



# Politecnico di Torino

**Politecnico di Torino**  
**ICT for smart mobility**  
**Laboratory Report 2**

**Group 3**

s306627 Canineo Komar Vitor  
s306667 Gamballi Lucca  
s308673 Seyedeh Sarah Ghiyasi

Melia Marco  
Vassio Luca

Academic Year 2023-2024

# 1 Prediction using ARIMA models

do not use previous and  
impersonal form

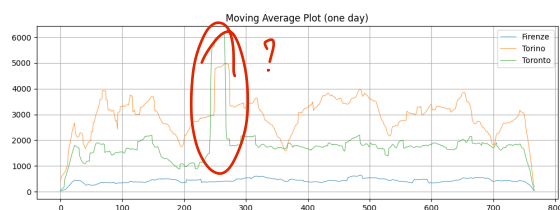
## 1.1 Extracted Period

In order to choose the month to be analyzed it was important to take into account the seasonality of certain phenomena such as holidays, for instance, December would not be a good candidate to be analyzed due to the Christmas and New Year holidays, for similar reasons August would be a bad candidate, in both cases the process would most likely be nonstationary, making the analysis less accurate. Therefore October was chosen as the month to be analyzed, because there aren't any major events during this month, and seems pretty regular in behavior.

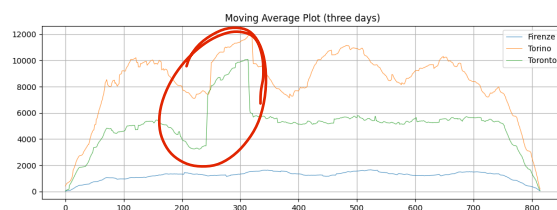
## 1.2 Missing Values and Filtering

There are currently 736 data points available out of a total of 736. Consequently, there are some missing values. Given that the ARIMA model cannot handle missing data, it is essential to impute these missing values. One effective technique involves propagating the values from the last available sample. Since the number of missing points is minimal, this approach should suffice. Also, it is important to point out that the data was filtered using the same process as in the previous lab. bookings with less than two minutes and more than 3 hours were removed as well as bookings that had the same initial and final address.

## 1.3 Stationarity



(a) Moving average using one day window



(b) Moving average using three days window

Figure 1: Moving average for different window size

To check stationarity a moving window approach was used, this can be seen on figure1, both graphs don't show any visible trend, therefore the process can be considered stationary.

but they show clear outliers! why?

## 1.4 Autocorrelation and Partial Autocorrelation

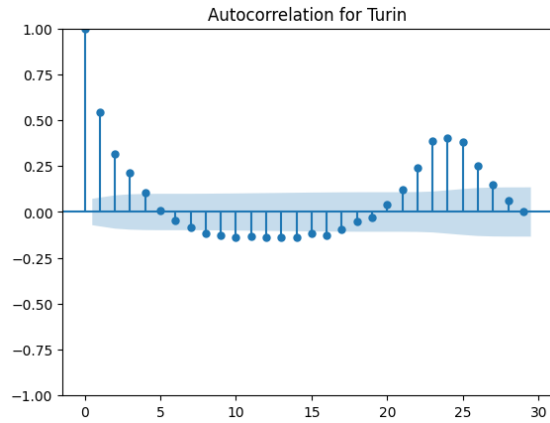
The Autocorrelation and Partial Autocorrelation functions can be seen on figures 2,3 and 4. In Turin, the Partial AutoCorrelation Function (PACF) exhibits a pronounced drop, reaching zero rapidly after  $q=1$ . Conversely, the AutoCorrelation Function (ACF) demonstrates a gradual decline, converging to zero around  $p=5$ .

The patterns observed in Florence closely mirror those in Turin. In Florence, the PACF experiences a significant decrease, reaching zero swiftly after  $q=1$ , while the ACF gradually diminishes, attaining zero around  $p=5$ . This parallel behavior suggests similarities in habits between the two Italian cities.

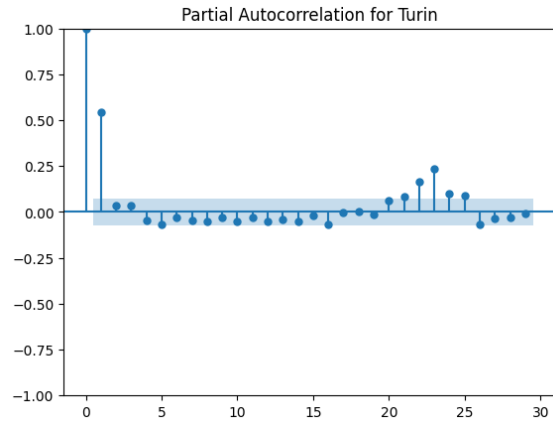
Toronto, however, deviates from the observed patterns in the Italian cities. In Toronto, the ACF follows a different trajectory, declining gently and reaching zero after  $p=2$ . This distinctive behavior in Toronto's ACF indicates dissimilarities in the time series dynamics compared to the Italian cities.

## 1.5 Training and Test Datasets

It was decided to use the last week of the month as test dataset, consequently the first three weeks were used as training. Considering the 31 days of October and (each with 24 hours), this division resulted in a

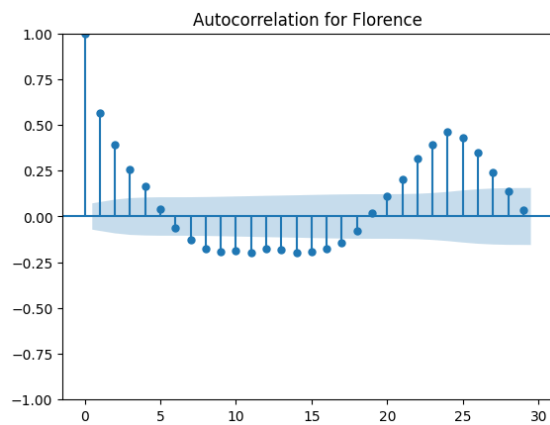


(a) Autocorrelation for Turin

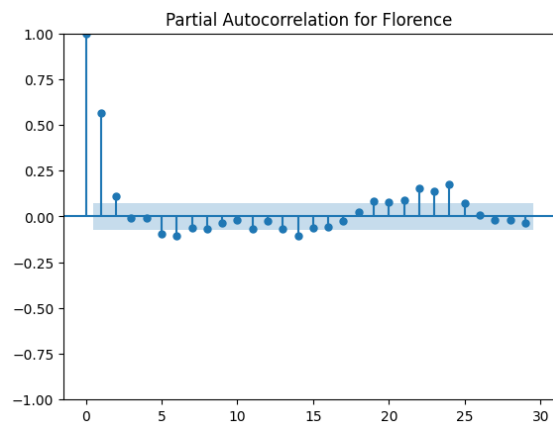


(b) Partial Autocorrelation for Turin

Figure 2: Autocorrelation and Partial Autocorrelation for Turin



(a) Autocorrelation for Florence



(b) Partial Autocorrelation for Florence

Figure 3: Autocorrelation and Partial Autocorrelation for Florence

training dataset with 504 samples and a test dataset with 240 samples. Which means that the train and test sets represent approximately, 70% and 30% of the data, respectively. This division of the data was adopted in order to use the entirety of available data, additionally it allows to keep a lengthy enough training dataset while not compromising the test set size. ✓

## 1.6 Model training

Considering the division of the dataset mentioned in Section 1.5, and the fact that the dataset can be considered a stationary process, an ARIMA model of order (2,0,2) was trained for the cities Torino, Firenze and Toronto. The results obtained on the test set can be seen in Table 1. Additionally, Figures 5, 6 and 7 represent, respectively, the comparison between the ground truth and the predicted value. In those figures the ground truth is shown in blue, whereas the predicted values are shown in red, as it is shown the results were very good, with the predicted values closely following the ground truth values. It is also interesting to point out that the predicted values seem shifted to the right in relation to the ground truth, for instance when observing large peaks, this can be explained by the model being influenced by last values, the model assumed that the future will be relatively similar to the past when large values occur this does not hold.

→ not really ... the <sup>2</sup> ena is sizeable  
compare with a baseline predictor (?)

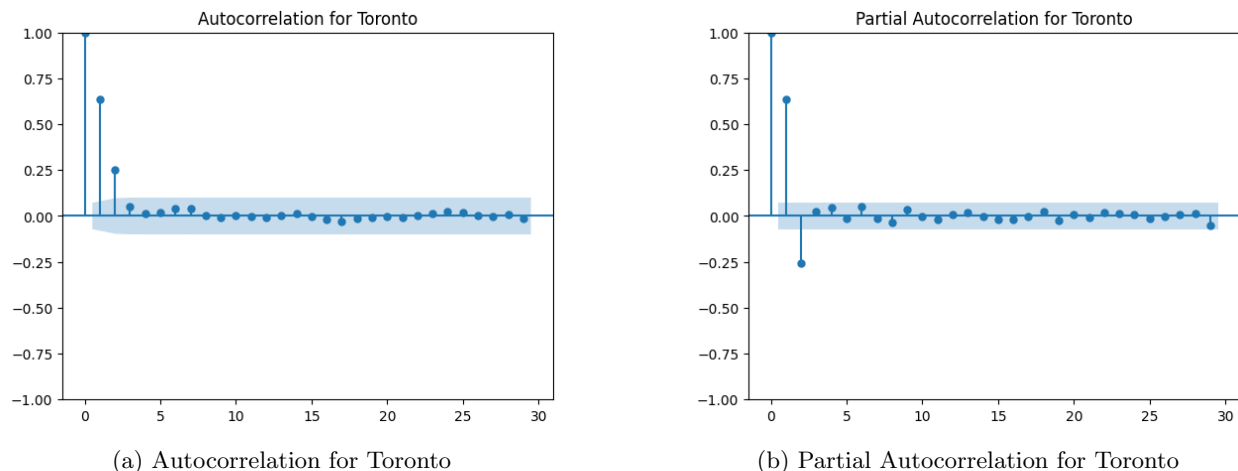


Figure 4: Autocorrelation and Partial Autocorrelation for Toronto

*absolute error?  
or percentage error?*

	Torino	Firenze	Toronto
MAPE	0.8799	0.5878	0.9005

Table 1: Mean absolute percentage error for each city.

*0,8% (?) are you  
sure? This look  
very low.*

## 1.7 Grid Search

### 1.7.1 Keep the number of training samples fixed, and do a grid search varying $(p, d, q)$ and observe how the error changes.

Considering the same size of training and test sets discussed in Section 1.5, values of  $p, q \in [2, 6]$  were considered to perform the Grid Search. Different values of  $d$  were not considered due to the fact that the process is stationary, therefore the  $d$  can be set to zero. Figure 8 shows the mean absolute percentage error (MAPE) obtained by each model trained and for each city, among Torino, Firenze and Toronto. In Figure 8, models of order that resulted in lower-upper (LU) decomposition errors are represented by  $-1$ .

*$d=0$*

Analysing Figure 8 and considering Firenze, it can be seen that there are a couple of equivalent models. Those are the models of orders  $(2, 0, 3)$ ,  $(2, 0, 4)$ ,  $(3, 0, 3)$  and  $(3, 0, 4)$ , all of which resulted in a MAPE value of approximately 0.499. Therefore the best option is the model of order  $(2, 0, 3)$  because it produces the same performance metric with less complexity.

The same reasoning can be applied to the city of Toronto. In Figure 8 it can be seen that even more models can be considered equivalent order. Those are the models of orders  $(2, 0, 2)$ ,  $(2, 0, 3)$ ,  $(2, 0, 6)$ ,  $(3, 0, 3)$ ,

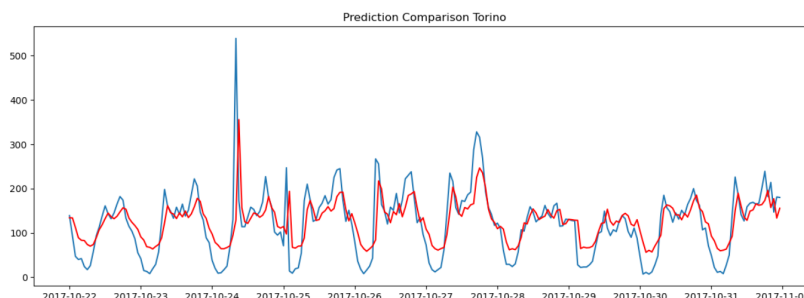


Figure 5: Torino prediction comparison.

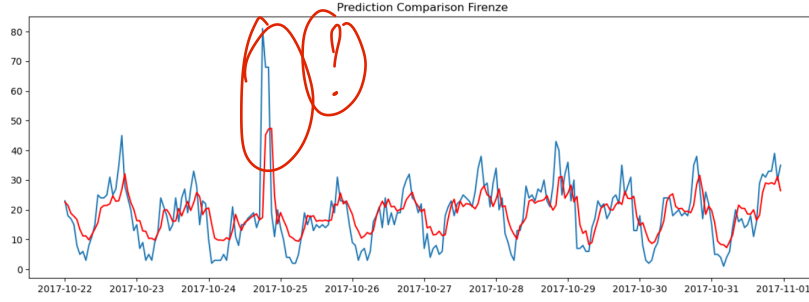


Figure 6: Firenze prediction comparison.

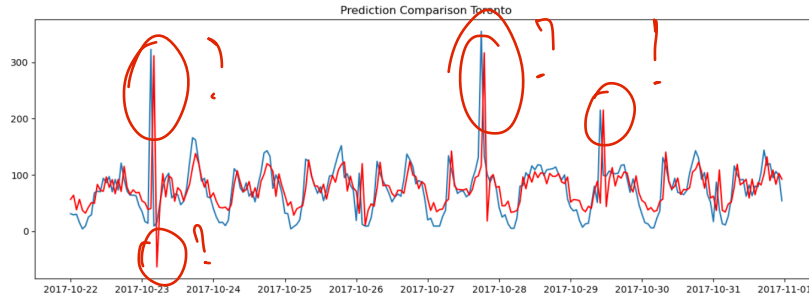


Figure 7: Toronto prediction comparison.

(4,0,3) and (4,0,4) all of them resulted in a MAPE value of approximately 0.90. Therefore, applying the logic used to decide the best model for Firenze, the best option for Toronto is the model of order (2,0,2).

Finally, for Torino there is one model that outstands the others, which is the one of order (3,0,5). This model produced a MAPE of approximately 0.62, whereas all other models are less optimal. The closest one stands at a MAPE value of 0.65, while all others have values greater than 0.67. These are very high! why?

### 1.7.2 Given the best parameter configuration, change N and the learning strategy

Knowing the best parameter configuration for each city, we proceeded to evaluate the effect of changing the size of the training set (represented by  $N$ ) and the learning strategy. Two learning strategies were considered, sliding window and expanding window. For what considers  $N$ , the array  $\mathbf{N} = [11 \ 28 \ \dots \ 504]$  represent the thirty different values adopted. A full representation of  $\mathbf{N}$  can be seen in Equation 1. Figure 9 shows

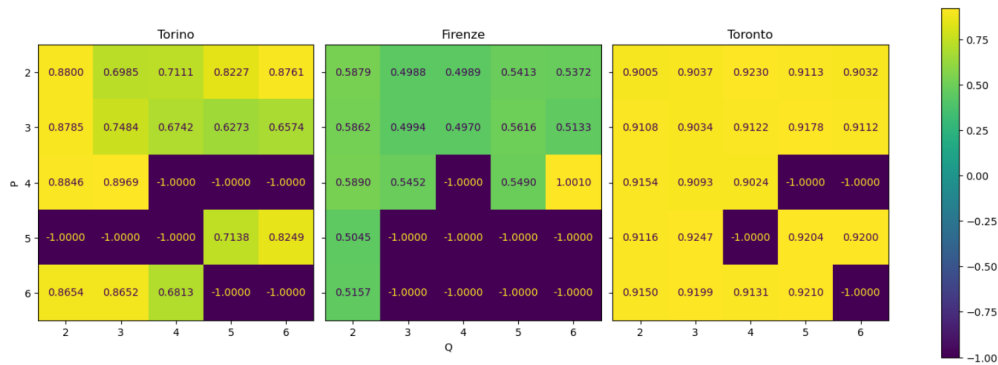


Figure 8: Grid Search results.

the result on the test set for Torino, Firenze and Toronto. It is worth mentioning that when evaluating the performance over all values of  $N$  adopted, the test set was fixed for each city.



Figure 9: Sliding and expanding windows comparison.

### 1.7.3 Compare results for the different cities. How the relative error changes w.r.t. the absolute number of rentals?

Analysing Figure 9a, it can be seen that for almost all values of  $N$  the expanding window learning strategy produces a greater error than the sliding one. The only values of  $N$  for which the sliding is performing worse are the ones that are closer to zero. Additionally, it can be seen that the smaller MAPE values occurs when considering the sliding window and when  $N$  is approximately 50. This result is an indicative that, for Torino, the optimal  $N$  value is in the vicinity of  $N = 50$  and sliding window is the best learning strategy.

Considering Figure 9b, it can be seen that for all values of  $N$  the sliding results in a MAPE value greater or equal than the ones that result from an expanding window strategy. Furthermore, it can be seen that the smaller MAPE values occurs when considering the expanding window and when  $N$  is close to 120. This result is an indicative that, for Firenze, the optimal  $N$  value is in the vicinity of  $N = 120$  and expanding window is the best learning strategy. Additionally, it can be seen that in general the difference between the two strategies is smaller in the case of Firenze than in the case of Torino.

From Figure 9c it can be seen that for values of  $N$  greater than 200 both strategies follow a similar behaviour. For these values of  $N$  the expanding window slightly outperforms the sliding window. For values smaller than 200, it can be seen that the curves do not behave in the same way. Moreover, it can be seen that the smaller MAPE values occurs when considering the expanding window and when  $N$  is, once again, close to 120. This result is an indicative that, for Toronto, the optimal  $N$  value is in the vicinity of  $N = 120$  and expanding window is the best learning strategy. Finally, when analysing the sliding window from Figures 9b and 9c, it can be seen that for certain values of  $N$  there is a peak in the MAPE value. Take  $N = 113$  in Figure 9c as an example. A possible explanation for this behavior is that the model has overfitting issues for those values of  $N$ .

- results are not consistent with the typical errors seen in the class and by other groups. check if there are any problems and what type of error do you see
  - no time horizon comparison
  - no baseline comparison
- RESUBMIT after fixing

## 2 Appendix

### 2.1 Values of N

Equation 1 shows all values adopted for N,

$$\mathbf{N} = \begin{bmatrix} 11 & 28 & 45 & 62 & 79 & 96 & 113 & 130 & 147 & 164 \\ 181 & 198 & 215 & 232 & 249 & 266 & 283 & 300 & 317 & 334 \\ 351 & 368 & 385 & 402 & 419 & 436 & 453 & 470 & 487 & 504 \end{bmatrix}. \quad (1)$$

### 2.2 Code

```
1 """ lab2
2
3     authors :
4         Lucca Gamballi
5         Sarah Ghiyasi Seyedeh
6         Vitor Canineo Komar
7
8     ARIMA models
9 """
10
11
12 import random as rnd
13 from statsmodels.graphics.tsaplots import plot_acf
14 from statsmodels.graphics.tsaplots import plot_pacf
15 from statsmodels.tsa.arima.model import ARIMA
16 import pandas as pd
17 import numpy as np
18 import matplotlib.pyplot as plt
19 from matplotlib import colors
20 import pymongo as pm
21 import pprint
22 from datetime import datetime
23 from sklearn.metrics import mean_squared_error
24 from sklearn.metrics import mean_absolute_error
25 from sklearn.metrics import r2_score
26 from sklearn.metrics import mean_absolute_percentage_error
27 from sklearn.metrics import ConfusionMatrixDisplay
28 import seaborn as sn
29
30 get_ipython().run_line_magic('matplotlib', 'inline')
31
32 client = pm.MongoClient('bigdatadb.polito.it',
33                         ssl=True,
34                         authSource = 'carsharing',
35                         username = 'ictts',
36                         password = 'Ict4SM22!',
37                         tlsAllowInvalidCertificates = True)
38 db = client['carsharing']
39
40 PermanentBookings = db['PermanentBookings']
41 ActiveBookings = db['ActiveBookings']
42
43 PermanentParkings = db['PermanentParkings']
44 ActiveParkings = db['ActiveParkings']
45
46 enjoy_PermanentBookings = db['enjoy_PermanentBookings']
47 enjoy_ActiveBookings = db['enjoy_ActiveBookings']
48
49 enjoy_PermanentParkings = db['enjoy_PermanentParkings']
50 enjoy_ActiveParkings = db['enjoy_ActiveParkings']
51
```

do not  
compare  
strings  
else locations

```

52 def make_hour_df_filter(collection, city):
53     aggregation = collection.aggregate([
54         {'$match' : {'city': city,
55                     'init_address' : {
56                         '$final_address' : {'$ne': '
57                         # 'init_address' not equal to 'final_address'
58                         }
59                     }},
60         {'$project' : {'duration': {'$subtract': [
61             '$final_time', '$init_time']},
62             'hour': {'$hour': '$init_date'},
63             'year': {'$year': '$init_date'},
64             'timestamp': {'$dateToString': {'format
65                 '":"%Y-%m-%d-%H"', "date": "$init_date"}}}}
66         },
67         {'$match': {
68             'duration': {
69                 '$gte': 60, #
70                 '$lte': 3600 * 3 #
71             }
72         },
73         {'$group' : {
74             '_id' : '$timestamp',
75             'count' : {'$sum': 1},
76             'avg' : {'$avg': 1}
77         }
78     })
79
80
81     hours = []
82     bookings = []
83     for el in aggregation:
84         hours.append(el['_id'])
85         bookings.append(el['count'])
86     df = pd.DataFrame({'hours': hours, 'count': bookings})
87     df['hours'] = pd.to_datetime(df['hours'], format="%Y-%m-%d-%H")
88     df = df.sort_values('hours')
89     mask = (df['hours'].dt.year == 2017) & (df['hours'].dt.month == 10)
90     df = df[mask]
91     return df
92
93 hourDfFirenze = make_hour_df_filter(enjoyPermanentBookings, 'Firenze')
94 hourDfTorino = make_hour_df_filter(PermanentBookings, 'Torino')
95 hourDfToronto = make_hour_df_filter(PermanentBookings, 'Toronto')
96
97
98 # ### 2 Checking for missing values
99 len(hourDfTorino)
100 # There are missing values
101
102 def fill_missing_values(df):
103     df = (df.set_index('hours')
104           .reindex(pd.date_range('2017-10-01', '2017-11-1', freq='h'))
105           .rename_axis(['hours'])).fillna(method='ffill').iloc[: -1]
106     return df
107
108 hourDfFirenze = fill_missing_values(hourDfFirenze)
109 hourDfTorino = fill_missing_values(hourDfTorino)
110 hourDfToronto = fill_missing_values(hourDfToronto)
111

```



```

112 plt.figure(figsize=(12, 4))
113 plt.plot(hourDfFirenze.index, hourDfFirenze['count'], label='Firenze', linewidth=0.6)
114 plt.plot(hourDfTorino.index, hourDfTorino['count'], label='Torino', linewidth=0.6)
115 plt.plot(hourDfToronto.index, hourDfToronto['count'], label='Toronto', linewidth=0.6)
116 plt.title('Bookings per Hour (Filtered)')
117 plt.tight_layout()
118 plt.legend()
119 plt.grid()
120 plt.show()
121
122 plt.figure(figsize=(12, 4))
123 plt.plot(np.convolve(hourDfFirenze['count'], np.ones(6)), label='Firenze', linewidth=0.6)
124 plt.plot(np.convolve(hourDfTorino['count'], np.ones(6)), label='Torino', linewidth=0.6)
125 plt.plot(np.convolve(hourDfToronto['count'], np.ones(6)), label='Toronto', linewidth=0.6)
126 plt.legend()
127 plt.title('Moving Average Plot (six hours)')
128 plt.grid()
129 plt.show()
130
131 plt.figure(figsize=(12, 4))
132 plt.plot(np.convolve(hourDfFirenze['count'], np.ones(24)), label='Firenze', linewidth=0.6)
133 plt.plot(np.convolve(hourDfTorino['count'], np.ones(24)), label='Torino', linewidth=0.6)
134 plt.plot(np.convolve(hourDfToronto['count'], np.ones(24)), label='Toronto', linewidth=0.6)
135 plt.legend()
136 plt.title('Moving Average Plot (one day)')
137 plt.grid()
138 plt.show()
139
140 plt.figure(figsize=(12, 4))
141 plt.plot(np.convolve(hourDfFirenze['count'], np.ones(24*3)), label='Firenze', linewidth=0.6)
142 plt.plot(np.convolve(hourDfTorino['count'], np.ones(24*3)), label='Torino', linewidth=0.6)
143 plt.plot(np.convolve(hourDfToronto['count'], np.ones(24*3)), label='Toronto', linewidth=0.6)
144 plt.legend()
145 plt.title('Moving Average Plot (three days)')
146 plt.grid()
147 plt.show()
148
149 plt.figure(figsize=(12, 4))
150 plt.plot(np.convolve(hourDfFirenze['count'], np.ones(24*7)), label='Firenze', linewidth=0.6)
151 plt.plot(np.convolve(hourDfTorino['count'], np.ones(24*7)), label='Torino', linewidth=0.6)
152 plt.plot(np.convolve(hourDfToronto['count'], np.ones(24*7)), label='Toronto', linewidth=0.6)
153 plt.legend()
154 plt.title('Moving Average Plot (seven days)')
155 plt.grid()
156 plt.show()
157
158 plot_acf(hourDfTorino['count'], title='Autocorrelation for Turin')
159
160 plot_pacf(hourDfTorino['count'], title='Partial Autocorrelation for Turin')
161
162 plot_acf(hourDfFirenze['count'], title='Autocorrelation for Florence')
163
164 plot_pacf(hourDfFirenze['count'], title='Partial Autocorrelation for Florence')
165
166 plot_acf(hourDfToronto['count'], title='Autocorrelation for Toronto')
167
168 plot_pacf(hourDfToronto['count'], title='Partial Autocorrelation for Toronto')
169
170
171 p=2; d=0; q=2
172 modelTorino = ARIMA(hourDfTorino['count'], order=(p,d,q))
173 modelTorino = modelTorino.fit()
174
175 modelFirenze = ARIMA(hourDfFirenze['count'], order=(p,d,q))
176 modelFirenze = modelFirenze.fit()

```

what does the convolve method do ?

```

177
178 modelToronto = ARIMA(hourDfToronto['count'], order=(p,d,q))
179 modelToronto = modelToronto.fit()
180
181 fig = plt.figure(figsize=(15,5))
182 plt.plot(hourDfTorino['count'])
183 plt.plot(modelTorino.fittedvalues, color='red')
184 plt.title('Prediction Comparison Torino')
185
186 fig = plt.figure(figsize=(15,5))
187 plt.plot(hourDfFirenze['count'])
188 plt.plot(modelFirenze.fittedvalues, color='red')
189 plt.title('Prediction Comparison Firenze')
190
191 fig = plt.figure(figsize=(15,5))
192 plt.plot(hourDfToronto['count'])
193 plt.plot(modelToronto.fittedvalues, color='red')
194 plt.title('Prediction Comparison Toronto')
195
196 mean_squared_error(hourDfTorino['count'], modelTorino.fittedvalues), mean_absolute_error(
    hourDfTorino['count'], modelTorino.fittedvalues)
197 mean_squared_error(hourDfFirenze['count'], modelFirenze.fittedvalues), mean_absolute_error(
    hourDfFirenze['count'], modelFirenze.fittedvalues)
198 mean_squared_error(hourDfToronto['count'], modelToronto.fittedvalues), mean_absolute_error(
    hourDfToronto['count'], modelToronto.fittedvalues)
199
200 def trainTestSplit(df, split):
201     train = df[0:split]
202     test = df[split:df.shape[0]]
203
204     return train, test
205
206 def walkFowardValidation(train, test, order, windowType='sliding'):
207     history = [x for x in train]
208     predictions = []
209
210     for i, testSample in enumerate(test):
211         model = ARIMA(history, order=order)
212         model_fit = model.fit()
213         output = model_fit.forecast()
214         yhat = output[0]
215         predictions.append(yhat)
216         history.append(testSample)
217         if windowType=='sliding':
218             history.pop(0)
219
220     return predictions
221
222 trainSamples = 3*7*24 #first 3 week of samples
223 N = trainSamples
224
225 trainTorino, testTorino = trainTestSplit(hourDfTorino['count'], trainSamples)
226 trainFirenze, testFirenze = trainTestSplit(hourDfFirenze['count'], trainSamples)
227 trainToronto, testToronto = trainTestSplit(hourDfToronto['count'], trainSamples)
228
229 print(hourDfTorino['count'].shape, hourDfFirenze['count'].shape, hourDfToronto['count'].
    shape)
230 print(trainTorino.shape, trainFirenze.shape, trainToronto.shape)
231 print(testTorino.shape, testFirenze.shape, testToronto.shape)
232
233 p=2; d=0; q=2
234 order = (p,d,q)
235 predsTorino = walkFowardValidation(trainTorino, testTorino, order, windowType='sliding')
236 predsFirenze = walkFowardValidation(trainFirenze, testFirenze, order, windowType='sliding')
237 predsToronto = walkFowardValidation(trainToronto, testToronto, order, windowType='sliding')

```

MAPE (?)

```

238
239 fig = plt.figure(figsize=(15,5))
240 plt.plot(testTorino, label='truth')
241 plt.plot(testTorino.index, predsTorino, color='red', label='prediction')
242 plt.title('Prediction Comparison Torino')
243
244 fig = plt.figure(figsize=(15,5))
245 plt.plot(testFirenze, label='truth')
246 plt.plot(testFirenze.index, predsFirenze, color='red', label='prediction')
247 plt.title('Prediction Comparison Firenze')
248
249 fig = plt.figure(figsize=(15,5))
250 plt.plot(testToronto, label='truth')
251 plt.plot(testToronto.index, predsToronto, color='red', label='prediction')
252 plt.title('Prediction Comparison Toronto')
253
254 print(mean_squared_error(testTorino, predsTorino), mean_absolute_error(testTorino,
    predsTorino), mean_absolute_percentage_error(testTorino, predsTorino))
255 print(mean_squared_error(testFirenze, predsFirenze), mean_absolute_error(testFirenze,
    predsFirenze), mean_absolute_percentage_error(testFirenze, predsFirenze))
256 print(mean_squared_error(testToronto, predsToronto), mean_absolute_error(testToronto,
    predsToronto), mean_absolute_percentage_error(testToronto, predsToronto))
257
258 P = [i for i in range(2,7)]
259 Q = [i for i in range(2,7)]
260 d = 0
261
262 resultsTorino = []
263 resultsFirenze = []
264 resultsToronto = []
265 for p in P:
266     for q in Q:
267         order = (p,d,q)
268         print(order)
269         predsTorino = walkFowardValidation(trainTorino, testTorino, order, windowType='sliding')
270         mse = mean_squared_error(testTorino, predsTorino)
271         mae = mean_absolute_error(testTorino, predsTorino)
272         mape = mean_absolute_percentage_error(testTorino, predsTorino)
273         resultsTorino.append((order, mse, mae, mape))
274
275         predsFirenze = walkFowardValidation(trainFirenze, testFirenze, order, windowType='
    sliding')
276         mse = mean_squared_error(testFirenze, predsFirenze)
277         mae = mean_absolute_error(testFirenze, predsFirenze)
278         mape = mean_absolute_percentage_error(testFirenze, predsFirenze)
279         resultsFirenze.append((order, mse, mae, mape))
280
281         predsToronto = walkFowardValidation(trainToronto, testToronto, order, windowType='
    sliding')
282         mse = mean_squared_error(testToronto, predsToronto)
283         mae = mean_absolute_error(testToronto, predsToronto)
284         mape = mean_absolute_percentage_error(testToronto, predsToronto)
285         resultsToronto.append((order, mse, mae, mape))
286
287 ## considering TORINO
288 ## restart from order 5,0,5 due to LU decomposition error on orders
289 ##     4,0,4 --- 4,0,5 --- 4,0,6
290 ##     5,0,2 --- 5,0,3 --- 5,0,4
291 ## orders 6,0,5 and 6,0,6 also had LU decomposition error
292
293 P = [i for i in range(4,7)]
294 Q = [i for i in range(2,7)]
295 d = 0
296
297 for p in P:

```

```

298     for q in Q:
299         if p == 4 and q < 7:
300             pass
301         elif p == 5 and q < 5:
302             pass
303         else:
304             order = (p,d,q)
305             print(order)
306             predsTorino = walkFowardValidation(trainTorino, testTorino, order, windowType='
sliding')
307             mse = mean_squared_error(testTorino, predsTorino)
308             mae = mean_absolute_error(testTorino, predsTorino)
309             mape = mean_absolute_percentage_error(testTorino, predsTorino)
310             resultsTorino.append((order, mse, mae, mape))
311
312     ## considering FIRENZE
313     ## restart from order 4,0,5 due to LU decomposition error on orders
314     ##     4,0,4
315     P = [i for i in range(4,7)]
316     Q = [i for i in range(2,7)]
317     d = 0
318
319     for p in P:
320         for q in Q:
321             if p == 4 and q < 5:
322                 pass
323             else:
324                 order = (p,d,q)
325                 print(order)
326                 predsFirenze = walkFowardValidation(trainFirenze, testFirenze, order, windowType='
sliding')
327                 mse = mean_squared_error(testFirenze, predsFirenze)
328                 mae = mean_absolute_error(testFirenze, predsFirenze)
329                 mape = mean_absolute_percentage_error(testFirenze, predsFirenze)
330                 resultsFirenze.append((order, mse, mae, mape))
331
332     ## considering FIRENZE
333     ## restart from order 6,0,2 due to LU decomposition error on orders
334     ##     5,0,3 --- 5,0,4 --- 5,0,5 --- 5,0,6
335     ## from 6,0,3 all orders result in LU decomposition error
336     P = [i for i in range(5,7)]
337     Q = [i for i in range(2,7)]
338     d = 0
339
340     for p in P:
341         for q in Q:
342             if p == 5 and q < 7:
343                 pass
344             else:
345                 order = (p,d,q)
346                 print(order)
347                 predsFirenze = walkFowardValidation(trainFirenze, testFirenze, order, windowType='
sliding')
348                 mse = mean_squared_error(testFirenze, predsFirenze)
349                 mae = mean_absolute_error(testFirenze, predsFirenze)
350                 mape = mean_absolute_percentage_error(testFirenze, predsFirenze)
351                 resultsFirenze.append((order, mse, mae, mape))
352
353     ## considering TORONTO
354     ## restart from order 4,0,6 due to LU decomposition error on orders
355     ##     4,0,5 --- 4,0,6
356     ##     5,0,2 --- 5,0,3 --- 5,0,4
357     P = [i for i in range(4,7)]
358     Q = [i for i in range(2,7)]
359     d = 0

```

```

360
361 for p in P:
362     for q in Q:
363         if p == 4 and q < 7:
364             pass
365         elif p == 5 and q < 5:
366             pass
367         else:
368             order = (p,d,q)
369             print(order)
370             predsToronto = walkFowardValidation(trainToronto, testToronto, order, windowType='
sliding')
371             mse = mean_squared_error(testToronto, predsToronto)
372             mae = mean_absolute_error(testToronto, predsToronto)
373             mape = mean_absolute_percentage_error(testToronto, predsToronto)
374             resultsToronto.append((order, mse, mae, mape))
375
376 ordersListsTorino = []
377 for index, element in enumerate(resultsTorino):
378     ordersListsTorino.append(element[0])
379
380 ordersListsFirenze = []
381 for element in resultsFirenze:
382     ordersListsFirenze.append(element[0])
383
384 ordersListsToronto = []
385 for element in resultsToronto:
386     ordersListsToronto.append(element[0])
387
388 metricTorino = np.zeros((5,5))
389 metricFirenze = np.zeros((5,5))
390 metricToronto = np.zeros((5,5))
391 elTorino = 0
392 elFirenze = 0
393 elToronto = 0
394 for p in range(2,7):
395     for q in range(2,7):
396         if (p,d,q) in ordersListsTorino:
397             metricTorino[p-2][q-2] = resultsTorino[elTorino][3]
398             elTorino += 1
399         else:
400             metricTorino[p-2][q-2] = -1
401
402         if (p,d,q) in ordersListsFirenze:
403             metricFirenze[p-2][q-2] = resultsFirenze[elFirenze][3]
404             elFirenze += 1
405         else:
406             metricFirenze[p-2][q-2] = -1
407
408         if (p,d,q) in ordersListsToronto:
409             metricToronto[p-2][q-2] = resultsToronto[elToronto][3]
410             elToronto += 1
411         else:
412             metricToronto[p-2][q-2] = -1
413
414 data = [
415     {"city": "Torino", "results": metricTorino},
416     {"city": "Firenze", "results": metricFirenze},
417     {"city": "Toronto", "results": metricToronto},
418 ]
419
420
421 f, axes = plt.subplots(1, 3, figsize=(15, 5), sharey='row')
422
423 for i, element in enumerate(data):

```

```

424     disp = ConfusionMatrixDisplay(element['results'], display_labels=[2,3,4,5,6])
425     disp.plot(ax=axes[i], values_format=".4f")
426     disp.ax_.set_title(element['city'])
427     disp.im_.colorbar.remove()
428     disp.ax_.set_xlabel('')
429     if i!=0:
430         disp.ax_.set_ylabel('')
431     else:
432         disp.ax_.set_ylabel('P')
433
434 f.text(0.4, 0.05, 'Q', ha='left')
435 plt.subplots_adjust(wspace=0.40, hspace=0.1)
436
437 f.tight_layout()
438 f.colorbar(disp.im_, ax=axes)
439 plt.show()
440
441
442 # from grid search
443 bestOrderTorino = (3,0,5)
444 bestOrderFirenze = (2,0,3)
445 bestOrderToronto = (2,0,2)
446
447 trainSamples = 3*7*24 #first 3 week of samples
448 N = trainSamples
449 trainTorino, fixedTest = trainTestSplit(hourDfTorino['count'], trainSamples) #set fixed test
450     set
451 mseSliding = []
452 mseExpanding = []
453
454 maeSliding = []
455 maeExpanding = []
456
457 mapeSliding = []
458 mapeExpanding = []
459 N = [i for i in range(trainSamples, 1, -trainSamples//(30))] #range(start, stop, step)
460 for n in N:
461     print(n)
462     trainTorino, testTorino = trainTestSplit(hourDfTorino['count'], n)
463     predsTorinoSliding = walkFowardValidation(trainTorino, fixedTest, bestOrderTorino,
464         windowType='sliding')
465     predsTorinoExpanding = walkFowardValidation(trainTorino, fixedTest, bestOrderTorino,
466         windowType='expanding')
467
468     mseSliding.append(mean_squared_error(fixedTest, predsTorinoSliding))
469     mseExpanding.append(mean_squared_error(fixedTest, predsTorinoExpanding))
470
471     maeSliding.append(mean_absolute_error(fixedTest, predsTorinoSliding))
472     maeExpanding.append(mean_absolute_error(fixedTest, predsTorinoExpanding))
473
474     mapeSliding.append(mean_absolute_percentage_error(fixedTest, predsTorinoSliding))
475     mapeExpanding.append(mean_absolute_percentage_error(fixedTest, predsTorinoExpanding))
476
477 fig = plt.figure(figsize=(6,4))
478 plt.plot(N, mseSliding, label='sliding')
479 plt.plot(N, mseExpanding, label='expanding')
480 plt.title('Sliding and expanding windows Comparison Torino')
481 plt.legend()
482 plt.xlabel('N')
483 plt.ylabel('mse')
484 plt.show()
485
486 fig = plt.figure(figsize=(6,4))
487 plt.plot(N, maeSliding, label='sliding')

```

```

486 plt.plot(N, maeExpanding, label='expanding')
487 plt.title('Sliding and expanding windows Comparison Torino')
488 plt.legend()
489 plt.xlabel('N')
490 plt.ylabel('mae')
491 plt.show()
492
493 fig = plt.figure(figsize=(6,4))
494 plt.plot(N, mapeSliding, label='sliding')
495 plt.plot(N, mapeExpanding, label='expanding')
496 plt.title('Sliding and expanding windows Comparison Torino')
497 plt.legend()
498 plt.xlabel('N')
499 plt.ylabel('mape')
500 plt.show()
501
502
503 # In[266]:
504
505
506 trainSamples = 3*7*24 #first 3 week of samples
507 N = trainSamples
508 trainFirenze, fixedTest = trainTestSplit(hourDfFirenze['count'], trainSamples) #set fixed
    test set
509
510 mseSliding = []
511 mseExpanding = []
512
513 maeSliding = []
514 maeExpanding = []
515
516 mapeSliding = []
517 mapeExpanding = []
518 N = [i for i in range(trainSamples, 1, -trainSamples//(30))]
519 for n in N:
520     print(n)
521     trainFirenze, testFirenze = trainTestSplit(hourDfFirenze['count'], n)
522     predsFirenzeSliding = walkFowardValidation(trainFirenze, fixedTest, (2,0,3), windowType='
        sliding')
523     predsFirenzeExpanding = walkFowardValidation(trainFirenze, fixedTest, (2,0,3), windowType='
        expanding')
524
525     mseSliding.append(mean_squared_error(fixedTest, predsFirenzeSliding))
526     mseExpanding.append(mean_squared_error(fixedTest, predsFirenzeExpanding))
527
528     maeSliding.append(mean_absolute_error(fixedTest, predsFirenzeSliding))
529     maeExpanding.append(mean_absolute_error(fixedTest, predsFirenzeExpanding))
530
531     mapeSliding.append(mean_absolute_percentage_error(fixedTest, predsFirenzeSliding))
532     mapeExpanding.append(mean_absolute_percentage_error(fixedTest, predsFirenzeExpanding))
533
534 fig = plt.figure(figsize=(6,4))
535 plt.plot(N, mseSliding, label='sliding')
536 plt.plot(N, mseExpanding, label='expanding')
537 plt.title('Sliding and expanding windows Comparison Firenze')
538 plt.legend()
539 plt.xlabel('N')
540 plt.ylabel('mse')
541 plt.show()
542
543 fig = plt.figure(figsize=(6,4))
544 plt.plot(N, maeSliding, label='sliding')
545 plt.plot(N, maeExpanding, label='expanding')
546 plt.title('Sliding and expanding windows Comparison Firenze')
547 plt.legend()

```

```

548 plt.xlabel('N')
549 plt.ylabel('mae')
550 plt.show()
551
552 fig = plt.figure(figsize=(6,4))
553 plt.plot(N, mapeSliding, label='sliding')
554 plt.plot(N, mapeExpanding, label='expanding')
555 plt.title('Sliding and expanding windows Comparison Firenze')
556 plt.legend()
557 plt.xlabel('N')
558 plt.ylabel('mape')
559 plt.show()
560
561 trainSamples = 3*7*24 #first 3 week of samples
562 N = trainSamples
563 trainToronto, fixedTest = trainTestSplit(hourDfToronto['count'], trainSamples) #set fixed
    test set
564
565 mseSliding = []
566 mseExpanding = []
567
568 maeSliding = []
569 maeExpanding = []
570
571 mapeSliding = []
572 mapeExpanding = []
573 N = [i for i in range(trainSamples, 1, -trainSamples//(30))] #range(start, stop, step)
574 for n in N:
575     trainToronto, testToronto = trainTestSplit(hourDfToronto['count'], n)
576     predsTorontoSliding = walkFowardValidation(trainToronto, fixedTest, bestOrderToronto,
        windowType='sliding')
577     predsTorontoExpanding = walkFowardValidation(trainToronto, fixedTest, bestOrderToronto,
        windowType='expanding')
578
579     mseSliding.append(mean_squared_error(fixedTest, predsTorontoSliding))
580     mseExpanding.append(mean_squared_error(fixedTest, predsTorontoExpanding))
581
582     maeSliding.append(mean_absolute_error(fixedTest, predsTorontoSliding))
583     maeExpanding.append(mean_absolute_error(fixedTest, predsTorontoExpanding))
584
585     mapeSliding.append(mean_absolute_percentage_error(fixedTest, predsTorontoSliding))
586     mapeExpanding.append(mean_absolute_percentage_error(fixedTest, predsTorontoExpanding))
587
588 fig = plt.figure(figsize=(6,4))
589 plt.plot(N, mseSliding, label='sliding')
590 plt.plot(N, mseExpanding, label='expanding')
591 plt.title('Sliding and expanding windows Comparison Toronto')
592 plt.legend()
593 plt.xlabel('N')
594 plt.ylabel('mse')
595 plt.show()
596
597 fig = plt.figure(figsize=(6,4))
598 plt.plot(N, maeSliding, label='sliding')
599 plt.plot(N, maeExpanding, label='expanding')
600 plt.title('Sliding and expanding windows Comparison Toronto')
601 plt.legend()
602 plt.xlabel('N')
603 plt.ylabel('mae')
604 plt.show()
605
606 fig = plt.figure(figsize=(6,4))
607 plt.plot(N, mapeSliding, label='sliding')
608 plt.plot(N, mapeExpanding, label='expanding')
609 plt.title('Sliding and expanding windows Comparison Toronto')

```



```
610 plt.legend()
611 plt.xlabel('N')
612 plt.ylabel('mape')
613 plt.show()
```