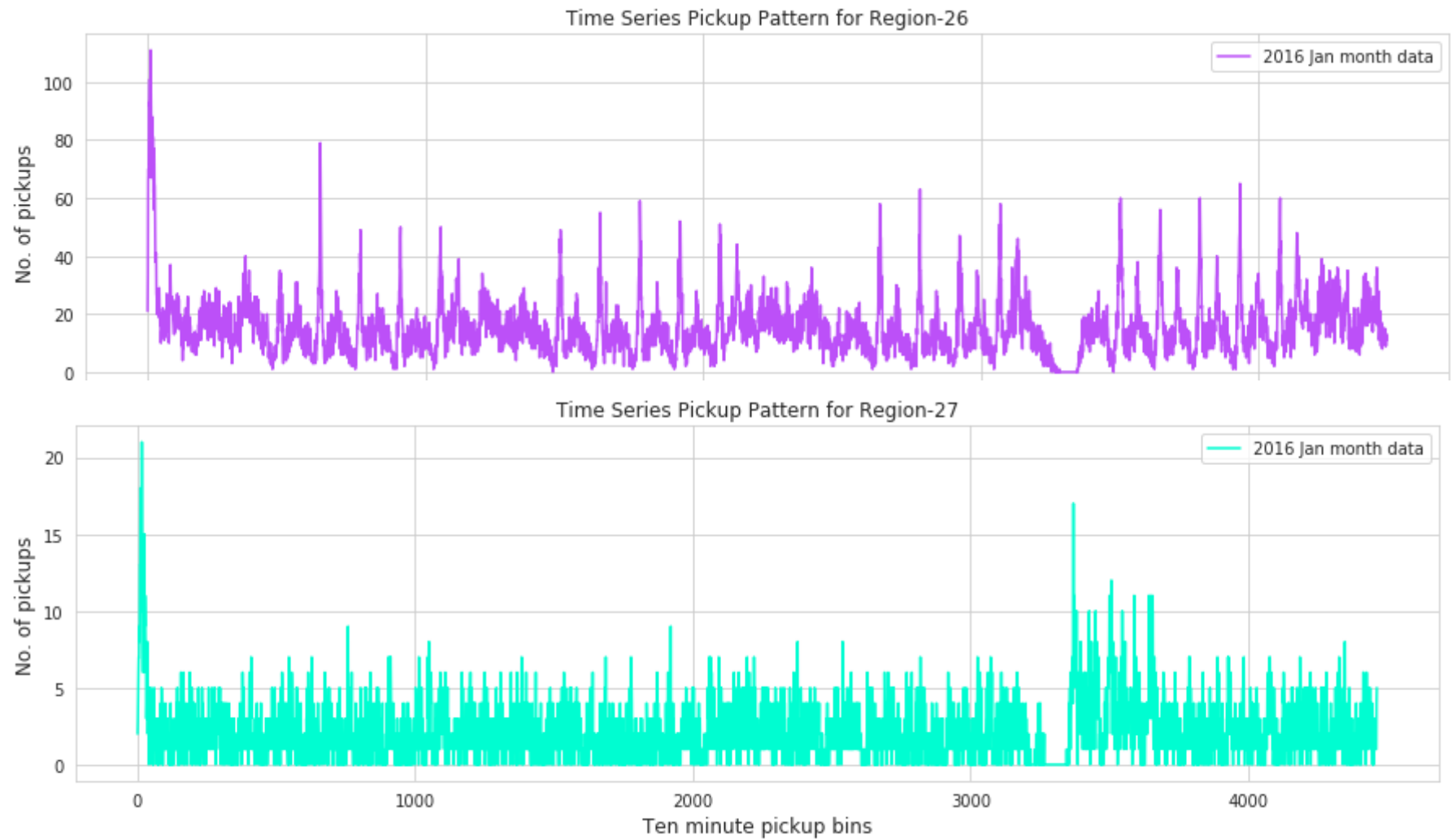


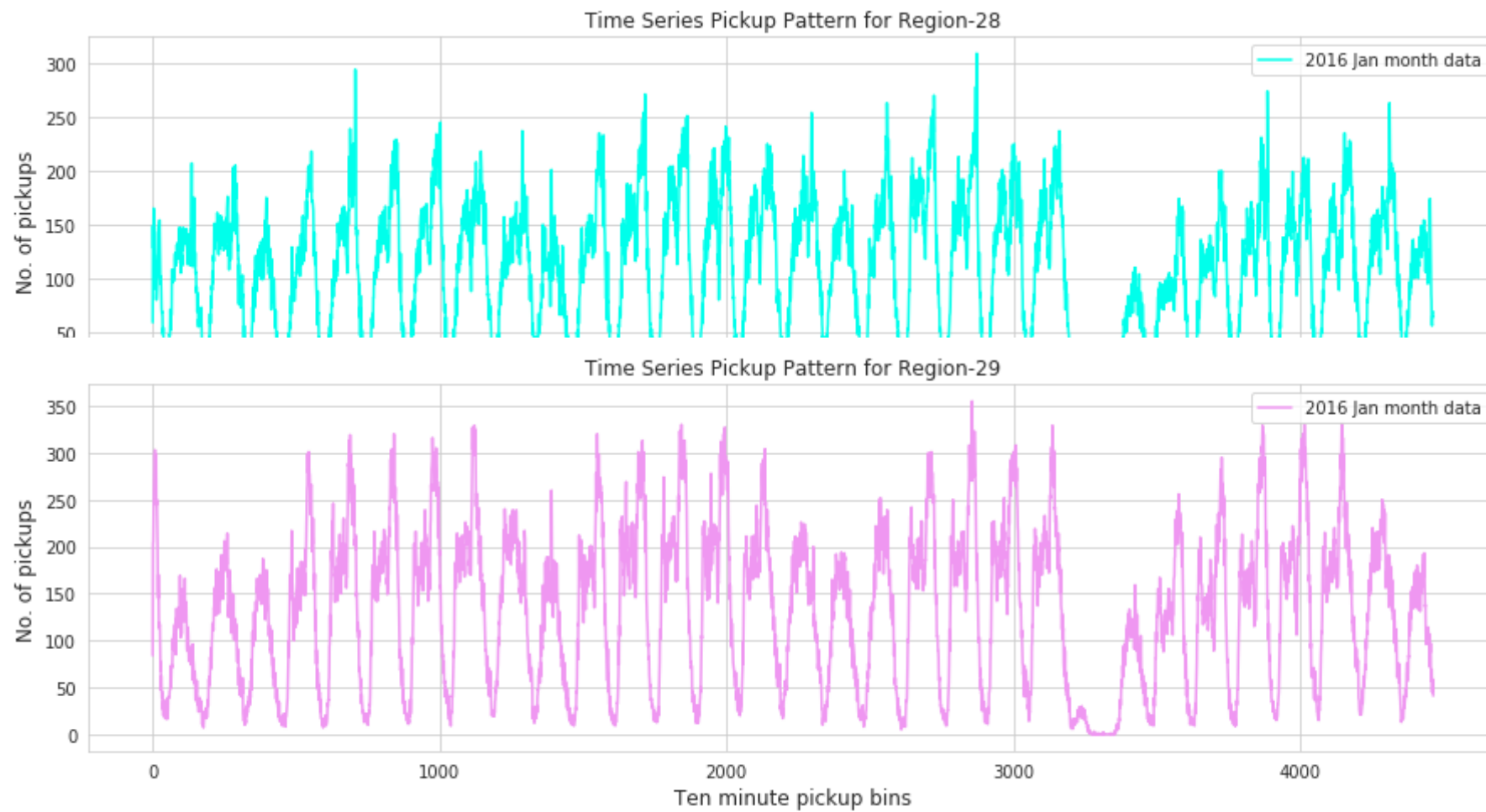




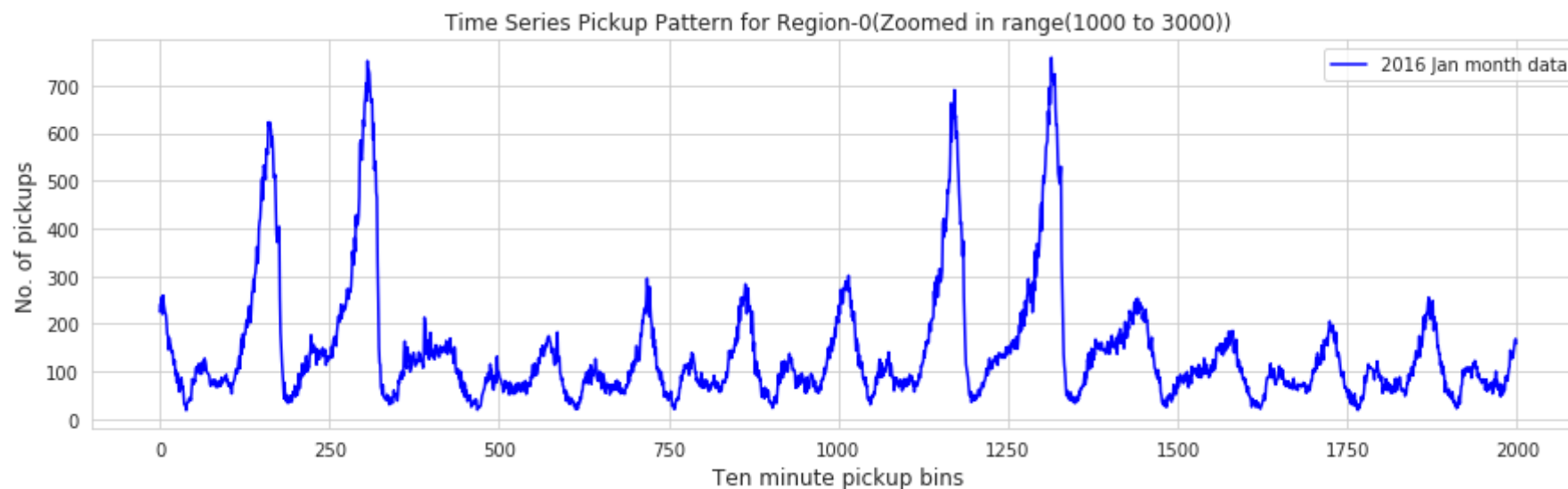








```
In [215]: 1 plt.figure(figsize=(15,4))
2 sns.set_style('whitegrid')
3 plt.title('Time Series Pickup Pattern for Region-0(Zoomed in range(1000 to 3000))', size=12)
4 plt.xlabel('Ten minute pickup bins', fontsize=12)
5 plt.ylabel('No. of pickups', fontsize=12)
6 plt.plot(np.arange(2000),regions_cum[0][1000:3000], color='blue',label='2016 Jan month data')
7 plt.legend()
8 plt.show()
```



Observation:

1. From the above time series for no. of pickups in 24 hours, we observed that it have repetative behaviour so we can use fourier transform to build new features.

[2.2.2.1] Fourier Transform:

Fourier Transform decompose a signal into sum of sine and cosine waves of different amplitudes and frequencies.

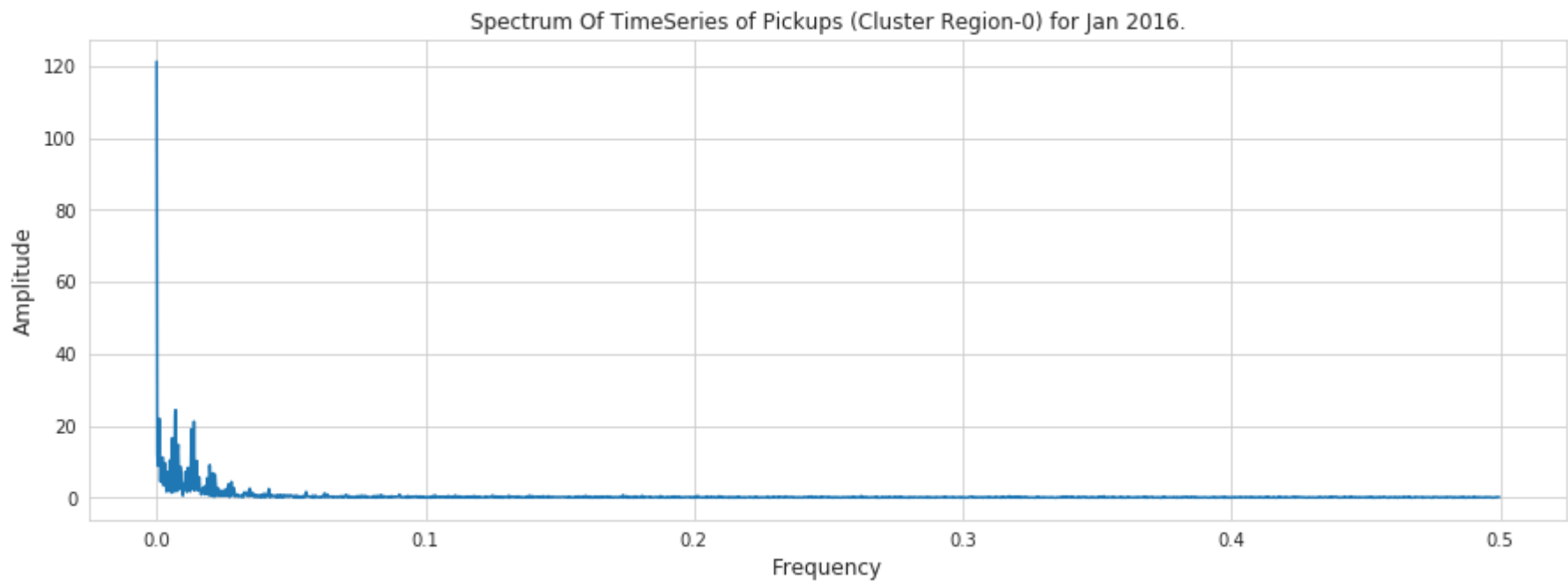
1. Here we are using DFT(using fast fourier transform algorithm) to create the following features:

- a. Top five amplitudes/peaks present in digital signal.
- b. Top five frequencies corresponds to top five amplitudes.
- c. Angle corresponds to top five amplitudes.

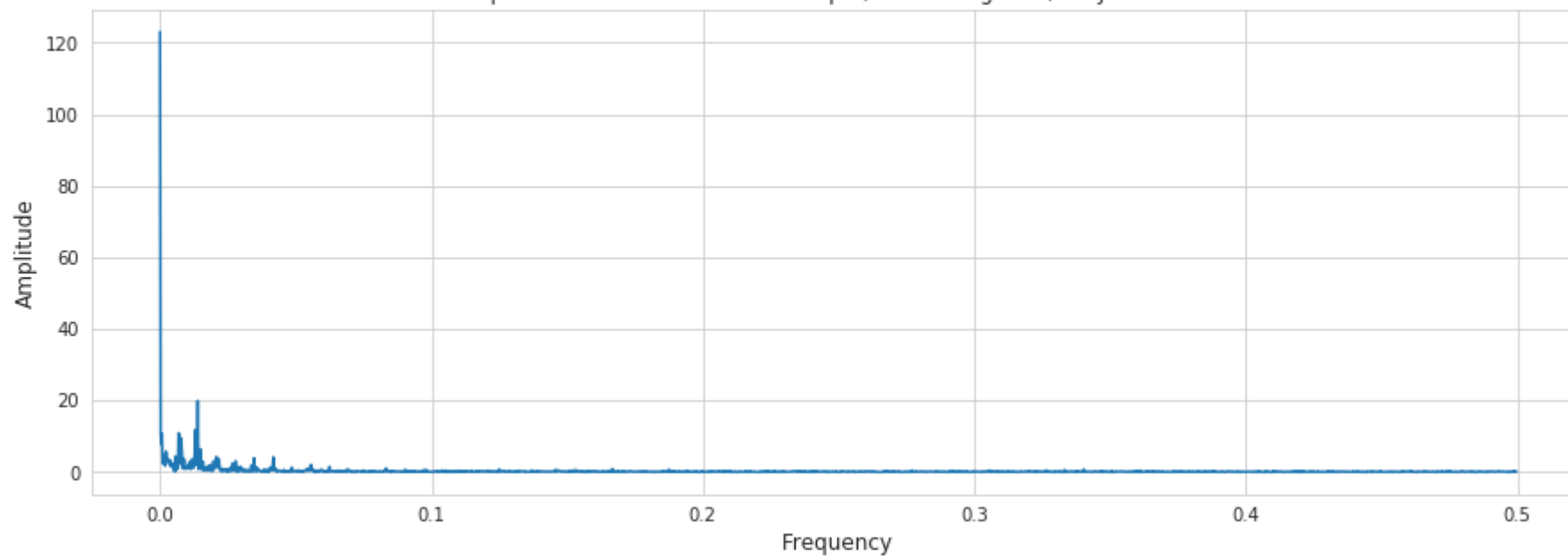
[a.]Spectrum of TimeSeries of Pickups of Region:

In [216]:

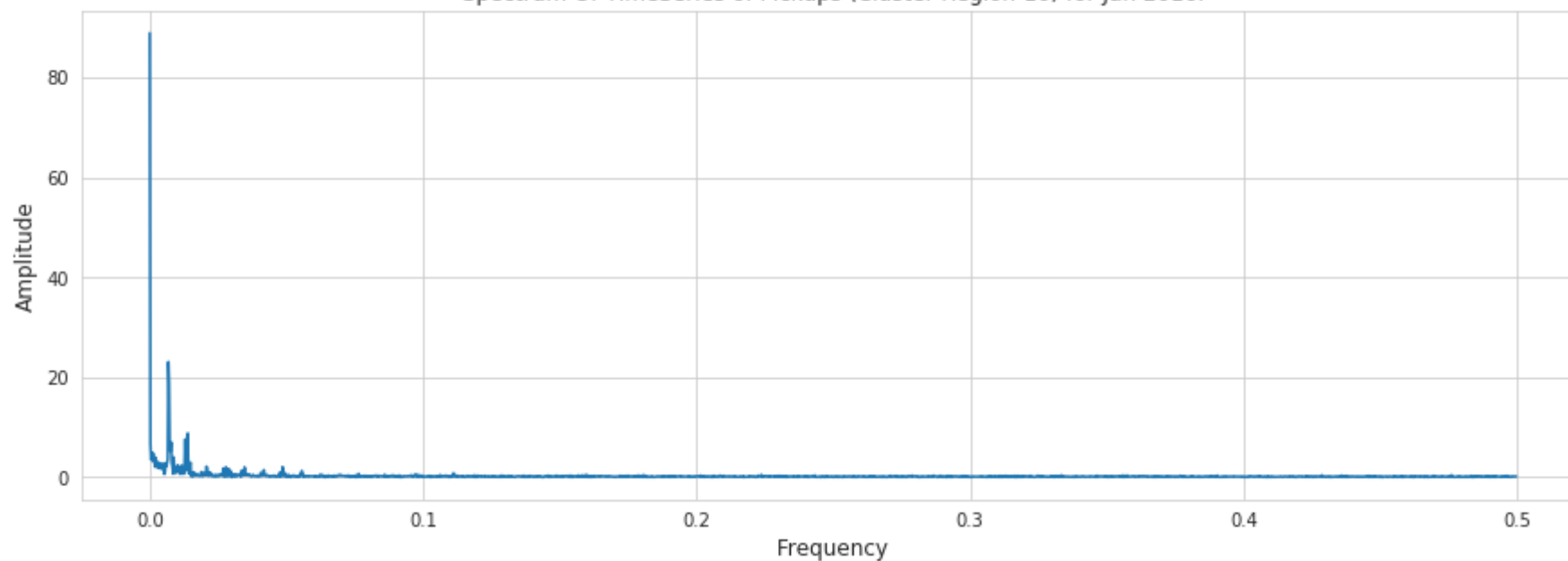
```
1 N=2**12
2 for i in range(0,30,5):
3     Y = np.fft.fft(np.array(regions_cum[i])[0:N])/N
4     # read more about the fftfreq: https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fftfreq.html
5     amp=abs(Y[0:int(N/2)])
6     freq = np.fft.fftfreq(N, 1)[0:int(N/2)]
7     n = len(freq)
8     plt.figure(figsize = (15, 5))
9     plt.plot(freq[:], amp[:])
10    plt.xlabel("Frequency",fontsize=12)
11    plt.ylabel("Amplitude",fontsize=12)
12    plt.title("Spectrum Of TimeSeries of Pickups (Cluster Region-%d) for Jan 2016."%i,size=12)
13    plt.show()
```



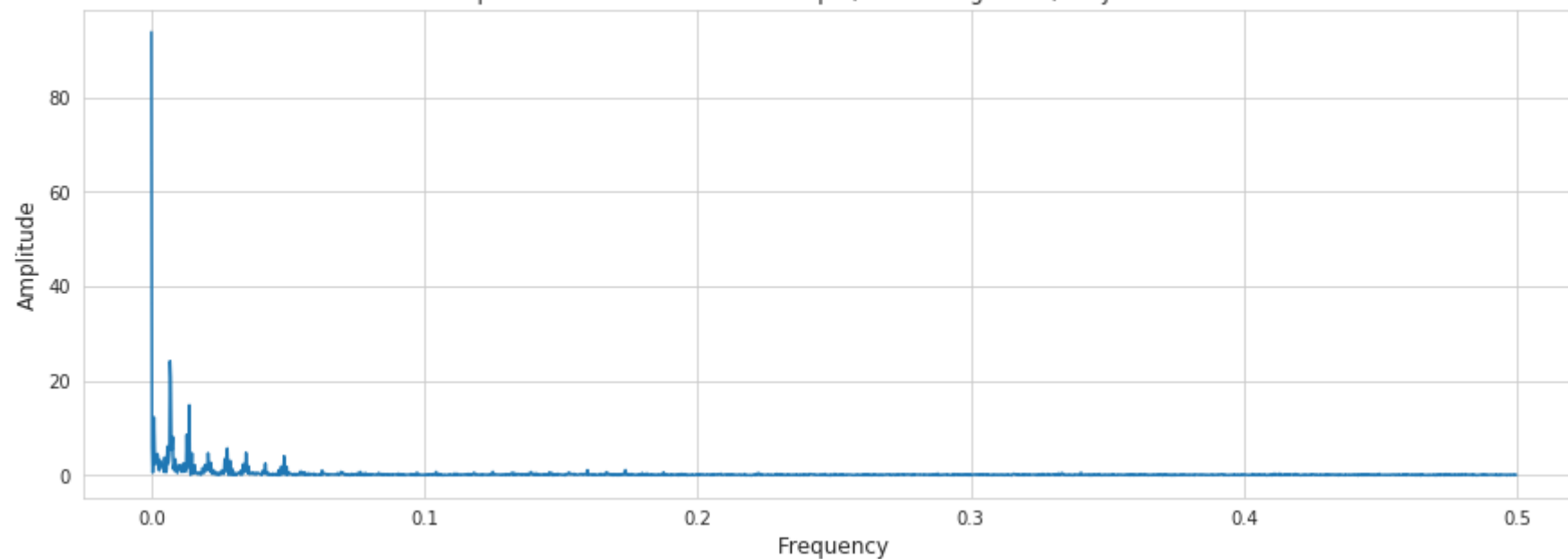
Spectrum Of TimeSeries of Pickups (Cluster Region-5) for Jan 2016.



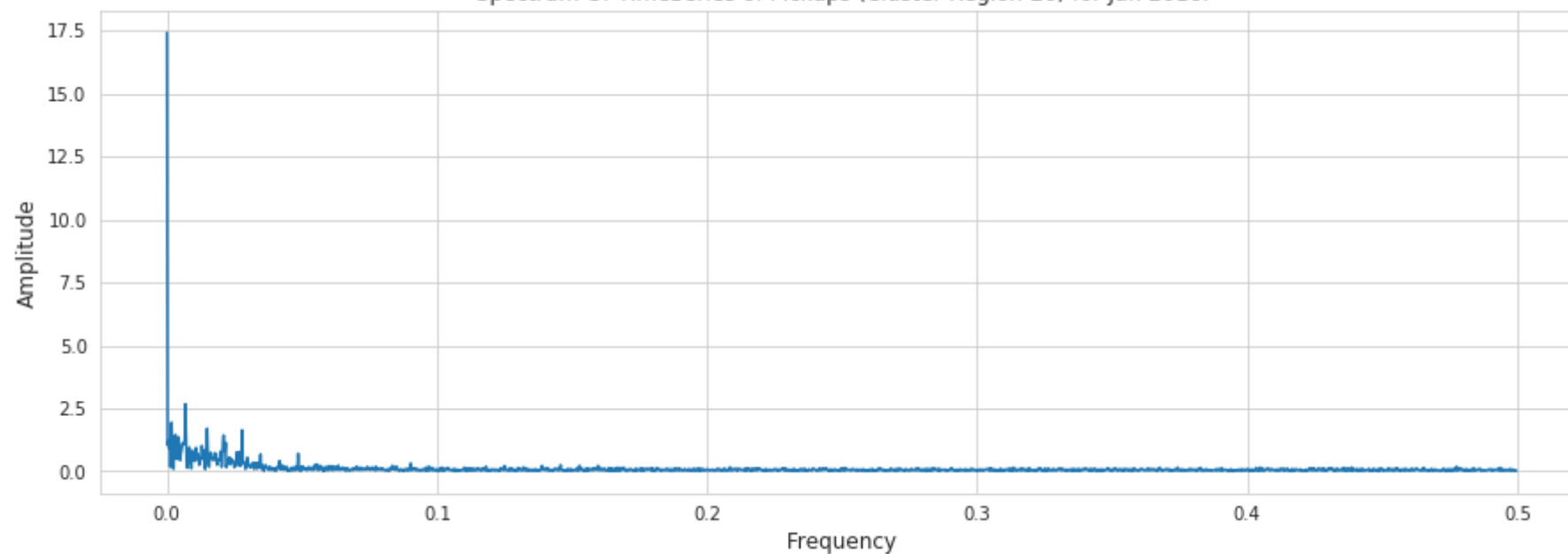
Spectrum Of TimeSeries of Pickups (Cluster Region-10) for Jan 2016.

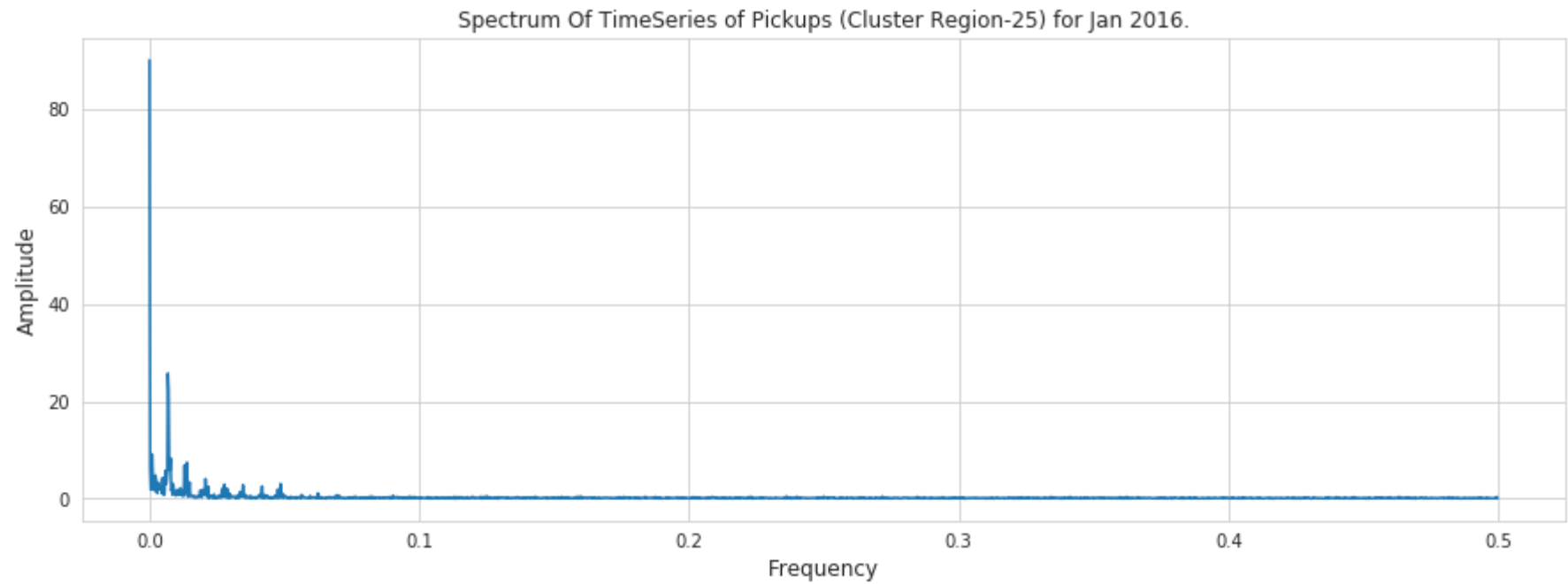


Spectrum Of TimeSeries of Pickups (Cluster Region-15) for Jan 2016.



Spectrum Of TimeSeries of Pickups (Cluster Region-20) for Jan 2016.





[b.]Feature engineering using DFT:

In [217]:

```

1  # getting peaks: https://blog.ytotech.com/2015/11/01/findpeaks-in-python/
2  # read more about fft function : https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fft.html
3  # The FFT algorithm is much more efficient if the number of data points is a power of 2 (128, 512, 1024, etc.)
4  '''
5  sno.    cluster    f1 , f2 , f3 , f4 , f5
6  1        0         x   y   z   a   b
7  2        0         x   y   z   a   b
8  3        0         x   y   z   a   b
9  4        0         x   y   z   a   b
10 5        0         x   y   z   a   b
11 .        .         .   .   .   .   .
12 .        .         .   .   .   .   .
13 .        .         .   .   .   .   .
14 .        .         .   .   .   .   .
15 4459     0         x   y   z   a   b
16 '''
17 N=2**12
18 amplitude_lists=[]
19 frequency_lists=[]
20 angle_lists=[]
21 for i in range(0,30):
22     # divided by N bcz numpy implementation of FFT doesnt consider 1/N
23     Y = np.fft.fft(np.array(regions_cum[i])[0:N])/N
24     # read more about the fftfreq: https://docs.scipy.org/doc/numpy/reference/generated/numpy.fft.fftfreq.html
25     freq = np.fft.fftfreq(N, 1)[1:int(N/2)]
26     # top 5 amplitude
27     amplitude=abs(Y)[1:int(N/2)]
28     # phase of signal
29     angle= np.angle(Y)[1:int(N/2)]
30
31     #TOP 5-MAX AMPLITUDE/PEAKS
32     top_5_amp = amplitude[np.argsort(amplitude)[::-1]][:5]
33     #ANGLE CORRESPONDS TO FOURIER COFFICIENT OF WHICH GIVE MAX AMPLITUDE
34     top_5_angle=angle[np.argsort(amplitude)[::-1]][:5]#np.argsort(angle)
35     #TOP-5 FREQ CORRESPONDING TO MAX AMP
36     top_5_freq = freq[np.argsort(amplitude)[::-1]][:5]
37     # for each cluster we have 5 freq, 5 ampli, 5 phase and these freq,amp,phase are same for all pickup_bin in a cl
38     #bcz we know that a wave have a freq and in our case that wave is TimeSeries of pickups for a cluster
39     for k in range(4459):
40         amplitude_lists.append(top_5_amp)
41         frequency_lists.append(top_5_freq)

```

```

42     angle_lists.append(top_5_angle)
43

```

Holt Winters Model:

1. Each Time series dataset can be decomposed into its components which are Trend, Seasonality and Residual.
2. Datasets which show a similar set of pattern after fixed intervals of a time period suffer from seasonality. In our case timeseries of pickups have some sort of seasonality we can observe from timeseries plot.
3. Holt winter takes into account both trend and seasonality to forecast future prices.

$\ell_t = \alpha(y_t - s_{t-L}) + (1-\alpha)(\ell_{t-1} + b_{t-1})$:level

$b_t = \beta(\ell_t - \ell_{t-1}) + (1-\beta)b_{t-1}$:trend

$s_t = \gamma(y_t - \ell_t) + (1-\gamma)s_{t-L}$:seasonal

$\hat{y}_{t+m} = \ell_t + mb_t + s_{t-L+1+(m-1) \bmod L}$:forecast

where: smoothing parameter $0 \leq \alpha, \beta, \gamma \leq 1$

4. The Holt-Winters seasonal method comprises the forecast equation and three smoothing equations — one for the level ℓ_t , one for trend b_t and one for the seasonal component denoted by s_t , with smoothing parameters α , β and γ .
5.
 - a. level: equation shows a weighted average between the seasonally adjusted observation and the non-seasonal forecast for time t .
 - b. trend: equation is identical to Holt's linear method.
 - c. seasonal: equation shows a weighted average between the current seasonal index, and the seasonal index of the same season last year (i.e., s time periods ago).

In [218]:

```

1  #reference links:
2  #https://www.analyticsvidhya.com/blog/2018/02/time-series-forecasting-methods/
3  #https://grisha.org/blog/2016/01/29/triple-exponential-smoothing-forecasting/
4
5  def initial_trend(series, slen):
6      sum = 0.0
7      for i in range(slen):
8          sum += float(series[i+slen] - series[i]) / slen
9      return sum / slen
10
11 def initial_seasonal_components(series, slen):
12     seasonals = {}
13     season_averages = []
14     n_seasons = int(len(series)/slen)
15     # compute season averages
16     for j in range(n_seasons):
17         season_averages.append(sum(series[slen*j:slen*j+slen])/float(slen))
18     # compute initial values
19     for i in range(slen):
20         sum_of_vals_over_avg = 0.0
21         for j in range(n_seasons):
22             sum_of_vals_over_avg += series[slen*j+i]-season_averages[j]
23         seasonals[i] = sum_of_vals_over_avg/n_seasons
24     return seasonals
25
26
27 def triple_exponential_smoothing(series, slen, alpha, beta, gamma, n_preds):
28     result = []
29     seasonals = initial_seasonal_components(series, slen)
30     for i in range(len(series)+n_preds):
31         if i == 0: # initial values
32             smooth = series[0]
33             trend = initial_trend(series, slen)
34             result.append(series[0])
35             continue
36         if i >= len(series): # we are forecasting
37             m = i - len(series) + 1
38             result.append((smooth + m*trend) + seasonals[i%slen])
39         else:
40             val = series[i]
41             last_smooth, smooth = smooth, alpha*(val-seasonals[i%slen]) + (1-alpha)*(smooth+trend)

```

```
42         trend = beta * (smooth-last_smooth) + (1-beta)*trend
43         seasonals[i%slen] = gamma*(val-smooth) + (1-gamma)*seasonals[i%slen]
44         result.append(smooth+trend+seasonals[i%slen])
45     return result
46
```

In [272]:

```
1  alpha = 0.9
2  beta = 0.2
3  gamma = 0.45
4  season_len = 30
5
6  holts_predicted_list = []
7  for r in range(0,30):
8      holts_predict_values = triple_exponential_smoothing(regions_cum[r][0:4464], season_len, alpha, beta, gamma, 0)
9      holts_predicted_list.append(holts_predict_values[5:])# first 5 used as initial seen points to predict nxt
10 len(holts_predicted_list[0])
```

Out[272]: 4459

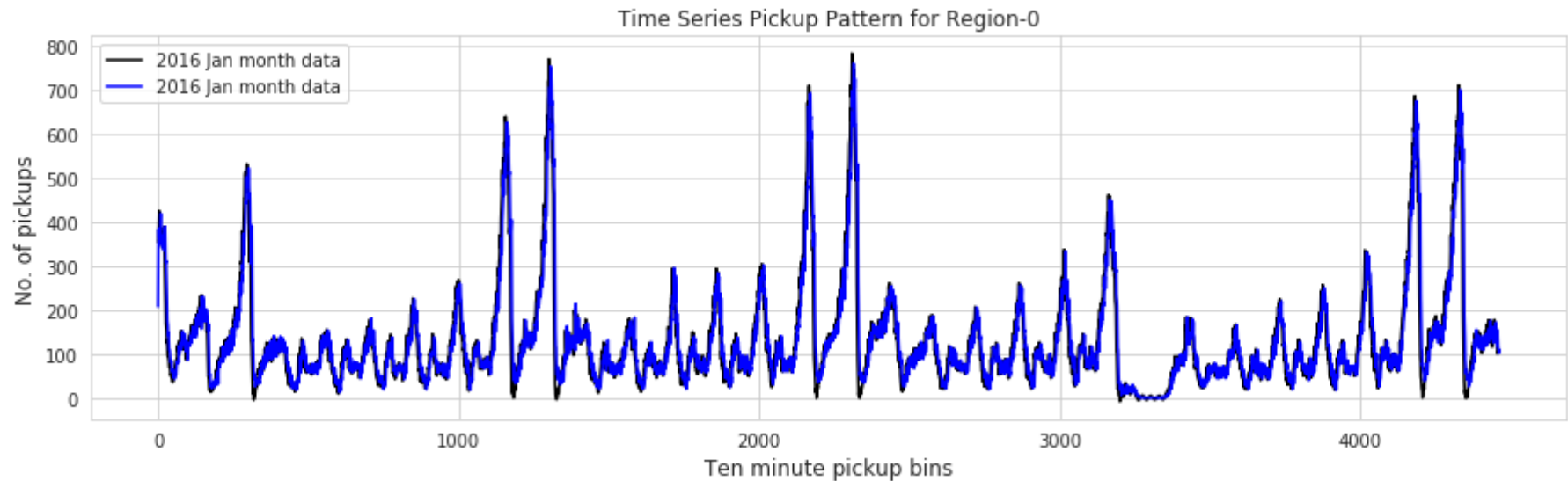
Sainity check for optimal alpha, optimal beta and gamma values by using timeseries:

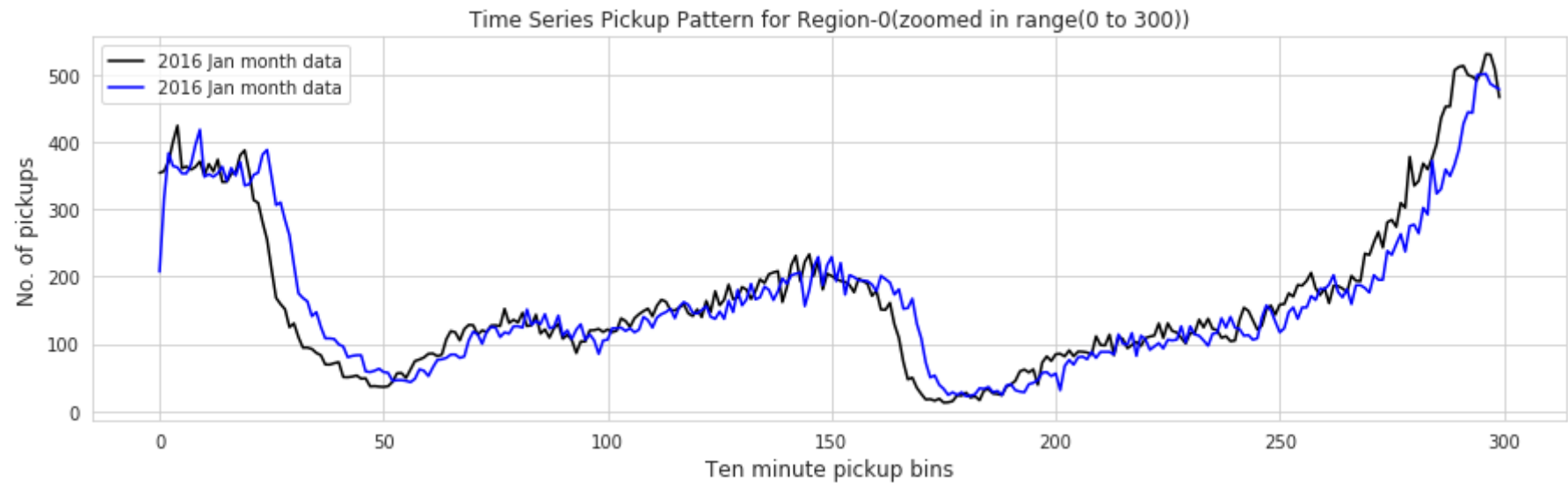
```

In [220]: 1 region_no=int(input('Enter region no.:'))
          2 #range_=int(input('Enter No. of smpls to be plot(<4459):'))
          3 for i in range(2):
          4     plt.figure(i,figsize=(15,4))
          5     sns.set_style('whitegrid')
          6     plt.xlabel('Ten minute pickup bins', fontsize=12)
          7     plt.ylabel('No. of pickups', fontsize=12)
          8
          9     if i ==1:
         10         plt.title('Time Series Pickup Pattern for Region-%d(zoomed in range(0 to 300))'%region_no, size=12)
         11         plt.plot(np.arange(300),predict_list_2[region_no][0:300], color='black',label='2016 Jan month data')
         12         plt.plot(np.arange(300),regions_cum[region_no][:300], color='blue',label='2016 Jan month data')
         13     else:
         14         plt.title('Time Series Pickup Pattern for Region-%d'%region_no, size=12)
         15         plt.plot(np.arange(4459),predict_list_2[region_no][0:4459], color='black',label='2016 Jan month data')
         16         plt.plot(np.arange(4464),regions_cum[region_no][:4464], color='blue',label='2016 Jan month data')
         17     plt.legend()
         18     plt.show()

```

Enter region no.:0





Observation:

1. After trying lot of values for alpha, beta and gamma:

- a. optimal $\alpha = .9$
- b. optimal $\beta = .2$
- c. optimal $\gamma = .45$

For these optimal values we observe that our prediction and actual timeseries are almost overlapped.

[2.3]Train-Test Split

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

In [221]:

```

1 # train, test split : 70% 30% split
2 # 1st 5-pickup_bins are used as given data we used it for predict nxt #pickups
3 # 4459(pickup_bins) * 30(clusters) = 133,770
4 # 4459 * 30 * .7 = 93639 = total train data
5 # 4459 * 30 * .3 = 40131 = total test data
6 # 4459 * .7 = 3121
7 # 4459 * .3 = 1337
8 # Before we start predictions using the tree based regression models we take 3 months of 2016 pickup data
9 # and split it such that for every region we have 70% data in train and 30% in test,
10 # ordered date-wise for every region
11 print("size of train data :", int(4459*30*0.7))
12 print("size of test data :", int(4459*30*0.3))

```

size of train data : 93639

size of test data : 40131

In [222]:

```

1 # extracting first 3121 timestamp values i.e 70% of 4459 (total timestamps) for our training data
2 train_features = [features[i*4459:(4459*i+3121)] for i in range(0,30)]
3 # temp = [0]*(12955 - 9068)
4 test_features = [features[(4459*(i))+3121:4459*(i+1)] for i in range(0,30)]

```

In [223]:

```

1 print("Number of data clusters",len(train_features), "Number of data points in trian data", len(train_features[0]),\
2      "Each data point contains", len(train_features[0][0]),"features")
3 print("Number of data clusters",len(train_features), "Number of data points in test data", len(test_features[0]),\
4      "Each data point contains", len(test_features[0][0]),"features")

```

Number of data clusters 30 Number of data points in trian data 3121 Each data point contains 5 features

Number of data clusters 30 Number of data points in test data 1338 Each data point contains 5 features

In [250]:

```

1 train_frequencies = [frequency_lists[i*4459:(4459*i+3121)] for i in range(30)]
2 test_frequencies = [frequency_lists[(i*4459)+3121:(4459*(i+1))] for i in range(30)]
3
4 train_amplitudes = [amplitude_lists[i*4459:(4459*i+3121)] for i in range(30)]
5 test_amplitudes = [amplitude_lists[(i*4459)+3121:(4459*(i+1))] for i in range(30)]
6
7 train_angles = [angle_lists[i*4459:(4459*i+3121)] for i in range(30)]
8 test_angles = [angle_lists[(i*4459)+3121:(4459*(i+1))] for i in range(30)]
9

```

```
In [251]: 1 train_amplitudes[0]
          2 len(test_angles[0])
```

Out[251]: 1338

```
In [273]: 1 # extracting first 9169 timestamp values i.e 70% of 13099 (total timestamps) for our training data
          2 tsne_train_flat_lat = [i[:3121] for i in tsne_lat] # we have 30 sublist containing 3121 value in each list
          3 tsne_train_flat_lon = [i[:3121] for i in tsne_lon]
          4 tsne_train_flat_weekday = [i[:3121] for i in tsne_weekday]
          5 tsne_train_flat_output = [i[:3121] for i in output]
          6 tsne_train_flat_exp_avg = [i[:3121] for i in predict_list]
          7 tsne_train_holts = [i[:3121] for i in holts_predicted_list]
```

```
In [274]: 1 # extracting the rest of the timestamp values i.e 30% of 12956 (total timestamps) for our test data
          2 tsne_test_flat_lat = [i[3121:] for i in tsne_lat]
          3 tsne_test_flat_lon = [i[3121:] for i in tsne_lon]
          4 tsne_test_flat_weekday = [i[3121:] for i in tsne_weekday]
          5 tsne_test_flat_output = [i[3121:] for i in output]
          6 tsne_test_flat_exp_avg = [i[3121:] for i in predict_list]
          7 tsne_test_holts = [i[3121:] for i in holts_predicted_list]
```



```
In [275]: 1 # the above contains values in the form of list of lists (i.e. list of values of each region), here we make all of t
2 train_new_features = []
3 test_new_features = []
4 train_amp = []
5 test_amp = []
6 train_freq = []
7 test_freq = []
8 train_ang = []
9 test_ang = []
10 train_holt = []
11 test_holt = []
12 #creating single list
13 for i in range(0,30):
14     train_new_features.extend(train_features[i])
15     test_new_features.extend(test_features[i])
16     train_amp.extend(train_amplitudes[i])
17     test_amp.extend(test_amplitudes[i])
18     train_freq.extend(train_frequencies[i])
19     test_freq.extend(test_frequencies[i])
20     train_ang.extend(train_angles[i])
21     test_ang.extend(test_angles[i])
22     train_holt.extend(tsne_train_holts[i])
23     test_holt.extend(tsne_test_holts[i])
```

```
In [276]: 1 len(train_holt)
```

Out[276]: 93630

```
In [277]: 1 train_data = np.hstack((train_freq, train_amp, train_ang, train_new_features))
2 test_data = np.hstack((test_freq, test_amp, test_ang, test_new_features))
```

In [279]:

```
1 # converting lists of lists into single list i.e flatten
2 # a = [[1,2,3,4],[4,6,7,8]]
3 # print(sum(a,[]))
4 # [1, 2, 3, 4, 4, 6, 7, 8]
5
6 tsne_train_lat = sum(tsne_train_flat_lat, [])
7 tsne_train_lon = sum(tsne_train_flat_lon, [])
8 tsne_train_weekday = sum(tsne_train_flat_weekday, [])
9 tsne_train_output = sum(tsne_train_flat_output, [])
10 tsne_train_exp_avg = sum(tsne_train_flat_exp_avg, [])
11 tsne_train_holts_ = sum(tsne_train_holts, [])
```

In [280]:

```
1 # converting lists of lists into sinle list i.e flatten
2 # a = [[1,2,3,4],[4,6,7,8]]
3 # print(sum(a,[]))
4 # [1, 2, 3, 4, 4, 6, 7, 8]
5
6 tsne_test_lat = sum(tsne_test_flat_lat, [])
7 tsne_test_lon = sum(tsne_test_flat_lon, [])
8 tsne_test_weekday = sum(tsne_test_flat_weekday, [])
9 tsne_test_output = sum(tsne_test_flat_output, [])
10 tsne_test_exp_avg = sum(tsne_test_flat_exp_avg, [])
11 tsne_test_holts_ = sum(tsne_test_holts, [])
```

In [281]:

```
1 # Preparing the data frame for our train data
2 columns = ['freq_1', 'freq_2', 'freq_3', 'freq_4', 'freq_5', \
3            'amp_1', 'amp_2', 'amp_3', 'amp_4', 'amp_5', \
4            'ang_1', 'ang_2', 'ang_3', 'ang_4', 'ang_5', \
5            'ft_5', 'ft_4', 'ft_3', 'ft_2', 'ft_1']
6            #'ang_1', 'ang_2', 'ang_3', 'ang_4', 'ang_5', \
7 x_train = pd.DataFrame(data=train_data, columns=columns)
8 x_train['lat'] = tsne_train_lat
9 x_train['lon'] = tsne_train_lon
10 x_train['weekday'] = tsne_train_weekday
11 x_train['exp_avg'] = tsne_train_exp_avg
12 x_train['holt_triplet_avg'] = tsne_train_holts_
13 print(x_train.shape)
```

(93630, 25)

In [282]: 1 x_train.head(5)

Out[282]:

	freq_1	freq_2	freq_3	freq_4	freq_5	amp_1	amp_2	amp_3	amp_4	amp_5	...	ft_5	ft_4	ft_3	ft_2	ft_1	
0	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	207.0	315.0	383.0	364.0	362.0	40.72
1	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	315.0	383.0	364.0	362.0	353.0	40.72
2	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	383.0	364.0	362.0	353.0	353.0	40.72
3	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	364.0	362.0	353.0	353.0	366.0	40.72
4	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	362.0	353.0	353.0	366.0	395.0	40.72

5 rows × 25 columns

```
In [283]: 1 # Preparing the data frame for our train data
2 x_test = pd.DataFrame(data=test_data, columns=columns)
3 x_test['lat'] = tsne_test_lat
4 x_test['lon'] = tsne_test_lon
5 x_test['weekday'] = tsne_test_weekday
6 x_test['exp_avg'] = tsne_test_exp_avg
7 x_test['holt_triplet_avg'] = tsne_test_holts_
8 print(x_test.shape)
```

(40140, 25)

In [284]:

```
1 x_test.head()
```

Out[284]:

	freq_1	freq_2	freq_3	freq_4	freq_5	amp_1	amp_2	amp_3	amp_4	amp_5	...	ft_5	ft_4	ft_3	ft_2	ft_1	
0	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	70.0	66.0	64.0	81.0	65.0	40.725
1	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	66.0	64.0	81.0	65.0	104.0	40.725
2	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	64.0	81.0	65.0	104.0	102.0	40.725
3	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	81.0	65.0	104.0	102.0	109.0	40.725
4	0.00708	0.006836	0.000977	0.013916	0.012939	24.4252	22.911823	22.06076	21.227861	19.011012	...	65.0	104.0	102.0	109.0	121.0	40.725

5 rows × 25 columns

[2.4]Utility functions required for regression models

[a.]Function for Standardizing data

In [285]:

```
1 def std_data(train,test,mean):
2     scaler=StandardScaler(with_mean=mean)
3     std_train=scaler.fit_transform(train)
4     std_test=scaler.transform(test)
5     return std_train, std_test
```

[b.]Function to print parameters summary:

In [286]:

```
1 def param_list(params, para_summ):
2     para_summ.clear_rows()
3     for name, val in zip(params.keys(),params.values()):
4         para_summ.add_row([name, val])
5     print(para_summ)
```

[c.]Function for feature importance:

```
In [287]: 1 def plot_importance(model, clf):
2         fig = plt.figure(figsize = (8, 6))
3         ax = fig.add_axes([0,0,1,1])
4         model.plot_importance(clf, ax = ax, height = 0.3)
5         plt.xlabel("F Score", fontsize = 15)
6         plt.ylabel("Features", fontsize = 15)
7         plt.title("Feature Importance", fontsize = 15)
8         plt.tick_params(labelsize = 15)
9
10        plt.show()
```

[d.]function for scorer for Grid search and Random search:

```
In [288]: 1  #(mean_absolute_error(y_true, y_pred))/(sum(y_true)/len(y_pred))*100
2 def mape_scorer(y_true, y_pred):
3      #to avoid division by zero error
4     eps=.000001
5     y_true, y_pred = np.array(y_true), np.array(y_pred)
6     mape=(mean_absolute_error(y_true, y_pred))/(sum(y_true)/len(y_pred))*100
7      #mape = (mean_absolute_error(y_true,y_pred)/sum(y_true))*(100/len(y_pred))
8     return mape
```

[e.]Function for hyperparam tuning for linear regression:

Linear-Regression

In [289]:

```

1  def linearRegressor(x_train,y_train,CV,params_):
2      clf=SGDRegressor(loss='squared_loss', penalty='l2')
3      model=GridSearchCV(clf,\
4                          param_grid=params_,\
5                          n_jobs=-1,\
6                          return_train_score=True,\
7                          scoring=make_scorer(mape_scorer, greater_is_better=False),\
8                          cv=CV
9      )
10     model.fit(x_train, y_train)
11     train_mape = model.cv_results_['mean_train_score']
12     cv_mape = model.cv_results_['mean_test_score']
13     #if len(param_name.split())==1:
14     for i in range(len(train_mape)):
15         print('Train MAPE: %.4f'%abs(train_mape[i]), 'CV MAPE: %.4f'%abs(cv_mape[i]))
16     print()
17     params_=model.best_params_['alpha']
18     return model, params_

```

[f.]Functions required for XGBOOST and RANDOMFOREST-REGRESSOR model :

XGBOOST (GBDT)

2. function for Tunning hypeparam:

In [290]:

```

1 def Ensemble_Regressor(x_train,y_train,CV,params_, tune_param, searchMethod, regressor_used):
2     #INITIALIZE GBDT CLASSIFIER
3     if regressor_used=='xgb':
4         reg=xgb.XGBRegressor(n_estimators=params_xgb['n_estimators'],\
5                               max_depth=params_xgb['max_depth'],\
6                               eta=.02,\
7                               reg_alpha=params_xgb['reg_alpha'],\
8                               min_child_weight=params_xgb['min_child_weight'],\
9                               gamma=params_xgb['gamma'],\
10                              subsample=params_xgb['subsample'],\
11                              colsample_bytree=params_xgb['colsample_bytree'],\
12                              booster='gbtree'
13                               )
14     elif regressor_used=='rf':
15         reg=RandomForestRegressor(n_estimators=params_rf['n_estimators'],\
16                                   max_depth=params_rf['max_depth'],\
17                                   max_features = 'sqrt'
18                                   )
19     # APPLY RANDOM OR GRID SEARCH FOR HYPERPARAMETER TUNNING
20     if searchMethod=='random':
21         model=RandomizedSearchCV(reg,\
22                                   n_jobs=-1,\
23                                   cv=CV,\
24                                   param_distributions=tune_param,\
25                                   n_iter=6,\
26                                   return_train_score=True,\
27                                   scoring=make_scorer(mape_scorer, greater_is_better=False))
28     elif searchMethod=='grid':
29         model= GridSearchCV(estimator=reg,\
30                              param_grid=tune_param,\
31                              scoring=make_scorer(mape_scorer, greater_is_better=False),\
32                              n_jobs=-1,\
33                              cv=CV,\
34                              return_train_score=True
35                              )
36     model.fit(x_train,y_train)
37     train_mape = model.cv_results_['mean_train_score']
38     cv_mape = model.cv_results_['mean_test_score']
39     #if len(param_name.split())==1:
40     for i in range(len(train_mape)):
41         print('Train MAPE: %.4f'%abs(train_mape[i]),'CV MAPE: %.4f'%abs(cv_mape[i]))

```

```
42     print()
43     return model
```

```
In [291]: 1 def tuneALL_PARAM_XGB(train, y_train, CV, params_range, params_, searchMethod, regressor_used):
2         for param_name, param_list in zip(params_range.keys(),params_range.values()):
3             tune_param={}
4             tune_param[param_name]=param_list
5             #TUNNING HYPERPARAM
6             print('Tunning {}:'.format(param_name.upper()))
7             model=Ensemble_Regressor(train, y_train, CV, params_, tune_param, searchMethod, regressor_used)
8             #UPDATE OTIMAL VALUE OF PARAMETER
9             params_[param_name] = model.best_params_[param_name]
10        return model, params_
```

3. function for measuring perfomance on test data:

In [292]:

```

1 def test_performance_xgb(x_train,y_train,x_test,y_test,params_, model_summary,model_use=None,summary=False,regressor
2     '''FUNCTION FOR TEST PERFORMANCE(PLOT ROC CURVE FOR BOTH TRAIN AND TEST) WITH OPTIMAL HYPERPARAM'''
3     #INITIALIZE GBDT WITH OPTIMAL VALUE OF HYPERPARAMS
4     if regressor_used=='xgb':
5         clf=xgb.XGBRegressor(n_estimators=params_['n_estimators'],\
6                               max_depth=params_['max_depth'],\
7                               eta=.02,\
8                               reg_alpha=params_['reg_alpha'],\
9                               min_child_weight=params_['min_child_weight'],\
10                              gamma=params_['gamma'],\
11                              subsample=params_['subsample'],\
12                              colsample_bytree=params_['colsample_bytree'],\
13                              booster='gbtree',\
14                              verbose=1)
15     elif regressor_used=='rf':
16         clf=RandomForestRegressor(n_estimators=params_['n_estimators'],\
17                                   max_depth=params_['max_depth'],\
18                                   max_features='sqrt'
19         )
20     elif regressor_used=='linear':
21         clf=SGDRegressor(loss='squared_loss',\
22                           alpha=params_,\
23                           shuffle=False
24         )
25
26     clf.fit(x_train, y_train)
27     #PREDICTION FOR TRAIN AND TEST
28     y_pred=clf.predict(x_test)
29     y_pred_tr=clf.predict(x_train)
30
31     #TEST MAPE
32     test_mape=mape_scorer(y_test, y_pred)#np.mean(np.abs((y_test - y_pred) / y_test)) * 100
33     #TRAIN MAPE
34     train_mape=mape_scorer(y_train, y_pred_tr)#np.mean(np.abs((y_train - y_pred_tr) / y_train)) * 100
35
36     print('FOR OPTIMAL PARAMETERS, TRAIN MAPE: %.5f, TEST MAPE: %.5f'%(train_mape, test_mape))
37     model_summ_local=PrettyTable()
38     model_summ_local.field_names=['Model', 'Train(MAPE)', 'Test(MAPE)']
39     model_summ_local.add_row([model_use, '%.5f'%train_mape, '%.5f'%test_mape])
40
41     if summary:

```

```
42     model_summary.add_row([model_use, '%.5f'%train_mape, '%.5f'%test_mape])
43     if regressor_used=='xgb':
44         plot_importance(xgb, clf)
45         return model_summ_local
46     elif regressor_used=='rf':
47         plt.figure(1,figsize=(11,10))
48         sns.set_style('whitegrid')
49         plt.title('Feature Importances',size=15)
50         plt.xlabel('Fscore',fontsize=15)
51         plt.ylabel('Features',fontsize=15)
52         plt.barh(range(len(clf.feature_importances_)), clf.feature_importances_,tick_label=train.columns)
53         plt.show()
54         return model_summ_local,clf
55     return model_summ_local,clf
```

=====

Initialization of common objects:

In [346]:

```

1  #OBJECT FOR TIMESERIES CROSS VALIDATION
2  TBS=TimeSeriesSplit(n_splits=5)
3
4  #CROSSVALIDATION ALGO TO BE USED
5  searchMethod=['random', 'grid']
6
7  #MODEL USED
8  model_name=['EWMA-PREVIOUS-DATA', 'LINEAR-REGRESSOR', 'RANDOM-FOREST-REGRESSOR', 'XGBOOST-REGRESSOR', 'HOLTS-WINTER_MC
9
10 #GLOBAL SUMMARY OF ALL THE MODELS
11 model_summary=PrettyTable()
12 model_summary.field_names=['Model', 'Train(MAPE)', 'Test(MAPE)']
13
14 #DICT OF PARAMETERS, INITIALLY SET REASONABLE VALUES FOR PARAMETER, AND AFTER TUNNING UPDATE VALUE WITH OPT. VALUE
15 params_xgb=OrderedDict([
16     ('n_estimators',128),
17     ('max_depth',5),
18     ('min_child_weight',1),
19     ('gamma',0),
20     ('subsample',.8),
21     ('colsample_bytree',.8),
22     ('reg_alpha',.1)]
23 )
24 params_rf=OrderedDict([
25     ('n_estimators',50),
26     ('max_depth',8),
27     ('min_samples_leaf',1)]
28 )
29
30 # DICT OF HYPERPARAMETER FOR XGBOOST WITH RANGE OF VALUES
31 params_range_xgb=OrderedDict([
32     ('n_estimators', [128,256,512,650]),\
33     ('max_depth', [5,7,9]),\
34     ('min_child_weight', [1,3,5]),\
35     ('gamma', [i/10.0 for i in range(0,5)]),\
36     ('subsample', [.6,.7,.8,.9]),\
37     ('colsample_bytree', [.6,.7,.8,.9]),\
38     ('reg_alpha',[0, 0.001, 0.005, 0.01, 0.05]))]
39
40 # DICT OF HYPERPARAMETER FOR RANDOM-FOREST WITH RANGE OF VALUES
41 params_range_rf=OrderedDict([

```

```

42     ('n_estimators', [20,40,80,128,256,512]),\
43     ('max_depth', [9,12,15,18,20,25,27]))
44
45 #PARAMETER SUMMARY
46 param_summ=PrettyTable()
47 param_summ.field_names=['Parameter', 'Value']

```

Applying Machine Learning Models:

0.Holt Winters Model:

In [347]:

```

1 y_test_act=tsne_test_output; y_test_pred=x_test['holt_triplet_avg']
2 y_train_act=tsne_train_output; y_train_pred=x_train['holt_triplet_avg']
3 mape_hwm_test=mape_scorer(y_test_act, y_test_pred)
4 mape_hwm_train=mape_scorer(y_train_act, y_train_pred)
5 #ADD TO GLOBAL SUMMARY
6 model_summ_local.add_row(['HOLTS-WINTER_MODEL', '%.4f'%mape_hwm_train, '%.4f'%mape_hwm_test])
7 model_summary.add_row(['HOLTS-WINTER_MODEL', '%.4f'%mape_hwm_train, '%.4f'%mape_hwm_test])
8 print('TRAIN MAPE: %.4f'%mape_hwm_train, 'TEST MAPE: %.4f'%mape_hwm_test)

```

TRAIN MAPE: 3.2515 TEST MAPE: 3.3127

[0.1] Summary:

In [348]:

```
1 print(model_summ_local)
```

```

+-----+-----+-----+
|      Model      | Train(MAPE) | Test(MAPE) |
+-----+-----+-----+
| HOLTS-WINTER_MODEL |    3.2515   |    3.3127   |
+-----+-----+-----+

```

1. EWMA(Exponential Weighted Mean Avg):

In [349]:

```

1 y_test_act=tsne_test_output; y_test_pred=x_test['exp_avg']
2 y_train_act=tsne_train_output; y_train_pred=x_train['exp_avg']
3 mape_ewma_test=mape_scorer(y_test_act, y_test_pred)
4 mape_ewma_train=mape_scorer(y_train_act, y_train_pred)
5 #ADD TO GLOBAL SUMMARY
6 model_summ_local.clear_rows()
7 model_summ_local.add_row([model_name[0], '%.4f'%mape_ewma_train, '%.4f'%mape_ewma_test])
8 model_summary.add_row([model_name[0], '%.4f'%mape_ewma_train, '%.4f'%mape_ewma_test])
9 print('TRAIN MAPE: %.4f'%mape_ewma_train, 'TEST MAPE: %.4f'%mape_ewma_test)

```

TRAIN MAPE: 12.1853 TEST MAPE: 12.0552

[1.1] Summary:

In [350]:

```
1 print(model_summ_local)
```

```

+-----+-----+-----+
|      Model      | Train(MAPE) | Test(MAPE) |
+-----+-----+-----+
| EWMA-PREVIOUS-DATA | 12.1853    | 12.0552    |
+-----+-----+-----+

```



2.Linear Regressor:

[2.1]Hyperparameter Tunning:

In [351]:

```

1  %%time
2  params_={'alpha':[10**x for x in range(-10,5)]}
3  # STANDARDIZED TRAIN AND TEST DATA
4  train, test = std_data(x_train, x_test, mean=True)
5  #train, test = x_train, x_test
6  y_train, y_test = tsne_train_output, tsne_test_output
7  print('HYPERPARAMETER:\n')
8  param_list(params_, param_summ)
9  print()
10 model, params_ = linearRegressor(train, y_train, TBS, params_)

```

HYPERPARAMETER:

Parameter	Value
alpha	[1e-10, 1e-09, 1e-08, 1e-07, 1e-06, 1e-05, 0.0001, 0.001, 0.01, 0.1, 1, 10, 100, 1000, 10000]

```

Train MAPE: 1.1597 CV MAPE: 23.9725
Train MAPE: 1.1565 CV MAPE: 23.9617
Train MAPE: 1.1577 CV MAPE: 23.9790
Train MAPE: 1.1570 CV MAPE: 23.9903
Train MAPE: 1.1553 CV MAPE: 23.9829
Train MAPE: 1.1604 CV MAPE: 23.9270
Train MAPE: 1.1629 CV MAPE: 23.9184
Train MAPE: 1.1939 CV MAPE: 23.1511
Train MAPE: 2.0853 CV MAPE: 18.3379
Train MAPE: 6.7949 CV MAPE: 8.9453
Train MAPE: 13.7218 CV MAPE: 14.9495
Train MAPE: 38.7721 CV MAPE: 51.1496
Train MAPE: 60.2381 CV MAPE: 81.0606
Train MAPE: 65.4838 CV MAPE: 87.7927
Train MAPE: 65.7183 CV MAPE: 88.4724

```

```

CPU times: user 2min 25s, sys: 1.31 s, total: 2min 26s
Wall time: 2min 26s

```

[2.2]Optimal value of parameters after tuning:

```
In [352]: 1 print('Optimal Value of Hyperparameters after Tunning:\n')
          2 print('Alpha: ',params_)
```

Optimal Value of Hyperparameters after Tunning:

Alpha: 0.1

[2.3]Test performance :

```
In [353]: 1 %%time
          2 model_summ_local,clf=test_performance_xgb(train, y_train, test, y_test,\
          3                                           params_, model_summary, model_name[1],\
          4                                           summary=True,regressor_used='linear')
          5
```

FOR OPTIMAL PARAMETERS, TRAIN MAPE: 8.68103, TEST MAPE: 9.01357

CPU times: user 336 ms, sys: 396 ms, total: 732 ms

Wall time: 193 ms

[2.4] Model Summary:

```
In [354]: 1 #MODEL SUMMARY
          2 print(model_summ_local)
```

Model	Train(MAPE)	Test(MAPE)
LINEAR-REGRESSOR	8.68103	9.01357

=====

◀ ▶

3.RandomForest Regressor:

[3.1]Hyperparameter Tunning:

In [355]:

```

1  %%time
2  #TRAIN AND TEST DATA FOR XGBOOST MODELS
3  train, test = x_train, x_test
4  y_train, y_test = tsne_train_output, tsne_test_output
5  print('HYPERPARAMETER:\n')
6  param_list(params_range_rf, param_summ)
7  print()
8  model, params_ = tuneALL_PARAM_XGB(train, y_train, TBS, params_range_rf, params_rf, searchMethod[0], 'rf')

```

HYPERPARAMETER:

Parameter	Value
n_estimators	[20, 40, 80, 128, 256, 512]
max_depth	[9, 12, 15, 18, 20, 25, 27]

Tunning N_ESTIMATORS:

```

Train MAPE: 5.1173 CV MAPE: 5.5600
Train MAPE: 5.1687 CV MAPE: 5.5899
Train MAPE: 5.2574 CV MAPE: 5.7211
Train MAPE: 5.1531 CV MAPE: 5.5405
Train MAPE: 5.1421 CV MAPE: 5.6154
Train MAPE: 5.1145 CV MAPE: 5.5298

```

Tunning MAX_DEPTH:

```

Train MAPE: 0.9335 CV MAPE: 2.8518
Train MAPE: 2.7170 CV MAPE: 3.7283
Train MAPE: 4.4137 CV MAPE: 4.9601
Train MAPE: 1.1338 CV MAPE: 2.9534
Train MAPE: 1.0073 CV MAPE: 2.8946
Train MAPE: 1.6384 CV MAPE: 3.1028

```

CPU times: user 4min 32s, sys: 1.82 s, total: 4min 34s

Wall time: 6min 49s

[3.2]Optimal value of parameters after tuning:

```
In [356]: 1 print('Optimal Value of Hyperparameters after Tunning:\n')
          2 param_list(params_,param_summ)
```

Optimal Value of Hyperparameters after Tunning:

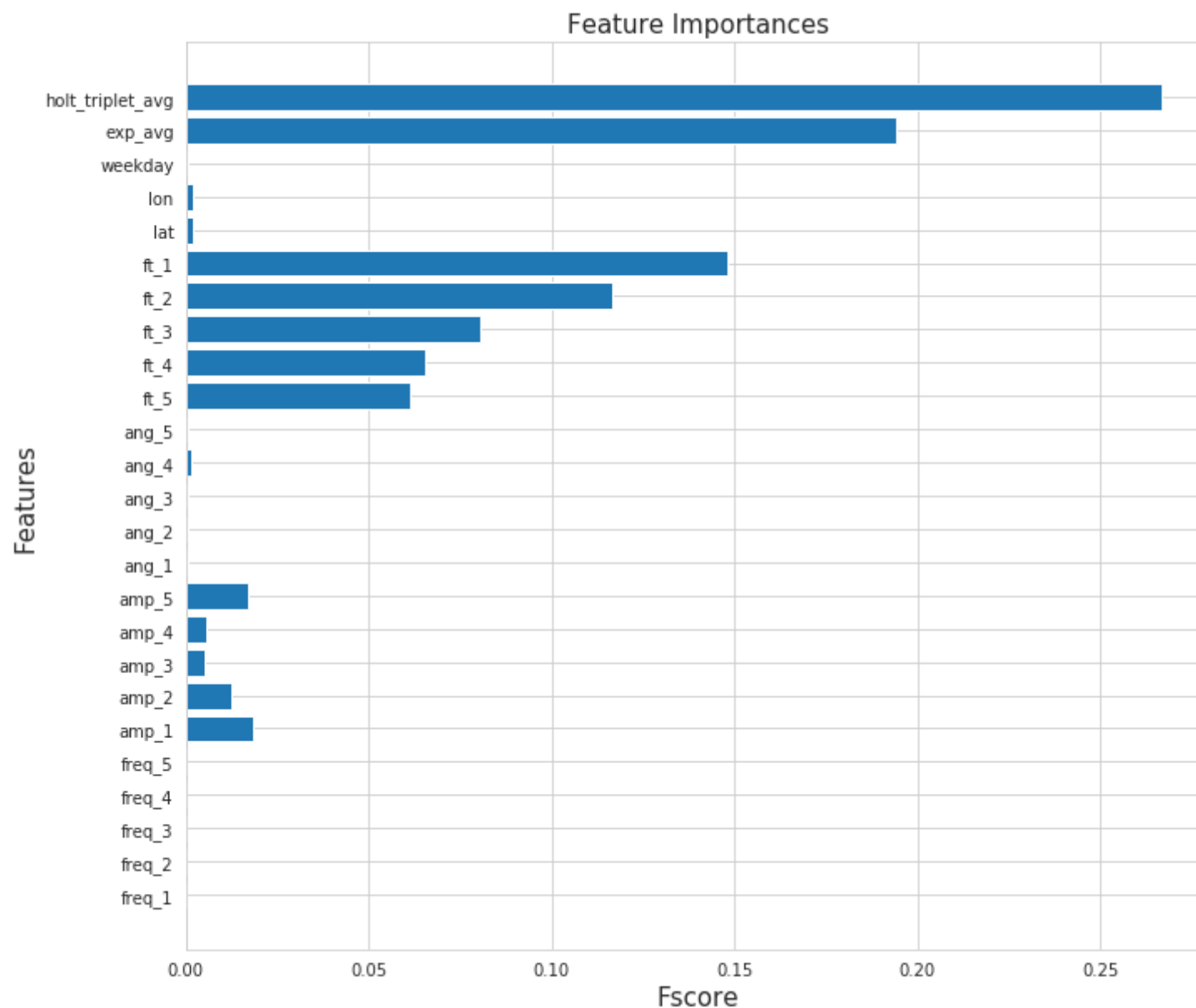
Parameter	Value
n_estimators	512
max_depth	27
min_samples_leaf	1

[3.3]Test performance :

In [357]:

```
1 %%time
2 model_summ_local,clf=test_performance_xgb(train, y_train, test, y_test,\
3                                           params_, model_summary, model_name[2],\
4                                           summary=True,regressor_used='rf')
5
```

FOR OPTIMAL PARAMETERS, TRAIN MAPE: 0.77336, TEST MAPE: 2.30103



CPU times: user 2min 12s, sys: 2.8 s, total: 2min 15s

Wall time: 2min 14s

[3.4] Model Summary:

In [358]:

```
1 #MODEL SUMMARY
2 print(model_summ_local)
```

```
+-----+-----+-----+
|      Model      | Train(MAPE) | Test(MAPE) |
+-----+-----+-----+
| RANDOM-FOREST-REGRESSOR | 0.77336 | 2.30103 |
+-----+-----+-----+
```

=====

◀

4. XGBOOST:

[4.1]Hyperparameter Tunning:

In [359]:

```

1 %%time
2 #TRAIN AND TEST DATA FOR XGBOOST MODELS
3 train, test = x_train, x_test
4 y_train, y_test = tsne_train_output, tsne_test_output
5 print('HYPERPARAMETER:\n')
6 param_list(params_range_xgb, param_summ)
7 print()
8 model, params_ = tuneALL_PARAM_XGB(train, y_train, TBS, params_range_xgb, params_xgb, searchMethod[0], 'xgb')

```

HYPERPARAMETER:

Parameter	Value
n_estimators	[128, 256, 512, 650]
max_depth	[5, 7, 9]
min_child_weight	[1, 3, 5]
gamma	[0.0, 0.1, 0.2, 0.3, 0.4]
subsample	[0.6, 0.7, 0.8, 0.9]
colsample_bytree	[0.6, 0.7, 0.8, 0.9]
reg_alpha	[0, 0.001, 0.005, 0.01, 0.05]

Tunning N_ESTIMATORS:

Train MAPE: 1.0634 CV MAPE: 1.2637
 Train MAPE: 0.9383 CV MAPE: 1.1785
 Train MAPE: 0.8361 CV MAPE: 1.1628
 Train MAPE: 0.7950 CV MAPE: 1.1622

Tunning MAX_DEPTH:

Train MAPE: 0.7950 CV MAPE: 1.1622
 Train MAPE: 0.5476 CV MAPE: 1.2088
 Train MAPE: 0.2923 CV MAPE: 1.2275

Tunning MIN_CHILD_WEIGHT:

Train MAPE: 0.7950 CV MAPE: 1.1622
 Train MAPE: 0.8068 CV MAPE: 1.1519
 Train MAPE: 0.8230 CV MAPE: 1.1513

Tunning GAMMA:

Train MAPE: 0.8230 CV MAPE: 1.1513

Train MAPE: 0.8215 CV MAPE: 1.1541
Train MAPE: 0.8212 CV MAPE: 1.1532
Train MAPE: 0.8212 CV MAPE: 1.1525
Train MAPE: 0.8225 CV MAPE: 1.1488

Tunning SUBSAMPLE:

Train MAPE: 0.8510 CV MAPE: 1.1552
Train MAPE: 0.8341 CV MAPE: 1.1542
Train MAPE: 0.8225 CV MAPE: 1.1488
Train MAPE: 0.8170 CV MAPE: 1.1575

Tunning COLSAMPLE_BYTREE:

Train MAPE: 0.8420 CV MAPE: 1.1772
Train MAPE: 0.8279 CV MAPE: 1.1651
Train MAPE: 0.8225 CV MAPE: 1.1488
Train MAPE: 0.8208 CV MAPE: 1.1560

Tunning REG_ALPHA:

Train MAPE: 0.8226 CV MAPE: 1.1504
Train MAPE: 0.8212 CV MAPE: 1.1517
Train MAPE: 0.8226 CV MAPE: 1.1488
Train MAPE: 0.8206 CV MAPE: 1.1531
Train MAPE: 0.8213 CV MAPE: 1.1547

CPU times: user 13min 4s, sys: 4.42 s, total: 13min 8s

Wall time: 21min 42s

[4.2]Optimal value of parameters after tunning:

```
In [360]: 1 print('Optimal Value of Hyperparameters after Tunning:\n')
          2 param_list(params_,param_summ)
          3
```

Optimal Value of Hyperparameters after Tunning:

Parameter	Value
n_estimators	650
max_depth	5
min_child_weight	5
gamma	0.4
subsample	0.8
colsample_bytree	0.8
reg_alpha	0.005

[4.3]Test performance :

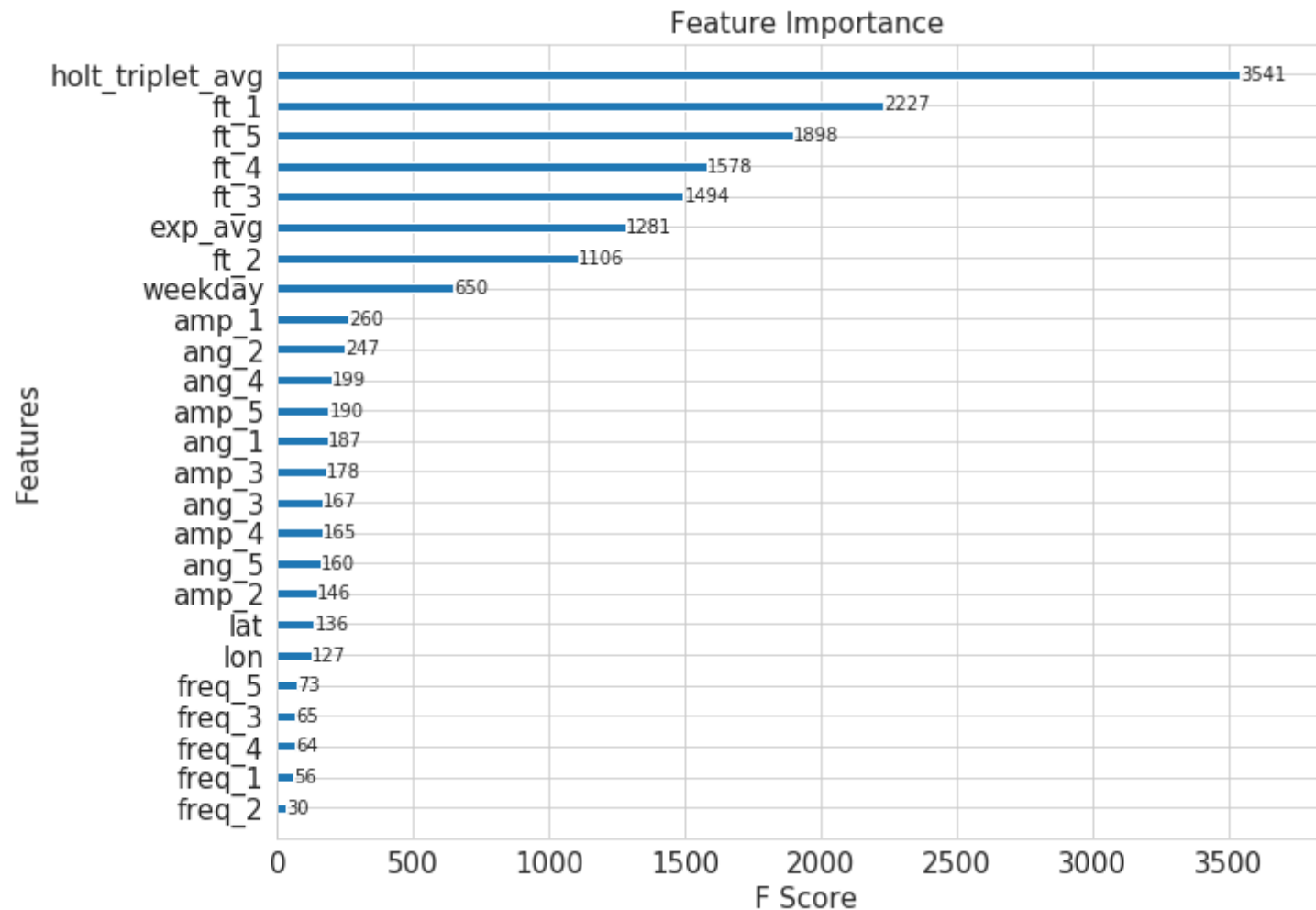
In [361]:

```

1 %%time
2 model_summ_local=test_performance_xgb(train, y_train, test, y_test,\
3                                     params_, model_summary, model_name[3],\
4                                     summary=True,regressor_used='xgb')

```

FOR OPTIMAL PARAMETERS, TRAIN MAPE: 0.87329, TEST MAPE: 1.13973



CPU times: user 1min 16s, sys: 1.89 s, total: 1min 18s

Wall time: 1min 15s

[4.4] Model Summary:

In [362]:

```
1 #MODEL SUMMARY
2 print(model_summ_local)
```

```
+-----+-----+-----+
|      Model      | Train(MAPE) | Test(MAPE) |
+-----+-----+-----+
| XGBOOST-REGRESSOR |    0.87329   |    1.13973   |
+-----+-----+-----+
```

**Conclusion:**

In [363]:

```
1 print(model_summary)
```

```
+-----+-----+-----+
|      Model      | Train(MAPE) | Test(MAPE) |
+-----+-----+-----+
|  HOLTS-WINTER_MODEL  |    3.2515   |    3.3127   |
|  EWMA-PREVIOUS-DATA  |   12.1853   |   12.0552   |
|   LINEAR-REGRESSOR   |    8.68103   |    9.01357   |
| RANDOM-FOREST-REGRESSOR |    0.77336   |    2.30103   |
|   XGBOOST-REGRESSOR   |    0.87329   |    1.13973   |
+-----+-----+-----+
```

Got best performance with model:

2. XGBOOST:

a. Train MAPE : 0.873

b. Test MAPE : 1.139

Procedure:

1. We have to solve the problem, where we have to predict no. of pickups at a given location of New York City in 10 minute interval.
2. We can pose this problem as:
 - a. Time-Series Forecasting(using past data predict future)
 - b. Regression
3. To proceed this problem we are using an approach where we divide the whole New York city into regions/area. So that we can predict the No. of pickups in that area and that area should be that much large only that a taxi/cab can move to that area in 10 minute of interval.
4. We used data:
 - a. Jan 2015 as training data
 - b. Jan 2016 as test data

Note: for Baseline models we use past data of 2016.
5. We proceed with loading the dataset and with basic details of the dataset:
 - a. How many datapoints present in dataset?
 - b. How many features present in dataset?
6. As we have raw data so we did data analysis and data cleaning:
 - a. Removal of outliers in dataset:
 - > Coordinates (longitude and latitude) lies outside of New York.
 - > According to guidelines trip taking more than 12 hrs are not allowed.
 - > Maximum cost of trip is 1000 etc.
7. Further analysis we divided New York city into regions and Whole January month into 10 minute bins(total bins=)
8. From Step 7 for every region, we got region-coordinates(centroid(lat., lon.) or label) and 10 minute time bin(index of 10 mint interval).
9. Training data(i.e. Jan 2015 data) for every data point we attach region label and time bin(10mint time bin).
10. Now we grouped training data based on region label and time bin to find no. of pickups in a region at particular time bin.(ex. reg=1, timebin=22, #pickups=95 and reg=1, timebin=25, #pickups=86)
11. But there is a problem, some of the timebins in a region have zero pickups and if we predict zero pickups for a cab it doesn't make any sense.
(reg=1 at timebin=89, #pickups=0 giving a information that zero pickups are there is not of any use for cab driver)
12. So to solve this problem we have 2-methods:

- a. fill zero: fill zero-pickups for a timebin not present for a cluster (use this approach for test data)
- b. fill smoothing : we fill with average of pickups from neighboring timebin(use this for training data bcz in this we look at future data)

13. We build Baseline models by using previous data or ratios features:

- a. -> SMA-Ratios(simple moving average)
 - > SMA-Predictions
- b. -> WMA-Ratios (weighted moving average)
 - > WMA-Predictions
- c. -> EWMA-Ratios(exponential weighted moving average)
 - > EWMA-Predictions

14. We prepare data for regression models:

- a. We used no. of pickups happened in previous 5 timebins as 5 new features.
- b. EWMA-Predictions model output as a feature for our regression model.
- c. HOLTS WINTER model o/p as a new feature for our regression model.
- d. Cluster/Region centroid latitude and longitude as a feature.
- e. day of week of pickup.

15. By plotting TimeSeries Pickup data we observe that our TimeSeries data have repetitive nature, Whenever a wave have repetitive nature we use Fourier transform to build new features.

16. We applied DFT(using FFT algorithm) and created new features:

- a. top 5 peaks present in Digital signal.
- b. frequencies corresponds to that peaks.
- c. angle corresponds to that peaks.

17. We build Regression models on the prepared data(with 24 feature):

- a. Linear Regression
- b. Random Forest Regression
- c. XGBOOST(GBDT)

Reference Links:

1. <https://www.analyticsvidhya.com/blog/2018/02/time-series-forecasting-methods/> (<https://www.analyticsvidhya.com/blog/2018/02/time-series-forecasting-methods/>)
2. <https://grisha.org/blog/2016/01/29/triple-exponential-smoothing-forecasting/> (<https://grisha.org/blog/2016/01/29/triple-exponential-smoothing-forecasting/>)
3. <https://www.appliedaicourse.com/> (<https://www.appliedaicourse.com/>)
4. <http://snowball.millersville.edu/~adecaria/ESCI386P/esci386-lesson17-Fourier-Transforms.pdf>
(<http://snowball.millersville.edu/~adecaria/ESCI386P/esci386-lesson17-Fourier-Transforms.pdf>)
5. <http://ataspinar.com/2018/04/04/machine-learning-with-signal-processing-techniques/> (<http://ataspinar.com/2018/04/04/machine-learning-with-signal-processing-techniques/>)