



Politecnico di Torino
ICT for Smart Mobility
Laboratory report # 2
Group 11

Berti Marco, s315819
Ravera Stefano, s318946
Scardi Alessia, s317628

Contents

1	Prediction using ARIMA models	2
1.1	Time-series extraction and filtering	2
1.2	Preprocessing	2
1.3	Time-series stationary properties	2
1.4	ACF & PACF	3
1.5	Training the model	4
1.6	Metrics	5
1.7	Impact of the parameters	6
1.7.1	Fixed N , tuning p and q	6
1.7.2	Best parameter configuration, tuning N and learning strategy	6
1.7.3	Results comparison	8
1.8	(Optional) Impact of the time horizon h of the prediction on the performance	8
2	Appendix	10
2.1	Figures	10
2.2	Tables	10
2.3	Code	13

1 Prediction using ARIMA models

ARIMA (AutoRegressive Integrated Moving Average) is a time series forecasting method that combines autoregression, differencing, and moving average components. It is used for predicting future points in a time series dataset.

Laboratory 2 is centered on the implementation of the ARIMA models, especially focusing on the hyperparameters tuning, to fit the preprocessed data from Laboratory 1 and perform predictions.

1.1 Time-series extraction and filtering

In order to use ARIMA models, the dataset should be made of time series. The data returned from LAB 1 are filtered considering a 30 days window, going from the 6-th of September to the 6-th of October 2017, chosen because we expect stationarity in each of the examined cities (Torino, Seattle, Berlin) for the cited time window. The stationarity is highly probable since September/October are not festivity months: it is reasonable to expect a high regularity in car rentals. ✓

1.2 Preprocessing

The available dataset is divided into hours, with the total of 720 points (30 days x 24 hours) in order to count the number of rentals recorded at each hour, and paying attention that the car was really moving. The resulting dataset is made of couples: datetime hour by hour and the count of rentals.

In order to work with float values, a random component is introduced in the number of rentals. how?

ARIMA models assume consistent time intervals between observations, not allowing missing samples: missing values are then defined as those whose count is not present for a specific hour. To address missing instances in the dataset, at first zeros are used, then an average methodology is employed. This involves the replacement of zero values with an average based on the available data from corresponding hours of the same day of the week within the time window of 30 days. The following step of the preprocessing part regards removing outliers from the dataset: a point is considered an outlier based on its distance from the just cited average, and replaced with the average if it is more than 100% or less than 30% of the average (the kept rentals should have a duration greater than 5 minutes and lower than 2 hours). why first insert 0? do you also check if the car moved?

The graph of the filtered system utilization over time is shown in fig. 1, showing the number of rentals per hour for the three cities.

1.3 Time-series stationary properties

Forecasting might be challenging due to the non-deterministic nature of time series data, meaning we cannot predict future occurrences with certainty. However, the task becomes more manageable if the time series is stationary. In such cases, predicting future outcomes

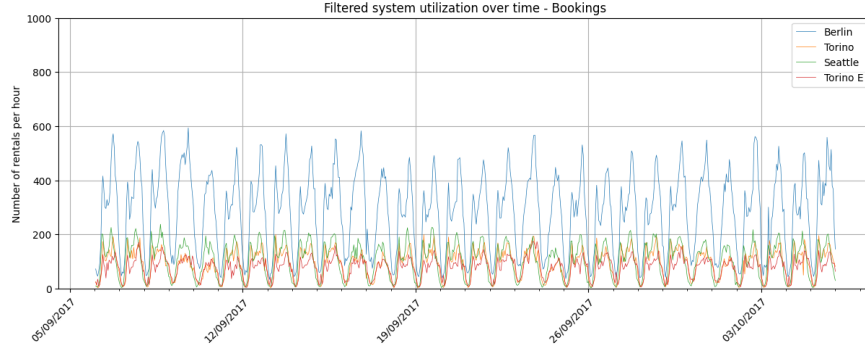


Figure 1: Filtered System Utilization over time for the Bookings.

involves assuming that the statistical properties of the time series will remain consistent over time.

Rolling statistics can be used to evaluate the stationary properties of the time series, which allows the setting of the parameter d - the number of differentiation steps to make the model stationary - of the ARIMA model: if the time series result to be stationary then d can be set to zero since no differentiation is needed.

By plotting the rolling average and the rolling standard deviation on a 7-days time window it can be shown that the time series of the three cities at ease are all stationary: fig. 2 shows the rolling statistics for the time series of Berlin (the other graphs are analogous and present in fig. 9, fig. 10, and fig. 11 in the Appendix section at 2.1). The resulting lines representing the rolling average and standard deviation are horizontal with respect to the time series, meaning that they do not change over time.

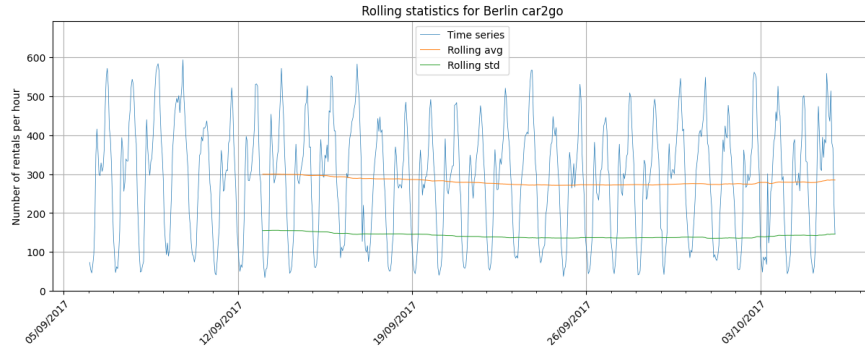


Figure 2: Rolling statistics with a 7-days window for Berlin's data.

1.4 ACF & PACF

Autocorrelation plots are also very useful to determine if a time series is stationary or not because it is simply possible to look at the decay of the time series, that, in the case of non-stationarity, is very slow. However the Autocorrelation Function (ACF) and the Partial Autocorrelation Function (PACF) of the time series can be computed to evaluate the

remaining parameters of the ARIMA model: p - the lag in the AR model - and q - the lag in the MA model. The following observations are made:

- PACF stops after lag p for Autoregressive (AR) models
- ACF stops after lag q for Moving average (MA) models
- ARMA models are the result of the 2 components of the AR and MA models, whose order is denoted by the parameters p and q . The AR component of the model is designed to capture the interdependence between the current value of the time series and its preceding values. Consequently, the hyper-parameter p is intricately linked to the count of hours that exhibit the cited correlation. The MA component of the model is designed to encapsulate the dependence of the current time series value on the noise.

For pure AR
for pure MA

?

The plots for the ACF and the PACF are shown in fig. 3 for Berlin and in the Appendix in fig. 12, fig. 13 and fig. 14 for the other cities. The initial guess for both the parameters p and q is 2 for all of them since the plots are very similar.

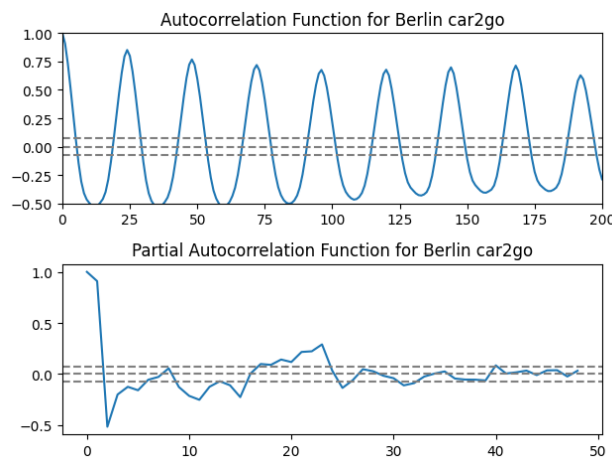


Figure 3: ACF and PACF for Berlin.

1.5 Training the model

The ARIMA model is trained and then used to predict new data. The number of past samples N used for the training set is equal to 360 (15 days), same as the number of samples used to test the model. The initial model employed a sliding window learning approach, involving the segmentation of a large dataset into smaller, overlapping subsets. Each of these subsets is used for training, allowing the model to iteratively update its knowledge as it traversed the dataset. In the testing phase, a distinct model was created for each test sample, enhancing adaptability to evolving data dynamics.

?

1.6 Metrics

Given N , (p, d, q) and a trained model, the error can be computed and evaluated through different metrics such as MPE (Mean Percentage Error) and MAPE (Mean Absolute Percentage Error) for each city.

MAPE is a widely used metric in forecasting and predictive modeling (including ARIMA) to evaluate the accuracy of a model's predictions. It provides a measure of how well a forecasting model is able to predict values in comparison to the actual observed values. MAPE measures the relative percentage difference between predicted and actual values. It is calculated as the average of the absolute percentage differences, as shown in the first expression of eq. 1

MPE is another metric used for evaluating the accuracy of forecasting models and it is a simpler variant of the MAPE, but it does not involve taking the absolute values of the percentage errors. It is evaluated as the second expression shown in eq. 1.

$$\text{MAPE} = \frac{100}{N} \sum_{t=1}^N \left| \frac{X_t - \hat{X}_t}{X_t} \right| \quad \text{MPE} = \frac{100}{N} \sum_{t=1}^N \left(\frac{X_t - \hat{X}_t}{X_t} \right) \quad (1)$$

where N is the number of observations, X_t represent the actual values and \hat{X}_t represent the predicted values.

? MPE is affected by the scale of the data. This is because MPE is based on the raw percentage errors without taking their absolute values. The sign and magnitude of the percentage errors contribute directly to the MPE calculation. As a result, datasets with different scales can lead to different MPE values, and comparing MPE across datasets may not be meaningful.

In contrast, MAPE addresses this issue by taking the absolute values of the percentage errors before averaging them. This ensures that the metric is scale-independent, making it more suitable for comparing forecasting performance across datasets with different scales.

In the case of ARIMA models, MPE provides insights into the average direction - if MPE is positive, it indicates that, on average, the ARIMA model tends to overestimate the actual values. This means that the model predictions are, on average, higher than the observed values - and magnitude - a larger absolute MPE value reflects a higher average percentage error since the magnitude of MPE indicates the average percentage by which the model's predictions deviate from the actual values - of errors, while the lower the MAPE, the better the ARIMA models' performance.

The results for the three cities for MPE and MAPE are shown in table 1. It is clearly shown that Berlin has the lowest MAPE parameter among the three cities, meaning that ARIMA model performs better on Berlin's time series. Also MPEs are negative for all the cities, meaning that the ARIMA model tends to underestimate the actual values, and the predictions are, on average, lower than the observed values. For the city of Berlin the MPE is also lower in magnitude with respect to the other cities, meaning that it is the city for which the predictions deviate the least from the actual values. ✓

1.7 Impact of the parameters

1.7.1 Fixed N , tuning p and q

The first strategy used to evaluate the impact of the parameters on the model is to fix the parameter N and perform a grid search on the parameters p and q . In grid search, a predefined grid of hyperparameter values is specified, and the model is trained and evaluated for each combination of these values. The hyperparameter tuning is performed in order to systematically find the best configuration of parameters that optimize the ARIMA model. The training approach used is the one described in section 1.5.

- N is fixed to 360 (15 days);
- Split training and test 50%/50% (360 hours of train data/360 hours of test data);
- Grid search with ranges $p = q = [0, 1, 2, 3, 4, 5]$ and $d = 0$;
- MPE is taken into account but MAPE is used as metric for evaluating the optimal parameters, since it is a scale-independent metric.

The heatmaps representing the MAPE for all the combination of parameters for Berlin and Torino are shown in fig. 4 while the heatmaps for Seattle and Torino_Enjoy are shown in the Appendix in fig. 15. The best combination of parameters for Berlin is $d = 0$ and $p = q = 2$ while for Torino the best choice is $d = 0$, $p = 2$ and $q = 5$ (same for Seattle and Torino_Enjoy).

1.7.2 Best parameter configuration, tuning N and learning strategy

The best values for the parameters p and q are listed in the previous section but they have been slightly changed for the sake of simplicity: since the parameters p and q determine the complexity of the ARIMA model, lower values of the parameters can be translated into a simpler model.

The best value of q is changed to $q = 3$ for Torino and $q = 2$ for Torino_Enjoy, since the difference in the MAPE when using the best parameters is about a few units, the choice of a simpler, more robust model is favoured to the choice of a more flexible but complex model.

Always keeping the last 10 days (240 test samples) of observations for predictions and the parameters p , q and d constant (Berlin: ARIMA(2,0,2), Torino: ARIMA(2,0,3), Seattle: ARIMA(2,0,5) and Torino_Enjoy: ARIMA(2,0,2)) it is possible to observe the outcome of the model when using different values of N and different learning strategies (sliding and expanding window).

	Berlin	Torino	Seattle	Torino_Enjoy
MAPE	15.947	39.688	34.180	32.272
MPE	-4.823	-24.093	-20.870	-15.254

Table 1: Evaluation metrics for the three cities with parameters $d=0$, $p=2$, $q=2$.

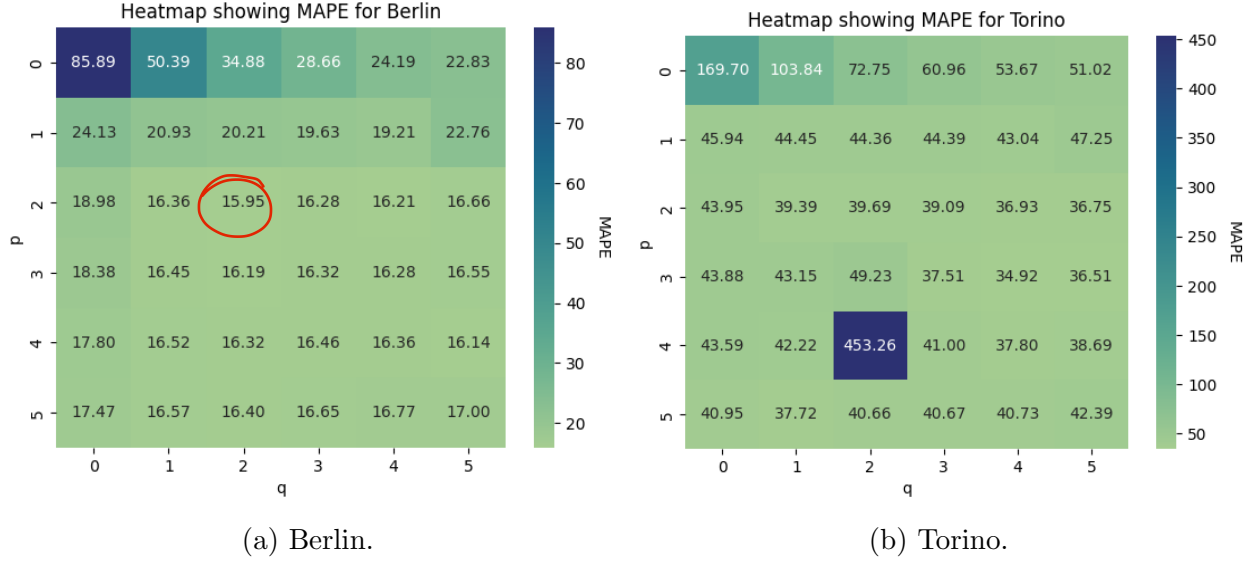


Figure 4: Heatmap showing MAPE.

The sliding window learning strategy has already been mentioned in section 1.5 while the expanding window learning strategy entails training a model on progressively larger data subsets.

Table 2 is referred to the city of Berlin and it shows the metrics MAPE and MPE for different values of N and for both the sliding and the expanding window learning strategy. The tables (ta. 19a, 19b and ta. 3) for the other cities are listed in the Appendix section.

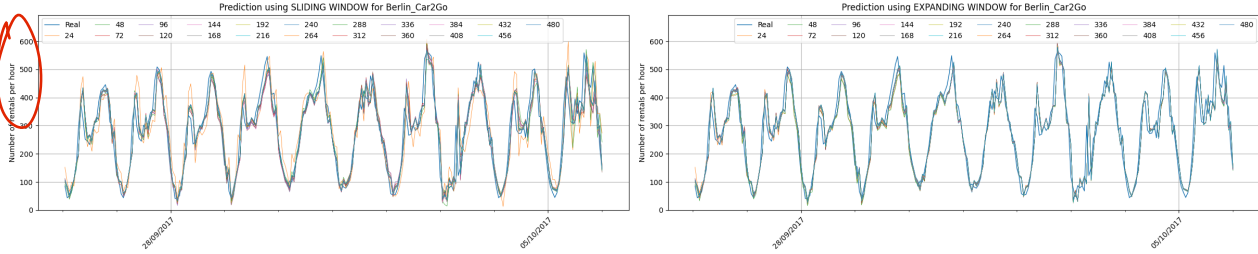
Figures 5a and 5b show respectively the graphical behaviour of the predicted time series when changing N through the sliding and the expanding window approach for Berlin (for other cities: fig. 16, fig. 17 and fig. 18 in Appendix).

	Sliding window		Expanding window	
N	MAPE	MPE	MAPE	MPE
24	27.524	-11.210	18.760	-3.593
48	18.583	-4.962	17.540	-1.135
72	17.820	-3.711	17.376	-3.253
96	18.226	-3.685	17.786	-4.345
120	18.416	-4.226	17.903	-4.342
144	18.319	-4.892	17.573	-4.431
168	18.107	-5.049	17.520	-4.079
192	17.785	-4.919	17.428	-3.647
216	17.452	-4.458	17.123	-3.100
240	17.359	-4.045	16.846	-3.757
264	17.542	-4.019	17.024	-4.570
288	17.634	-4.207	17.118	-4.936
312	17.514	-4.458	17.286	-5.282
336	17.550	-4.636	17.376	-5.446
360	17.456	-4.841	17.398	-5.398
384	17.284	-4.910	17.184	-5.038
408	17.069	-4.924	17.325	-5.436
432	17.028	-4.969	17.602	-5.884
456	17.177	-5.137	17.715	-6.142
480	17.348	-5.413	17.838	-6.345

better plot these results

Table 2: Parameter N tuning for Berlin with $d = 0$ and $p = q = 2$.

too small



(a) Predictions of the number of rentals per hour while changing N , using the sliding window approach, for Berlin. (b) Predictions of the number of rentals per hour while changing N , using the expanding window approach, for Berlin.

Figure 5: Berlin.

1.7.3 Results comparison

The influence of N on prediction accuracy is consistent across all cities. Larger values of N generally lead to improved forecasting performance, as demonstrated by the decrease in MAPE across different cities. This common trend suggests that capturing a longer historical context enhances model accuracy universally.

When comparing the different learning approaches for the three cities it occurs that the sliding window consistently outperforms the expanding window in Berlin, but this pattern is not universally observed across all cities. Seattle, for instance, exhibits closer performance between the two strategies, emphasizing the importance of considering city-specific characteristics in model selection.

The trade-off between model complexity and accuracy is evident in all cities. Choosing simpler models ($p = q = 2$) generally results in a marginal increase in MAPE, emphasizing the practical importance of simpler and more robust models.

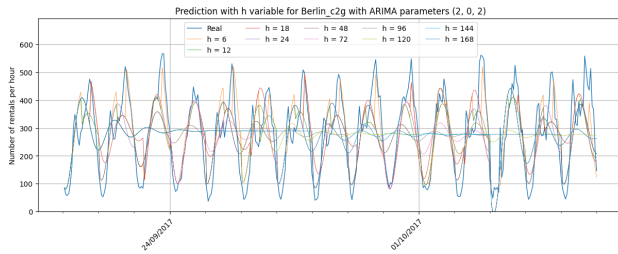
Considering the relative error (MAPE) in the context of absolute numbers of rentals is crucial for a comprehensive evaluation: in cities with smaller absolute rental numbers, higher relative errors may be observed.

1.8 (Optional) Impact of the time horizon h of the prediction on the performance

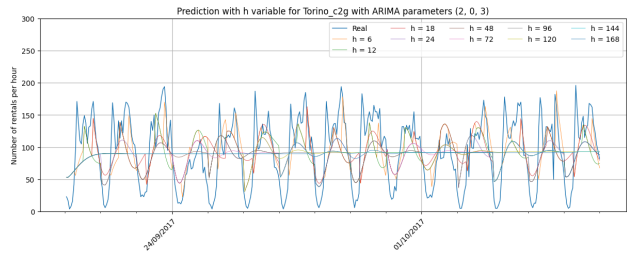
The time horizon h of an ARIMA model refers to the number of future time points to be forecast. In ARMA models, when the differencing component is equal to zero, the time horizon is the number of future time points to be predicted, it determines how far into the future the model will project.

When the number of points to be predicted at each iteration is increased ($h = 6$ (hours), 12, 18, 24, 48, 72, 96, 120, 144, 168 (7 days) points). While increasing h it can be noticed that the accuracy of the prediction decreases until becoming a horizontal time for Berlin and Torino. Fig. 6a shows the trend of the number of rentals per hour in Berlin through time, when changing the parameter h and keeping the ARIMA parameter described in sec. 1.7.2, and fig. 6b and 7b show a similar behaviour for the city of Torino. The just described trend

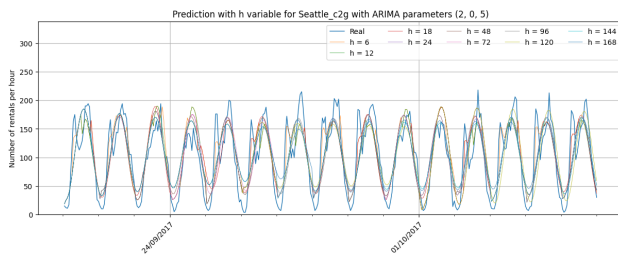
is not present in the time series regarding the city of Seattle, as shown in fig. 7a.



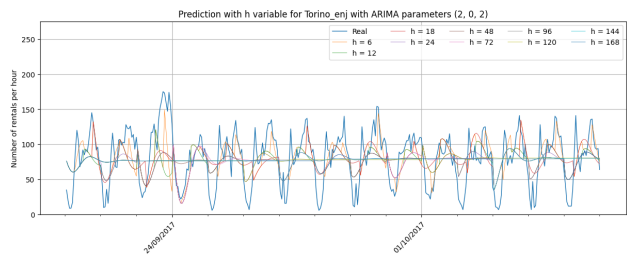
(a) Number of rentals per hour through time for Berlin with h variable.



(b) Number of rentals per hour through time for Torino with h variable.



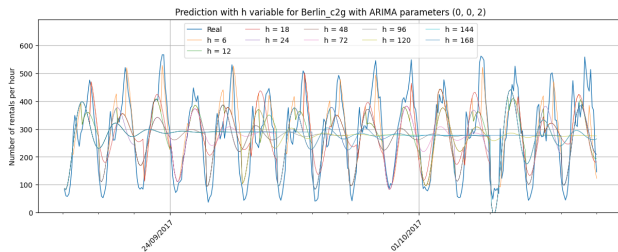
(a) Number of rentals per hour through time for Seattle with h variable.



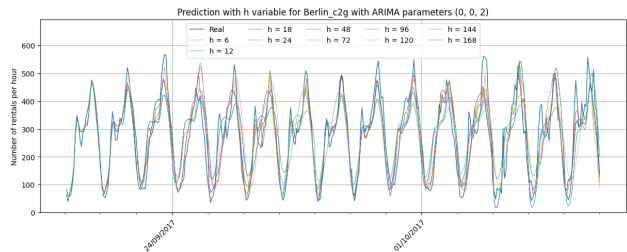
(b) Number of rentals per hour through time for Torino (Enjoy) with h variable.

The parameter h is then changed again as described above, but also the parameter p , the number of past data used for predictions, is changed for the analysis: $p=4, 8, 12, 24, 48$. The trade-off between the two quantities is what stands out from the graphs: when using a low value for the parameter p the short term accuracy is high, while the long term one is poor, also because of the presence of high values of h , resulting in a worsening of the performance when increasing.

When the chosen value for p is increased, the behaviour of the time series at the changing of the parameter h is similar to their real behaviour. The graphs for Berlin with the parameter $p=4$ and $p=24$ are respectively shown in fig. 8a and fig. 8b.



(a) Predictions with variable h and $p=4$ for Berlin.



(b) Predictions with variable h and $p=48$ for Berlin

— well done and complete
— missing baseline comparison

5/5

2 Appendix

2.1 Figures

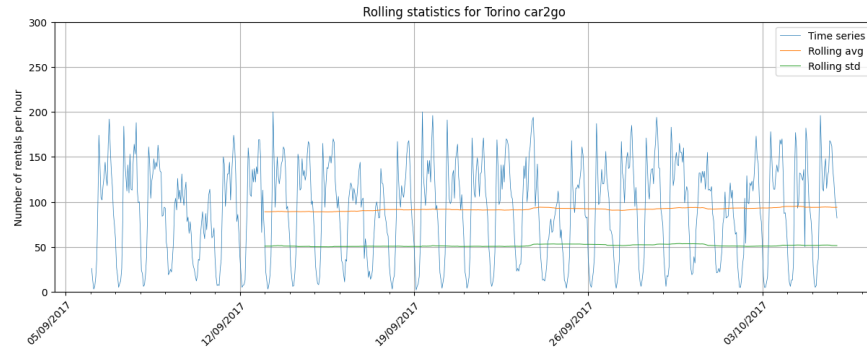


Figure 9: Rolling statistics with a 7-days window for Torino's data.

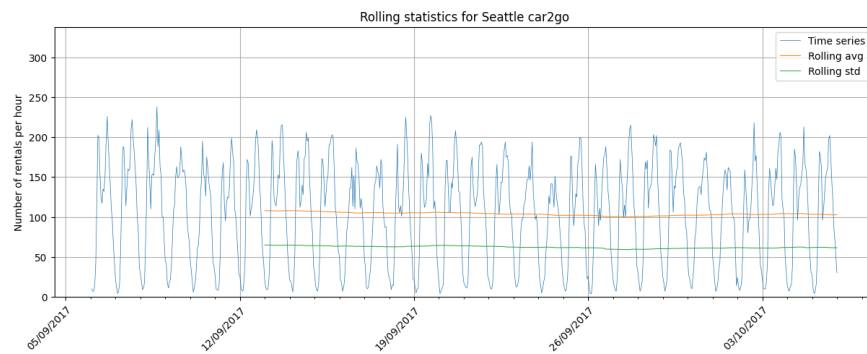


Figure 10: Rolling statistics with a 7-days window for Seattle's data.

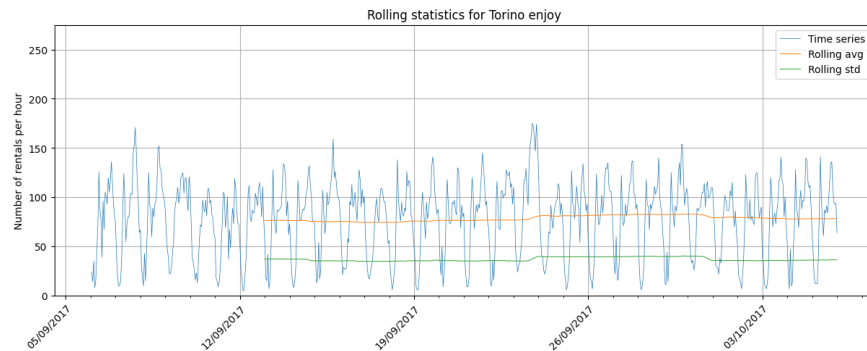


Figure 11: Rolling statistics with a 7-days window for Torino (Enjoy)'s data.

2.2 Tables

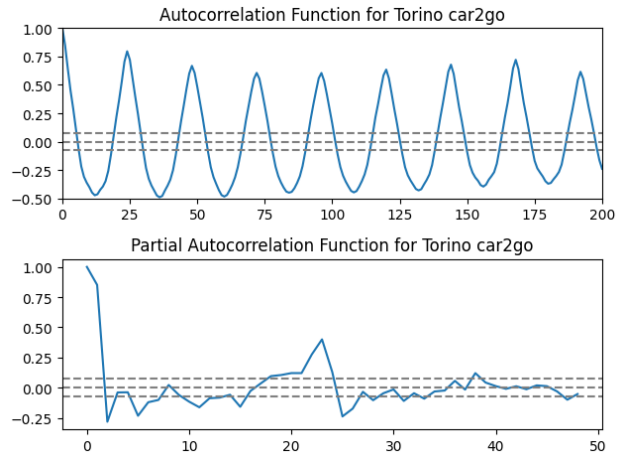


Figure 12: ACF and PACF for Berlin.

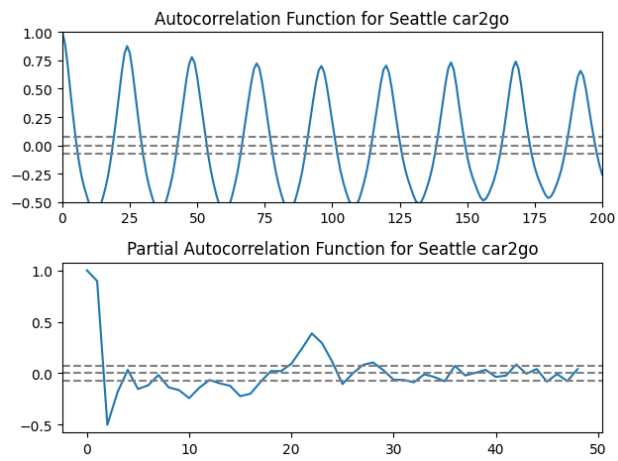


Figure 13: ACF and PACF for Berlin.

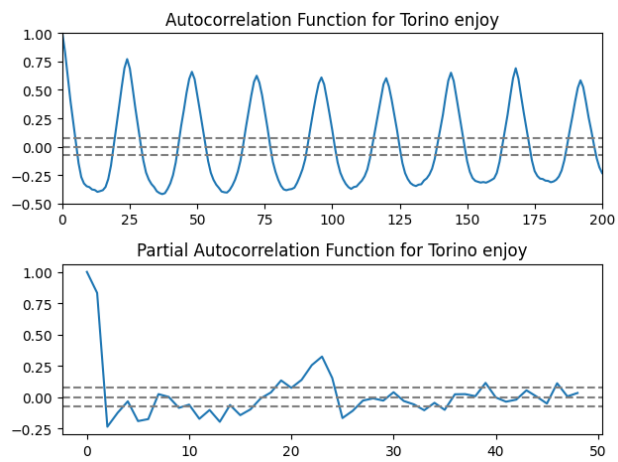


Figure 14: ACF and PACF for Berlin.

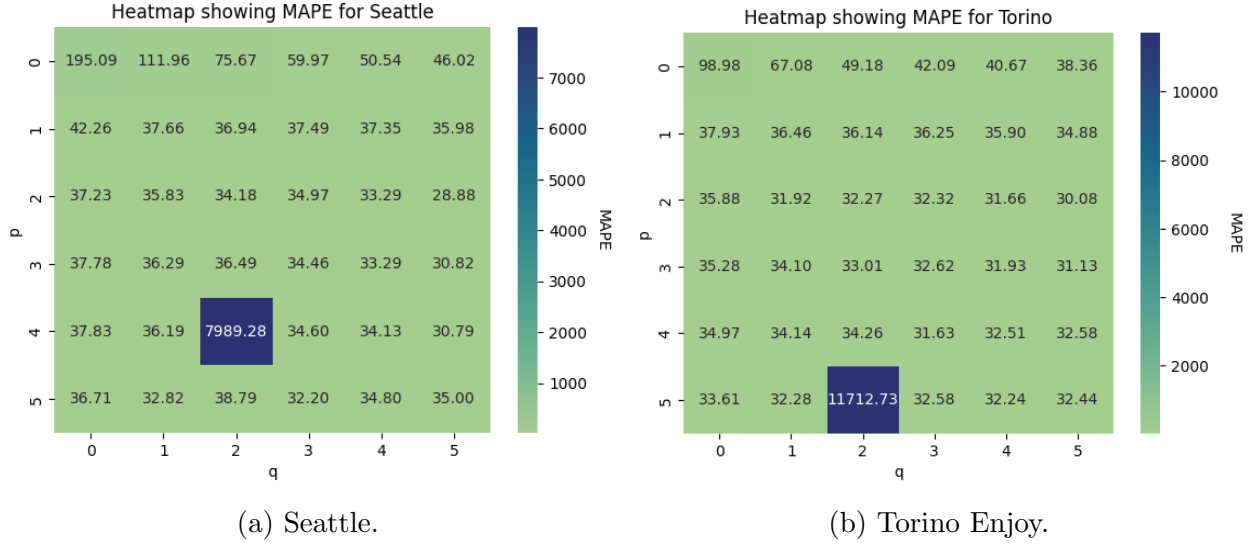


Figure 15: Heatmap showing MAPE.

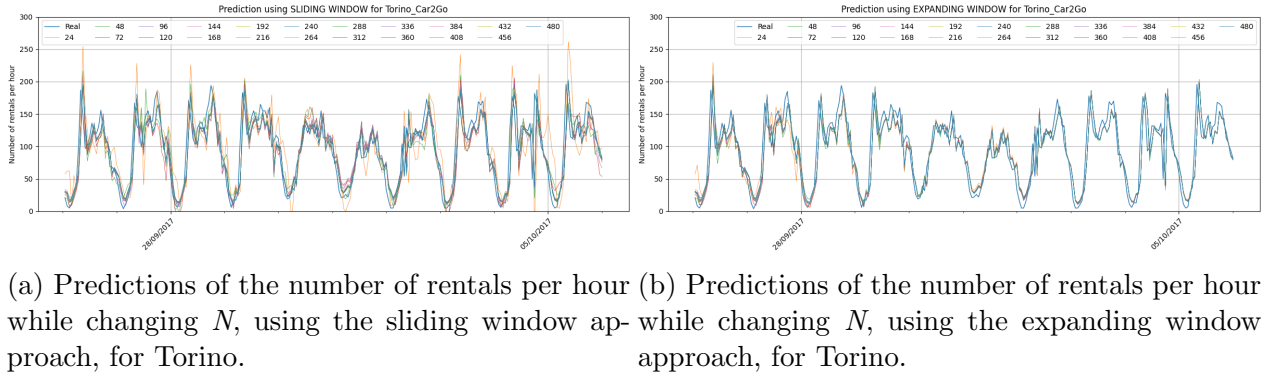


Figure 16: Torino.

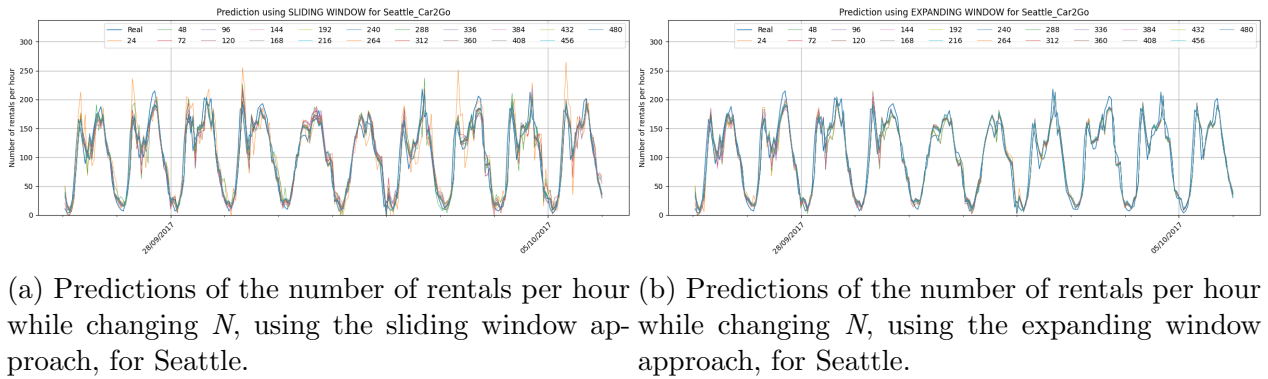
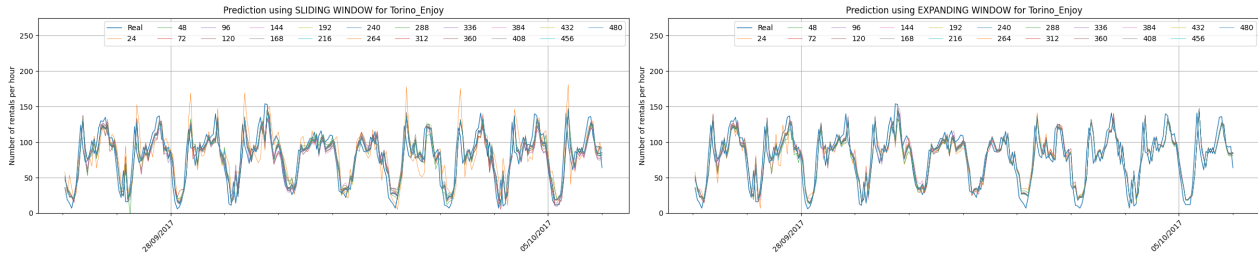


Figure 17: Seattle.



(a) Predictions of the number of rentals per hour while changing N , using the sliding window approach, for Torino_Enjoy. (b) Predictions of the number of rentals per hour while changing N , using the expanding window approach, for Torino_Enjoy.

Figure 18: Torino.

2.3 Code

```

1  # %%
2  import pymongo as pm # import MongoClient only
3  import pprint # prettyprinting for json objects
4  import datetime
5  import pytz
6  import matplotlib as mpl
7  import numpy as np
8  import matplotlib.pyplot as plt
9  import matplotlib.ticker as ticker
10 import matplotlib.dates as md
11 import pandas as pd
12 from statsmodels.tsa.stattools import acf, pacf
13 from statsmodels.tsa.arima.model import ARIMA
14 import numpy as np
15 import warnings
16
17 # %%
18 client = pm.MongoClient('bigdatadb.polito.it', ssl=True, authSource =
19     ↪ 'carsharing', username = 'ictts', password = 'Ict4SM22!',
20     ↪ tlsAllowInvalidCertificates=True)
21 db = client["carsharing"] # Choose the DB to use
22
23 c2g_PermanentBookings = db["PermanentBookings"]
24 enj_PermanentBookings = db["enjoy_PermanentBookings"]
25
26 # %%
27 our_cities = ["Berlin", "Torino", "Seattle"]
28
29 c2g_dates_PB_Berlin = []
30 c2g_count_PB_Berlin = []

```

	Sliding window		Expanding window	
N	MAPE	MPE	MAPE	MPE
24	72.190	-43.668	44.068	-25.895
48	42.281	-19.140	36.114	-17.953
72	41.745	-14.853	33.702	-14.023
96	36.414	-16.570	35.843	-18.381
120	37.661	-19.521	37.217	-20.559
144	39.600	-22.504	38.809	-22.445
168	39.721	-23.445	39.153	-23.073
192	39.026	-22.879	39.610	-23.581
216	38.292	-22.149	38.626	-22.225
240	37.562	-21.263	36.668	-19.919
264	37.219	-20.788	36.794	-20.085
288	37.420	-21.002	36.924	-20.416
312	38.253	-21.815	37.619	-21.237
336	38.698	-22.447	38.239	-22.062
360	38.708	-22.464	38.265	-22.108
384	38.165	-21.975	37.579	-21.092
408	37.486	-21.135	36.219	-19.433
432	37.153	-20.728	36.220	-19.408
456	37.498	-20.990	37.008	-20.412
480	38.845	-22.478	38.947	-22.358

(a) Parameter N tuning for Torino with $d = 0$, $p = 2$ and $q = 3$.

	Sliding window		Expanding window	
N	MAPE	MPE	MAPE	MPE
24	47.217	-22.375	34.199	-17.035
48	41.870	-20.662	29.166	-9.378
72	36.628	-16.558	25.154	-5.328
96	33.132	-16.218	28.003	-12.346
120	31.751	-14.786	27.802	-12.940
144	31.982	-17.526	29.638	-15.007
168	32.722	-17.821	30.768	-17.411
192	30.801	-15.264	30.489	-16.899
216	29.182	-14.988	28.881	-14.954
240	28.840	-14.846	28.019	-14.124
264	28.228	-14.213	29.800	-16.616
288	28.442	-14.900	29.918	-17.206
312	28.266	-14.950	30.054	-17.141
336	28.555	-15.125	30.692	-17.638
360	29.992	-16.871	30.715	-17.633
384	30.167	-16.932	29.610	-16.267
408	29.097	-16.067	29.601	-16.393
432	29.407	-16.540	29.704	-16.939
456	28.965	-16.059	30.514	-17.886
480	29.449	-16.482	30.817	-17.967

(b) Parameter N tuning for Seattle with $d = 0$, $p = 2$ and $q = 5$.

```

29 c2g_dates_PB_Torino = []
30 c2g_count_PB_Torino = []
31 c2g_dates_PB_Seattle = []
32 c2g_count_PB_Seattle = []
33
34 for city in our_cities:
35
36     # PermanentBookings
37     if city == "Berlin" or city == "Torino":
38         start_date = datetime.datetime(2017, 9, 6, 0, 0, 0, tzinfo =
39             ↪ datetime.timezone.utc).timestamp() - 3600
40         end_date = datetime.datetime(2017, 10, 6, 0, 0, 0, tzinfo =
41             ↪ datetime.timezone.utc).timestamp() - 3600
42     else: #city == "Seattle"
43         start_date = datetime.datetime(2017, 9, 6, 0, 0, 0, tzinfo =
44             ↪ datetime.timezone.utc).timestamp() + 28800
45         end_date = datetime.datetime(2017, 10, 6, 0, 0, 0, tzinfo =
46             ↪ datetime.timezone.utc).timestamp() + 28800
47
48     c2g_new_cursor_PB = c2g_PermanentBookings.aggregate([
49         { "$match": { "city": city,
50                     "init_time": {"$gte": start_date,

```

```

47         "$lt": end_date}}},
48     { "$project": { "_id": 0,
49                     "init_address": 1,
50                     "final_address": 1,
51                     "duration": { "$divide": [ {"$subtract":
52                                     ↪ ["$final_time", "$init_time"]}, 60]},
53                     "unique_timestamp": { "$dateFromString": {
54                                     ↪ "dateString": { "$dateToString": { "date":
55                                     ↪ "$init_date", "format": "%d-%m-%Y-%H"}}}},
56                     "equal": {"$eq": ["$init_address",
57                                     ↪ "$final_address"]}}}},
58     { "$match": { "equal": False,
59                 "duration": {"$gte": 5, "$lt": 120}}},
60     {"$group": {"_id": "$unique_timestamp",
61                "count": {"$sum": 1}}},
62     {"$project": {"_id": 1,
63                  "count": 1}},
64     {"$sort": {"_id": 1}}))
65
66 for elem in c2g_new_cursor_PB:
67     if city == "Berlin":
68         c2g_dates_PB_Berlin.append(elem["_id"])
69         c2g_count_PB_Berlin.append(elem["count"])
70     elif city == "Torino":
71         c2g_dates_PB_Torino.append(elem["_id"])
72         c2g_count_PB_Torino.append(elem["count"])
73     else:
74         c2g_dates_PB_Seattle.append(elem["_id"])
75         c2g_count_PB_Seattle.append(elem["count"])
76
77 # Introducing a random component to make the counts float
78 for i in c2g_count_PB_Berlin:
79     i = i + np.random.normal(0,1)
80 for i in c2g_count_PB_Torino:
81     i = i + np.random.normal(0,1)
82 for i in c2g_count_PB_Seattle:
83     i = i + np.random.normal(0,1)
84
85 # Creating the new lists of dates containing the missing values/dates
86 c2g_dates_PB_Berlin_UPDATED = []
87 c2g_dates_PB_Torino_UPDATED = []
88 c2g_dates_PB_Seattle_UPDATED = []
89
90 start = datetime.datetime(2017, 9, 6, 1, 0, 0)
91 end = datetime.datetime(2017, 10, 6, 0, 0, 0)

```



```

88 delta = datetime.timedelta(hours=1)
89
90 while start <= end:
91     c2g_dates_PB_Berlin_UPDATED.append(start)
92     c2g_dates_PB_Torino_UPDATED.append(start)
93     c2g_dates_PB_Seattle_UPDATED.append(start)
94     start += delta
95
96 # Creating the new lists of counts inserting 0 in the missing values
97 for i in range(len(c2g_dates_PB_Berlin_UPDATED)-1): #It could be also
98     ↪ c2g_dates_PB_Torino_UPDATED or c2g_dates_PB_Seattle_UPDATED
99
100     if c2g_dates_PB_Berlin_UPDATED[i] not in c2g_dates_PB_Berlin:
101         c2g_count_PB_Berlin.insert(i, 0)
102     if c2g_dates_PB_Torino_UPDATED[i] not in c2g_dates_PB_Torino:
103         c2g_count_PB_Torino.insert(i, 0)
104     if c2g_dates_PB_Seattle_UPDATED[i] not in c2g_dates_PB_Seattle:
105         c2g_count_PB_Seattle.insert(i, 0)
106
107 # %%
108 # DOING THE SAME THING FOR TORINO ENJOY
109
110 # PermanentBookings
111 enj_dates_PB_Torino = []
112 enj_count_PB_Torino = []
113
114 enj_new_cursor_PB = enj_PermanentBookings.aggregate([
115     { "$match": { "city": "Torino",
116         "init_time": {"$gte": datetime.datetime(2017, 9, 6,
117             ↪ 0, 0, 0).replace(tzinfo =
118             ↪ datetime.timezone.utc).timestamp() - 3600,
119             "$lt": datetime.datetime(2017, 10, 6,
120                 ↪ 0, 0, 0).replace(tzinfo =
121                 ↪ datetime.timezone.utc).timestamp()
122                 ↪ - 3600}}},
123     { "$project": { "_id": 0,
124         "init_address": 1,
125         "final_address": 1,
126         "duration": { "$divide": [ {"$subtract":
127             ↪ ["$final_time", "$init_time"]}, 60]},
128         "unique_timestamp": { "$dateFromString": {
129             ↪ "dateString": { "$dateToString": { "date":
130                 ↪ "$init_date", "format": "%d-%m-%Y-%H"}}}},
131         "equal": {"$eq": ["$init_address",
132             ↪ "$final_address"]}}},

```

```

123         { "$match": { "equal": False,
124                       "duration": {"$gte": 5, "$lt": 120}}},
125         {"$group": {"_id": "$unique_timestamp",
126                    "count": {"$sum": 1}}},
127         {"$project": {"_id": 1,
128                      "count": 1}},
129         {"$sort": {"_id": 1}}])
130
131 for elem in enj_new_cursor_PB:
132     enj_dates_PB_Torino.append(elem["_id"])
133     enj_count_PB_Torino.append(elem["count"])
134
135 # Introducing a random component to make the counts float
136 for i in enj_count_PB_Torino:
137     i = i + np.random.normal(0,1)
138
139 # Creating the new lists of dates containing the missing values/dates
140 enj_dates_PB_Torino_UPDATED = []
141
142 start_enj = datetime.datetime(2017, 9, 6, 1, 0, 0)
143 end_enj = datetime.datetime(2017, 10, 6, 0, 0, 0)
144 delta_enj = datetime.timedelta(hours=1)
145
146 while start_enj <= end_enj:
147     enj_dates_PB_Torino_UPDATED.append(start_enj)
148     start_enj += delta_enj
149
150 # Creating the new lists of counts inserting 0 in the missing values
151 for i in range(len(enj_dates_PB_Torino_UPDATED)-1): #It could be also
152     ↪ c2g_dates_PB_Torino_UPDATED or c2g_dates_PB_Seattle_UPDATED
153     if enj_dates_PB_Torino_UPDATED[i] not in enj_dates_PB_Torino:
154         enj_count_PB_Torino.insert(i, 0)
155
156 # %%
157 def substituteByItsMean(list_dates, list_counts):
158     for i in range(len(list_dates)-1): #Loop over dates length
159         if list_counts[i] == 0: # If this was a missing value
160             day = list_dates[i].weekday()
161             hour = list_dates[i].hour
162             sum = 0
163             count = 0
164             #print("\n" + str(i), " ", str(day), " ", str(hour))
165             for j in range(len(list_dates)-1):
166                 if list_dates[j].weekday() == day and list_dates[j].hour ==
167                     ↪ hour and list_dates[j] != list_dates[i]:

```

```

166         sum += list_counts[j]
167         count += 1
168         #print(str(list_dates[j].weekday()) + " " +
        ↪ str(list_dates[j].hour) + " " +
        ↪ str(list_counts[j]))
169     avg = sum/count
170     rounded_avg = round(avg)
171     #print(str(avg) + " " + str(rounded_avg))
172     list_counts[i] = rounded_avg
173     #print(list_counts[i])
174
175 def deleteOutliers(list_dates, list_counts):
176     for i in range(len(list_dates)-1): #Loop over dates length
177         day = list_dates[i].weekday()
178         hour = list_dates[i].hour
179         sum = 0
180         count = 0
181         #print("\n" + str(i), " ", str(day), " ", str(hour))
182         for j in range(len(list_dates)-1):
183             if list_dates[j].weekday() == day and list_dates[j].hour ==
            ↪ hour and list_dates[j] != list_dates[i]:
184                 sum += list_counts[j]
185                 count += 1
186                 #print(str(list_dates[j].weekday()) + " " +
            ↪ str(list_dates[j].hour) + " " + str(list_counts[j]))
187         avg = sum/count
188         rounded_avg = round(avg)
189         #print(str(avg) + " " + str(rounded_avg))
190         if list_counts[i] > 2*rounded_avg or list_counts[i] <
            ↪ 0.3*rounded_avg:
191             list_counts[i] = rounded_avg
192         #print(list_counts[i])
193
194 # %%
195 # METHOD TO REPLACE ZERO VALUES WITH THE MEAN OF THE SAME DAYS OF THE
    ↪ WEEK AT THE SAME HOURS
196 substituteByItsMean(c2g_dates_PB_Berlin_UPDATED, c2g_count_PB_Berlin)
197 substituteByItsMean(c2g_dates_PB_Torino_UPDATED, c2g_count_PB_Torino)
198 substituteByItsMean(c2g_dates_PB_Seattle_UPDATED, c2g_count_PB_Seattle)
199
200 deleteOutliers(c2g_dates_PB_Berlin_UPDATED, c2g_count_PB_Berlin)
201 deleteOutliers(c2g_dates_PB_Torino_UPDATED, c2g_count_PB_Torino)
202 deleteOutliers(c2g_dates_PB_Seattle_UPDATED, c2g_count_PB_Seattle)
203
204 substituteByItsMean(enj_dates_PB_Torino_UPDATED, enj_count_PB_Torino)

```

```

205
206 deleteOutliers(enj_dates_PB_Torino_UPDATED, enj_count_PB_Torino)
207
208 # %%
209 # Valutare se inserire il plot nel report!!!!!!
210 # fig1, ax1 = plt.subplots(layout="constrained")
211 # fig1.set_figwidth(12)
212 # ax1.plot(c2g_dates_PB_Berlin_UPDATED, c2g_count_PB_Berlin,
213 ↪ label="Berlin", linewidth=0.5)
214 # ax1.plot(c2g_dates_PB_Torino_UPDATED, c2g_count_PB_Torino,
215 ↪ label="Torino", linewidth=0.5)
216 # ax1.plot(c2g_dates_PB_Seattle_UPDATED, c2g_count_PB_Seattle,
217 ↪ label="Seattle", linewidth=0.5)
218 # ax1.plot(enj_dates_PB_Torino_UPDATED, enj_count_PB_Torino,
219 ↪ label="Torino E", linewidth=0.5)
220 # ax1.xaxis.set_major_locator(md.DayLocator(interval=7))
221 # ax1.xaxis.set_major_formatter(md.DateFormatter("%d/%m/%Y"))
222 # ax1.xaxis.set_minor_locator(md.DayLocator())
223 # plt.grid()
224 # plt.ylabel("Number of rentals per hour", rotation=90)
225 # plt.xticks(rotation=45, ha="right", rotation_mode="anchor")
226 # plt.ylim([0,2000])
227 # plt.title("Filtered system utilization over time - Bookings")
228 # plt.legend()
229 # #plt.savefig("/Users/marcoberti/Desktop/POLITO/ICT4SS/Secondo_anno/ICT_
230 ↪ for_Smart_Mobility/Mellia/Utilization_PB_cities_FILTERED.png")
231 # plt.show()
232
233 # %%
234 # Creating a dataframe for each city
235
236 c2g_Berlin = pd.DataFrame({"Date": c2g_dates_PB_Berlin_UPDATED,
237 ↪ "Rentals": c2g_count_PB_Berlin})
238 c2g_Torino = pd.DataFrame({"Date": c2g_dates_PB_Torino_UPDATED,
239 ↪ "Rentals": c2g_count_PB_Torino})
240 c2g_Seattle = pd.DataFrame({"Date": c2g_dates_PB_Seattle_UPDATED,
241 ↪ "Rentals": c2g_count_PB_Seattle})
242 enj_Torino = pd.DataFrame({"Date": enj_dates_PB_Torino_UPDATED,
243 ↪ "Rentals": enj_count_PB_Torino})
244
245 dataframes = [c2g_Berlin, c2g_Torino, c2g_Seattle, enj_Torino]
246 cities = ["Berlin car2go", "Torino car2go", "Seattle car2go", "Torino
247 ↪ enjoy"]
248
249 # %%

```

```

244 def plotRollingStats(df, city):
245     fig, ax = plt.subplots(layout="constrained")
246     fig.set_figwidth(12)
247     ax.plot(df["Date"], df["Rentals"], label="Time series", linewidth=0.5)
248     ax.plot(df["Date"], df["Rolling_avg"], label="Rolling avg",
249           ↪ linewidth=0.7)
250     ax.plot(df["Date"], df["Rolling_std"], label="Rolling std",
251           ↪ linewidth=0.7)
252     ax.xaxis.set_major_locator(md.DayLocator(interval=7))
253     ax.xaxis.set_major_formatter(md.DateFormatter("%d/%m/%Y"))
254     ax.xaxis.set_minor_locator(md.DayLocator())
255     plt.grid()
256     plt.ylabel("Number of rentals per hour", rotation=90)
257     plt.xticks(rotation=45, ha="right", rotation_mode="anchor")
258     plt.ylim([0, max(df["Rentals"])+100])
259     plt.title("Rolling statistics for " + city)
260     plt.legend()
261     #plt.savefig("/Users/marcoberti/Desktop/POLITO/ICT4SS/Secondo_anno/IC
262     ↪ T_for_Smart_Mobility/Mellia/Utilization_PB_cities_FILTERED.png")
263     plt.show()
264
265 def plot_ACF_PACF(df, city):
266     lag_acf = acf(df["Rentals"], nlags=200)
267     lag_pacf = pacf(df["Rentals"], nlags=48)
268
269     #Plot ACF:
270     plt.subplot(211)
271     plt.plot(lag_acf)
272     plt.axis([0, 200, -.5, 1])
273     plt.axhline(y=0, linestyle='--', color='gray')
274     plt.axhline(y=-1.96/np.sqrt(len(df["Rentals"])), linestyle='--', color='g
275     ↪ ray')
276     plt.axhline(y=1.96/np.sqrt(len(df["Rentals"])), linestyle='--', color='gr
277     ↪ ay')
278     plt.title('Autocorrelation Function for ' + city)
279
280     #Plot PACF:
281     plt.subplot(212)
282     plt.plot(lag_pacf)
283     plt.axhline(y=0, linestyle='--', color='gray')
284     plt.axhline(y=-1.96/np.sqrt(len(df["Rentals"])), linestyle='--', color='g
285     ↪ ray')
286     plt.axhline(y=1.96/np.sqrt(len(df["Rentals"])), linestyle='--', color='gr
287     ↪ ay')
288     plt.title('Partial Autocorrelation Function for ' + city)

```

```

282     plt.tight_layout()
283
284     plt.show()
285
286     # %%
287     # Adding "Rolling avg" and "Rolling var" for each time series..
288     for i in dataframes:
289         i['Rolling_avg'] = i.Rentals.rolling(168).mean() # 168 means 168 hours
290         ↪ in a week (= 24*7)
291         i['Rolling_std'] = i.Rentals.rolling(168).std()
292
293     # %%
294     # Plot rolling statistics for each city
295     for i in range(len(dataframes)):
296         plotRollingStats(dataframes[i], cities[i])
297
298     # %%
299     #Plot ACF and PACF for each city
300     for i in range(len(dataframes)):
301         plot_ACF_PACF(dataframes[i], cities[i])
302
303     # %%
304     def meanPercentageError(real, predicted):
305         sum = 0
306         for i in range(len(real)):
307             diff = real[i]-predicted[i]
308             div = diff/real[i]
309             sum += div
310         return sum/len(real)*100
311
312     def meanAbsolutePercentageError(real, predicted):
313         sum = 0
314         for i in range(len(real)):
315             diff = real[i]-predicted[i]
316             div = diff/real[i]
317             sum += abs(div)
318         return sum/len(real)*100
319
320     # %%
321     def slidingWindow(df, p, d, q, N, searchGrid, city):
322
323         train_len = N # 24 * 15 = 360
324         test_len = 720 - N # 720 - 360 = 360
325
326         prediction = np.zeros((test_len))

```

```

326
327     for i in range(0, test_len):
328         history = df["Rentals"][0 + i : train_len + i]
329         arima_model = ARIMA(history, order=(p,d,q),
330             ↪ enforce_stationarity=False)
331         fitted_model = arima_model.fit(method="statespace")
332         predicted_value = fitted_model.forecast()
333         prediction[i] = predicted_value.values[0]
334
335     mape = meanAbsolutePercentageError(df["Rentals"][test_len:].to_numpy(),
336         ↪ prediction)
337     mpe = meanPercentageError(df["Rentals"][test_len:].to_numpy(),
338         ↪ prediction)
339
340     print("\nMetrics for " + city + " with hyperparameters p=" + str(p) + "
341         ↪ and q=" + str(q))
342     print('MAPE: %.3f -- MPE: %.3f ' % (mape, mpe))
343
344     if searchGrid == False:
345         fig, ax = plt.subplots(layout="constrained")
346         fig.set_figwidth(12)
347         ax.plot(df["Date"][test_len:].to_numpy(),
348             ↪ df["Rentals"][test_len:].to_numpy(), label="Real values",
349             ↪ linewidth=0.5)
350         ax.plot(df["Date"][test_len:].to_numpy(), prediction,
351             ↪ label="Predicted values", linewidth=0.5)
352         ax.xaxis.set_major_locator(md.DayLocator(interval=7))
353         ax.xaxis.set_major_formatter(md.DateFormatter("%d/%m/%Y"))
354         ax.xaxis.set_minor_locator(md.DayLocator())
355         plt.grid()
356         plt.ylabel("Number of rentals per hour", rotation=90)
357         plt.xticks(rotation=45, ha="right", rotation_mode="anchor")
358         plt.ylim([0,max(df["Rentals"][test_len:].to_numpy())+100])
359         plt.title("Real and Predicted system utilization over time for " +
360             ↪ city)
361         plt.legend()
362         #plt.savefig("/Users/marcoberti/Desktop/POLITO/ICT4SS/Secondo_ann_
363             ↪ o/ICT_for_Smart_Mobility/Mellia/Utilization_PB_cities_FILTERE
364             ↪ D.png")
365         plt.show()
366
367     return (p, d, q), mape, mpe
368
369 # %%

```

```

360 # First attempt with p=2, d=0 (we've seen that the time series is
    ↳ stationary), and q=2
361 # Split training and test 50%/50% (360 hours of data/360 hours of data)
362 warnings.filterwarnings("ignore", category=UserWarning)
363 for i in range(len(dataframes)):
364     (p, d, q), mape, mpe = slidingWindow(dataframes[i], 2, 0, 2, 360,
    ↳ False, cities[i])
365
366 # %%
367 def gridSearch(df, city):
368
369     # Grid search with ranges p = q = [0,1,2,3,4,5] and d = 0
370     triplet_list = [] #saving (p,d,q)
371     mape_list = [] #saving MAPE
372     mpe_list = [] # saving MPE
373     for j in range(0, 6): #p ... the 6 is not included
374         for k in range(0, 6): #q ... the 6 is not included
375             pdq, mape, mpe = slidingWindow(df, j, 0, k, 360, True, city)
376             triplet_list.append(pdq)
377             mape_list.append(mape)
378             mpe_list.append(mpe)
379     result = pd.DataFrame({"p, d, q": triplet_list,
380                           "MAPE": mape_list,
381                           "MPE": mpe_list})
382     return result
383
384 # %%
385 warnings.filterwarnings("ignore", category=UserWarning)
386 for i in range(len(dataframes)):
387     output = gridSearch(dataframes[i], cities[i])
388     print(output)
389     if cities[i] == "Berlin car2go":
390         output.to_csv("Berlin_c2g_GridSearch.csv", index=False) # Best
    ↳ result: (2,0,2)
391
392     elif cities[i] == "Torino car2go":
393         output.to_csv("Torino_c2g_GridSearch.csv", index=False) # Best
    ↳ result: (2,0,5) but (2,0,3) can be used
394
395     elif cities[i] == "Seattle car2go":
396         output.to_csv("Seattle_c2g_GridSearch.csv", index=False) # Best
    ↳ result (2,0,5)
397
398     else: # cities[i] == "Torino enjoy":

```



```

399         output.to_csv("Torino_enj_GridSearch.csv", index=False) # Best
        ↪ result (2,0,5) but (2,0,2) can be used
400
401
402 # %%
403 def expandingWindow_N_variable(df, p, d, q, N_list):
404     # N_list is the list containing the various size of the training set
405     # test_len is the fixed time window over which we predict
406
407     train_len = 24 * 20
408     test_len = 24 * 10 # Keep always the last 10 days of data for
        ↪ predictions
409
410     prediction = np.zeros((len(N_list), test_len))
411
412     for i in range(len(N_list)):
413         print("\n\n" + str(N_list[i]) + "\n\n")
414         for j in range(0, test_len):
415             history = df["Rentals"][train_len - N_list[i] : train_len + j]
416             arima_model = ARIMA(history, order=(p,d,q),
                ↪ enforce_stationarity=False)
417             fitted_model = arima_model.fit(method="statespace")
418             predicted_value = fitted_model.forecast()
419             prediction[i, j] = predicted_value.values[0]
420
421     return prediction
422
423 def slidingWindow_N_variable(df, p, d, q, N_list):
424     # N_list is the list containing the various size of the training set
425     # test_len is the fixed time window over which we predict
426
427     train_len = 24 * 20
428     test_len = 24 * 10 # Keep always the last 10 days of data for
        ↪ predictions
429
430     prediction = np.zeros((len(N_list), test_len))
431
432     for i in range(len(N_list)):
433         print("\n\n" + str(N_list[i]) + "\n\n")
434         for j in range(0, test_len):
435             history = df["Rentals"][train_len - N_list[i] + j : train_len +
                ↪ j]
436             arima_model = ARIMA(history, order=(p,d,q),
                ↪ enforce_stationarity=False)
437             fitted_model = arima_model.fit(method="statespace")

```

```

438         predicted_value = fitted_model.forecast()
439         prediction[i , j] = predicted_value.values[0]
440
441     return prediction
442
443
444 # %%
445 N_list = np.arange(24, 504, 24) # 504 and not 480 because the last element
    ↪ is not included
446
447 # %%
448 warnings.filterwarnings("ignore", category=UserWarning)
449
450 p = 2; d = 0; q = 2      # Best results from grid search for Berlin c2g
451 print("Sliding Window for Berlin Car2Go with ARIMA parameters (%i,%i,%i)"
    ↪ % (p, d, q))
452 pred_SLIDING_Berlin_c2g = slidingWindow_N_variable(c2g_Berlin, p, d, q,
    ↪ N_list)
453
454 p = 2; d = 0; q = 3      # Best results from grid search for Berlin c2g
455 print("Sliding Window for Torino Car2Go with ARIMA parameters (%i,%i,%i)"
    ↪ % (p, d, q))
456 pred_SLIDING_Torino_c2g = slidingWindow_N_variable(c2g_Torino, p, d, q,
    ↪ N_list)
457
458 p = 2; d = 0; q = 5      # Best results from grid search for Berlin c2g
459 print("Sliding Window for Seattle Car2Go with ARIMA parameters (%i,%i,%i)"
    ↪ % (p, d, q))
460 pred_SLIDING_Seattle_c2g = slidingWindow_N_variable(c2g_Seattle, p, d, q,
    ↪ N_list)
461
462 p = 2; d = 0; q = 2      # Best results from grid search for Berlin c2g
463 print("Sliding Window for Torino Enjoy with ARIMA parameters (%i,%i,%i)" %
    ↪ (p, d, q))
464 pred_SLIDING_Torino_enj = slidingWindow_N_variable(enj_Torino, p, d, q,
    ↪ N_list)
465
466 # %%
467 def plot_sli_or_exp_windows(df, pred, city, sli_or_exp):
468     fig, ax = plt.subplots(layout="constrained")
469     fig.set_figwidth(12)
470     ax.plot(df["Date"][480:], df["Rentals"][480:], label="Real",
    ↪ linewidth=1)
471
472     for j in range(len(N_list)):

```

```

473         ax.plot(df["Date"][480:], pred[j], label = N_list[j], linewidth=0.5)
474
475     ax.xaxis.set_major_locator(md.DayLocator(interval=7))
476     ax.xaxis.set_major_formatter(md.DateFormatter("%d/%m/%Y"))
477     ax.xaxis.set_minor_locator(md.DayLocator())
478     plt.grid()
479     plt.ylabel("Number of rentals per hour", rotation=90)
480     plt.xticks(rotation=45, ha="right", rotation_mode="anchor")
481     plt.ylim([0,max(df["Rentals"])+100])
482     if sli_or_exp == 0:
483         plt.title("Prediction using SLIDING WINDOW for " + city)
484     else: #sli_or_exp == 1
485         plt.title("Prediction using EXPANDING WINDOW for " + city)
486     plt.legend(ncols=11)
487     if sli_or_exp == 0:
488         plt.savefig("/Users/marcoberti/Desktop/POLITO/ICT4SS/Secondo_anno/I_
         ↪ CT_for_Smart_Mobility/Mellia/Lab_2/SLIDING_W_N_Variable_" +
         ↪ city + ".png")
489     else: #sli_or_exp == 1
490         plt.savefig("/Users/marcoberti/Desktop/POLITO/ICT4SS/Secondo_anno/I_
         ↪ CT_for_Smart_Mobility/Mellia/Lab_2/EXPANDING_W_N_Variable_" +
         ↪ city + ".png")
491     plt.show()
492
493 # %%
494 plot_sli_or_exp_windows(c2g_Berlin, pred_SLIDING_Berlin_c2g,
495     ↪ "Berlin_Car2Go", 0)
496 plot_sli_or_exp_windows(c2g_Torino, pred_SLIDING_Torino_c2g,
497     ↪ "Torino_Car2Go", 0)
498 plot_sli_or_exp_windows(c2g_Seattle, pred_SLIDING_Seattle_c2g,
499     ↪ "Seattle_Car2Go", 0)
500 plot_sli_or_exp_windows(enj_Torino, pred_SLIDING_Torino_enj,
501     ↪ "Torino_Enjoy", 0)
502
503 # %%
504 warnings.filterwarnings("ignore", category=UserWarning)
505
506 p = 2; d = 0; q = 2      # Best results from grid search for Berlin c2g
507 print("Expanding Window for Berlin Car2Go with ARIMA parameters (%i,%i,%i)"
508     ↪ % (p, d, q))
509 pred_EXPANDING_Berlin_c2g = expandingWindow_N_variable(c2g_Berlin, p, d, q,
510     ↪ N_list)
511
512 p = 2; d = 0; q = 3      # Best results from grid search for Berlin c2g

```

```

507 print("Expanding Window for Torino Car2Go with ARIMA parameters (%i,%i,%i)"
      ↪ % (p, d, q))
508 pred_EXPANDING_Torino_c2g = expandingWindow_N_variable(c2g_Torino, p, d, q,
      ↪ N_list)
509
510 p = 2; d = 0; q = 5      # Best results from grid search for Berlin c2g
511 print("Expanding Window for Seattle Car2Go with ARIMA parameters
      ↪ (%i,%i,%i)" % (p, d, q))
512 pred_EXPANDING_Seattle_c2g = expandingWindow_N_variable(c2g_Seattle, p, d,
      ↪ q, N_list)
513
514 p = 2; d = 0; q = 2      # Best results from grid search for Berlin c2g
515 print("Expanding Window for Torino Enjoy with ARIMA parameters (%i,%i,%i)"
      ↪ % (p, d, q))
516 pred_EXPANDING_Torino_enj = expandingWindow_N_variable(enj_Torino, p, d, q,
      ↪ N_list)
517
518 # %%
519 plot_sli_or_exp_windows(c2g_Berlin, pred_EXPANDING_Berlin_c2g,
      ↪ "Berlin_Car2Go", 1)
520 plot_sli_or_exp_windows(c2g_Torino, pred_EXPANDING_Torino_c2g,
      ↪ "Torino_Car2Go", 1)
521 plot_sli_or_exp_windows(c2g_Seattle, pred_EXPANDING_Seattle_c2g,
      ↪ "Seattle_Car2Go", 1)
522 plot_sli_or_exp_windows(enj_Torino, pred_EXPANDING_Torino_enj,
      ↪ "Torino_Enjoy", 1)
523
524 # %%
525 def compute_MAPE_and_MPE_per_Prediction(df, predicted, N_list):
526     real = df["Rentals"][480:].to_numpy()
527     MAPE_list = []
528     MPE_list = []
529     for row in predicted: #predicted is the bidimensional array
530         mape = meanAbsolutePercentageError(real, row)
531         mpe = meanPercentageError(real, row)
532         MAPE_list.append(mape)
533         MPE_list.append(mpe)
534     result = pd.DataFrame({"Train": N_list,
535                           "MAPE": MAPE_list,
536                           "MPE": MPE_list})
537     return result
538
539 # %%
540 stats_SLIDING_Berlin_c2g = compute_MAPE_and_MPE_per_Prediction(c2g_Berlin,
      ↪ pred_SLIDING_Berlin_c2g, N_list)

```

```

541 stats_SLIDING_Berlin_c2g.to_csv("Stats_SLIDING_Berlin_c2g.csv", index=False)
542 stats_SLIDING_Torino_c2g = compute_MAPE_and_MPE_per_Prediction(c2g_Torino,
    ↪ pred_SLIDING_Torino_c2g, N_list)
543 stats_SLIDING_Torino_c2g.to_csv("Stats_SLIDING_Torino_c2g.csv", index=False)
544 stats_SLIDING_Seattle_c2g =
    ↪ compute_MAPE_and_MPE_per_Prediction(c2g_Seattle,
    ↪ pred_SLIDING_Seattle_c2g, N_list)
545 stats_SLIDING_Seattle_c2g.to_csv("Stats_SLIDING_Seattle_c2g.csv",
    ↪ index=False)
546 stats_SLIDING_Torino_enj = compute_MAPE_and_MPE_per_Prediction(enj_Torino,
    ↪ pred_SLIDING_Torino_enj, N_list)
547 stats_SLIDING_Torino_enj.to_csv("Stats_SLIDING_Torino_enj.csv", index=False)
548
549 stats_EXPANDING_Berlin_c2g =
    ↪ compute_MAPE_and_MPE_per_Prediction(c2g_Berlin,
    ↪ pred_EXPANDING_Berlin_c2g, N_list)
550 stats_EXPANDING_Berlin_c2g.to_csv("Stats_EXPANDING_Berlin_c2g.csv",
    ↪ index=False)
551 stats_EXPANDING_Torino_c2g =
    ↪ compute_MAPE_and_MPE_per_Prediction(c2g_Torino,
    ↪ pred_EXPANDING_Torino_c2g, N_list)
552 stats_EXPANDING_Torino_c2g.to_csv("Stats_EXPANDING_Torino_c2g.csv",
    ↪ index=False)
553 stats_EXPANDING_Seattle_c2g =
    ↪ compute_MAPE_and_MPE_per_Prediction(c2g_Seattle,
    ↪ pred_EXPANDING_Seattle_c2g, N_list)
554 stats_EXPANDING_Seattle_c2g.to_csv("Stats_EXPANDING_Seattle_c2g.csv",
    ↪ index=False)
555 stats_EXPANDING_Torino_enj =
    ↪ compute_MAPE_and_MPE_per_Prediction(enj_Torino,
    ↪ pred_EXPANDING_Torino_enj, N_list)
556 stats_EXPANDING_Torino_enj.to_csv("Stats_EXPANDING_Torino_enj.csv",
    ↪ index=False)
557
558 # %%
559 warnings.filterwarnings("ignore", category=UserWarning)
560
561 #This method can be used at first for the optimal arima parameters for
    ↪ each city, and then varying p!
562 def predict_more_values(df, p, d, q, N, h):
563
564     train_len = N # 24 * 15 = 360
565     test_len = 720 - N # 720 - 360 = 360
566
567     prediction = np.zeros((test_len))

```

```

568
569     for i in range(0, test_len, h):
570         if train_len+h < len(df["Rentals"]):
571             history = df["Rentals"][i : train_len + i]
572         else:
573             history = df["Rentals"][i :]
574         arima_model = ARIMA(history, order=(p,d,q),
575             → enforce_stationarity=False)
576         fitted_model = arima_model.fit(method="statespace")
577         predicted_value = fitted_model.forecast(steps=h)
578         if i+h < test_len:
579             prediction[i:i+h] = predicted_value.values
580         else:
581             prediction[i:] = predicted_value.values[:test_len%i]
582
583     return prediction, h
584
585 # %%
586 for i in range(len(dataframes)):
587     if cities[i] == "Berlin car2go":
588         p = 2; d = 0; q = 2
589         h6, n6 = predict_more_values(dataframes[i], p, d, q, 360, 6)
590         h12, n12 = predict_more_values(dataframes[i], p, d, q, 360, 12)
591         h18, n18 = predict_more_values(dataframes[i], p, d, q, 360, 18)
592         h24, n24 = predict_more_values(dataframes[i], p, d, q, 360, 24)
593         h48, n48 = predict_more_values(dataframes[i], p, d, q, 360, 48)
594         h72, n72 = predict_more_values(dataframes[i], p, d, q, 360, 72)
595         h96, n96 = predict_more_values(dataframes[i], p, d, q, 360, 96)
596         h120, n120 = predict_more_values(dataframes[i], p, d, q, 360, 120)
597         h144, n144 = predict_more_values(dataframes[i], p, d, q, 360, 144)
598         h168, n168 = predict_more_values(dataframes[i], p, d, q, 360, 168)
599         print("Berlin car2go DONE")
600
601         h_list = [h6, h12, h18, h24, h48, h72, h96, h120, h144, h168]
602         n_list = [n6, n12, n18, n24, n48, n72, n96, n120, n144, n168]
603
604         listone_Berlin_c2g = [h_list, n_list]
605
606     elif cities[i] == "Torino car2go":
607         p = 2; d = 0; q = 3
608         h6, n6 = predict_more_values(dataframes[i], p, d, q, 360, 6)
609         h12, n12 = predict_more_values(dataframes[i], p, d, q, 360, 12)
610         h18, n18 = predict_more_values(dataframes[i], p, d, q, 360, 18)
611         h24, n24 = predict_more_values(dataframes[i], p, d, q, 360, 24)

```

```

612     h48, n48 = predict_more_values(dataframes[i], p, d, q, 360, 48)
613     h72, n72 = predict_more_values(dataframes[i], p, d, q, 360, 72)
614     h96, n96 = predict_more_values(dataframes[i], p, d, q, 360, 96)
615     h120, n120 = predict_more_values(dataframes[i], p, d, q, 360, 120)
616     h144, n144 = predict_more_values(dataframes[i], p, d, q, 360, 144)
617     h168, n168 = predict_more_values(dataframes[i], p, d, q, 360, 168)
618     print("Torino car2go DONE")
619
620     h_list = [h6, h12, h18, h24, h48, h72, h96, h120, h144, h168]
621     n_list = [n6, n12, n18, n24, n48, n72, n96, n120, n144, n168]
622
623     listone_Torino_c2g = [h_list, n_list]
624
625     elif cities[i] == "Seattle car2go":
626         p = 2; d = 0; q = 5
627         h6, n6 = predict_more_values(dataframes[i], p, d, q, 360, 6)
628         h12, n12 = predict_more_values(dataframes[i], p, d, q, 360, 12)
629         h18, n18 = predict_more_values(dataframes[i], p, d, q, 360, 18)
630         h24, n24 = predict_more_values(dataframes[i], p, d, q, 360, 24)
631         h48, n48 = predict_more_values(dataframes[i], p, d, q, 360, 48)
632         h72, n72 = predict_more_values(dataframes[i], p, d, q, 360, 72)
633         h96, n96 = predict_more_values(dataframes[i], p, d, q, 360, 96)
634         h120, n120 = predict_more_values(dataframes[i], p, d, q, 360, 120)
635         h144, n144 = predict_more_values(dataframes[i], p, d, q, 360, 144)
636         h168, n168 = predict_more_values(dataframes[i], p, d, q, 360, 168)
637         print("Seattle car2go DONE")
638
639         h_list = [h6, h12, h18, h24, h48, h72, h96, h120, h144, h168]
640         n_list = [n6, n12, n18, n24, n48, n72, n96, n120, n144, n168]
641
642         listone_Seattle_c2g = [h_list, n_list]
643
644     else: # cities[i] == "Torino enjoy":
645         p = 2; d = 0; q = 2
646         h6, n6 = predict_more_values(dataframes[i], p, d, q, 360, 6)
647         h12, n12 = predict_more_values(dataframes[i], p, d, q, 360, 12)
648         h18, n18 = predict_more_values(dataframes[i], p, d, q, 360, 18)
649         h24, n24 = predict_more_values(dataframes[i], p, d, q, 360, 24)
650         h48, n48 = predict_more_values(dataframes[i], p, d, q, 360, 48)
651         h72, n72 = predict_more_values(dataframes[i], p, d, q, 360, 72)
652         h96, n96 = predict_more_values(dataframes[i], p, d, q, 360, 96)
653         h120, n120 = predict_more_values(dataframes[i], p, d, q, 360, 120)
654         h144, n144 = predict_more_values(dataframes[i], p, d, q, 360, 144)
655         h168, n168 = predict_more_values(dataframes[i], p, d, q, 360, 168)
656         print("Torino enjoy DONE")

```

```

657
658     h_list = [h6, h12, h18, h24, h48, h72, h96, h120, h144, h168]
659     n_list = [n6, n12, n18, n24, n48, n72, n96, n120, n144, n168]
660
661     listone_Torino_enj = [h_list, n_list]
662
663
664     # %%
665     def plot_different_h(df, listone, city, p_variable, p_val):
666
667         if p_variable == False:
668             if city == "Berlin_c2g":
669                 p = 2; d = 0; q = 2
670             elif city == "Torino_c2g":
671                 p = 2; d = 0; q = 3
672             elif city == "Seattle_c2g":
673                 p = 2; d = 0; q = 5
674             else: #city == "Torino_enj"
675                 p = 2; d = 0; q = 2
676
677             fig, ax = plt.subplots(layout="constrained")
678             fig.set_figwidth(12)
679             ax.plot(df["Date"][360:], df["Rentals"][360:], label="Real",
680                 ↪ linewidth=1)
681
682             for j in range(len(listone[0])):
683                 ax.plot(df["Date"][360:], listone[0][j], label = f"h =
684                 ↪ {listone[1][j]}", linewidth=0.5)
685
686             ax.xaxis.set_major_locator(md.DayLocator(interval=7))
687             ax.xaxis.set_major_formatter(md.DateFormatter("%d/%m/%Y"))
688             ax.xaxis.set_minor_locator(md.DayLocator())
689             plt.grid()
690             plt.ylabel("Number of rentals per hour", rotation=90)
691             plt.xticks(rotation=45, ha="right", rotation_mode="anchor")
692             plt.ylim([0, max(df["Rentals"])+100])
693             plt.title("Prediction with h variable for " + city + " with ARIMA
694             ↪ parameters (%i, %i, %i)" % (p, d, q))
695             plt.legend(ncols=5)
696             plt.savefig("/Users/marcoberti/Desktop/POLITO/ICT4SS/Secondo_anno/I
697             ↪ CT_for_Smart_Mobility/Mellia/Lab_2/Plot_h_variable_p=" + str(p)
698             ↪ + "_" + city + ".png")
699             plt.show()
700
701         else: # p_variable == True

```



```

697     if city == "Berlin_c2g":
698         d = 0; q = 2
699     elif city == "Torino_c2g":
700         d = 0; q = 3
701     elif city == "Seattle_c2g":
702         d = 0; q = 5
703     else: #city == "Torino_enj"
704         d = 0; q = 2
705
706     fig, ax = plt.subplots(layout="constrained")
707     fig.set_figwidth(12)
708     ax.plot(df["Date"][360:], df["Rentals"][360:], label="Real",
709             ↪ linewidth=1)
710
711     for j in range(len(listone[0])):
712         ax.plot(df["Date"][360:], listone[0][j], label = f"h =
713             ↪ {listone[1][j]}", linewidth=0.5)
714
715     ax.xaxis.set_major_locator(md.DayLocator(interval=7))
716     ax.xaxis.set_major_formatter(md.DateFormatter("%d/%m/%Y"))
717     ax.xaxis.set_minor_locator(md.DayLocator())
718     plt.grid()
719     plt.ylabel("Number of rentals per hour", rotation=90)
720     plt.xticks(rotation=45, ha="right", rotation_mode="anchor")
721     plt.ylim([0, max(df["Rentals"])+100])
722     plt.title("Prediction with h variable for " + city + " with ARIMA
723             ↪ parameters (%i, %i, %i)" % (p_val, d, q))
724     plt.legend(ncols=5)
725     plt.savefig("/Users/marcoberti/Desktop/POLITO/ICT4SS/Secondo_anno/I
726             ↪ CT_for_Smart_Mobility/Mellia/Lab_2/Plot_h_variable_p=" +
727             ↪ str(p_val) + "_" + city + ".png")
728     plt.show()
729
730 # %%
731 plot_different_h(c2g_Berlin, listone_Berlin_c2g, "Berlin_c2g", False, 0)
732 plot_different_h(c2g_Torino, listone_Torino_c2g, "Torino_c2g", False, 0)
733 plot_different_h(c2g_Seattle, listone_Seattle_c2g, "Seattle_c2g", False, 0)
734 plot_different_h(enj_Torino, listone_Torino_enj, "Torino_enj", False, 0)
735
736 # %%
737 warnings.filterwarnings("ignore", category=UserWarning)
738
739 p_list = [4, 8, 12, 24, 48]

```

```

737 for p in p_list:
738     for i in range(len(dataframes)):
739         if cities[i] == "Berlin car2go":
740             d = 0; q = 2
741             h6, n6 = predict_more_values(dataframes[i], p, d, q, 360, 6)
742             h12, n12 = predict_more_values(dataframes[i], p, d, q, 360, 12)
743             h18, n18 = predict_more_values(dataframes[i], p, d, q, 360, 18)
744             h24, n24 = predict_more_values(dataframes[i], p, d, q, 360, 24)
745             h48, n48 = predict_more_values(dataframes[i], p, d, q, 360, 48)
746             h72, n72 = predict_more_values(dataframes[i], p, d, q, 360, 72)
747             h96, n96 = predict_more_values(dataframes[i], p, d, q, 360, 96)
748             h120, n120 = predict_more_values(dataframes[i], p, d, q, 360,
749                 ↪ 120)
750             h144, n144 = predict_more_values(dataframes[i], p, d, q, 360,
751                 ↪ 144)
752             h168, n168 = predict_more_values(dataframes[i], p, d, q, 360,
753                 ↪ 168)
754             print("Berlin car2go DONE")
755
756             h_list = [h6, h12, h18, h24, h48, h72, h96, h120, h144, h168]
757             n_list = [n6, n12, n18, n24, n48, n72, n96, n120, n144, n168]
758
759             listone_Berlin_c2g = [h_list, n_list]
760
761         elif cities[i] == "Torino car2go":
762             d = 0; q = 3
763             h6, n6 = predict_more_values(dataframes[i], p, d, q, 360, 6)
764             h12, n12 = predict_more_values(dataframes[i], p, d, q, 360, 12)
765             h18, n18 = predict_more_values(dataframes[i], p, d, q, 360, 18)
766             h24, n24 = predict_more_values(dataframes[i], p, d, q, 360, 24)
767             h48, n48 = predict_more_values(dataframes[i], p, d, q, 360, 48)
768             h72, n72 = predict_more_values(dataframes[i], p, d, q, 360, 72)
769             h96, n96 = predict_more_values(dataframes[i], p, d, q, 360, 96)
770             h120, n120 = predict_more_values(dataframes[i], p, d, q, 360,
771                 ↪ 120)
772             h144, n144 = predict_more_values(dataframes[i], p, d, q, 360,
773                 ↪ 144)
774             h168, n168 = predict_more_values(dataframes[i], p, d, q, 360,
775                 ↪ 168)
776             print("Torino car2go DONE")
777
778             h_list = [h6, h12, h18, h24, h48, h72, h96, h120, h144, h168]
779             n_list = [n6, n12, n18, n24, n48, n72, n96, n120, n144, n168]
780
781             listone_Torino_c2g = [h_list, n_list]

```

```

776
777 elif cities[i] == "Seattle car2go":
778     d = 0; q = 5
779     h6, n6 = predict_more_values(dataframes[i], p, d, q, 360, 6)
780     h12, n12 = predict_more_values(dataframes[i], p, d, q, 360, 12)
781     h18, n18 = predict_more_values(dataframes[i], p, d, q, 360, 18)
782     h24, n24 = predict_more_values(dataframes[i], p, d, q, 360, 24)
783     h48, n48 = predict_more_values(dataframes[i], p, d, q, 360, 48)
784     h72, n72 = predict_more_values(dataframes[i], p, d, q, 360, 72)
785     h96, n96 = predict_more_values(dataframes[i], p, d, q, 360, 96)
786     h120, n120 = predict_more_values(dataframes[i], p, d, q, 360,
787         ↪ 120)
788     h144, n144 = predict_more_values(dataframes[i], p, d, q, 360,
789         ↪ 144)
790     h168, n168 = predict_more_values(dataframes[i], p, d, q, 360,
791         ↪ 168)
792     print("Seattle car2go DONE")
793
794     h_list = [h6, h12, h18, h24, h48, h72, h96, h120, h144, h168]
795     n_list = [n6, n12, n18, n24, n48, n72, n96, n120, n144, n168]
796
797     listone_Seattle_c2g = [h_list, n_list]
798
799 else: # cities[i] == "Torino enjoy":
800     d = 0; q = 2
801     h6, n6 = predict_more_values(dataframes[i], p, d, q, 360, 6)
802     h12, n12 = predict_more_values(dataframes[i], p, d, q, 360, 12)
803     h18, n18 = predict_more_values(dataframes[i], p, d, q, 360, 18)
804     h24, n24 = predict_more_values(dataframes[i], p, d, q, 360, 24)
805     h48, n48 = predict_more_values(dataframes[i], p, d, q, 360, 48)
806     h72, n72 = predict_more_values(dataframes[i], p, d, q, 360, 72)
807     h96, n96 = predict_more_values(dataframes[i], p, d, q, 360, 96)
808     h120, n120 = predict_more_values(dataframes[i], p, d, q, 360,
809         ↪ 120)
810     h144, n144 = predict_more_values(dataframes[i], p, d, q, 360,
811         ↪ 144)
812     h168, n168 = predict_more_values(dataframes[i], p, d, q, 360,
813         ↪ 168)
814     print("Torino enjoy DONE")
815
816     h_list = [h6, h12, h18, h24, h48, h72, h96, h120, h144, h168]
817     n_list = [n6, n12, n18, n24, n48, n72, n96, n120, n144, n168]
818
819     listone_Torino_enj = [h_list, n_list]

```

```
815 plot_different_h(c2g_Berlin, listone_Berlin_c2g, "Berlin_c2g", True, 0)
816 plot_different_h(c2g_Torino, listone_Torino_c2g, "Torino_c2g", True, 0)
817 plot_different_h(c2g_Seattle, listone_Seattle_c2g, "Seattle_c2g", True,
    ↪ 0)
818 plot_different_h(enj_Torino, listone_Torino_enj, "Torino_enj", True, 0)
819
820
821
822
```

	Sliding window		Expanding window	
N	MAPE	MPE	MAPE	MPE
24	53.239	-32.106	40.273	-22.496
48	36.803	-18.796	35.416	-16.758
72	35.394	-17.269	34.490	-17.422
96	34.775	-15.660	34.949	-19.622
120	35.490	-17.607	34.528	-19.774
144	36.627	-19.173	36.680	-20.740
168	37.062	-20.548	35.405	-20.238
192	36.743	-21.317	34.634	-19.389
216	35.306	-20.092	33.963	-18.184
240	35.485	-20.084	33.656	-17.561
264	34.229	-18.550	33.962	-17.892
288	34.268	-18.505	33.898	-17.880
312	34.420	-18.763	34.058	-18.066
336	34.486	-18.771	33.933	-17.844
360	34.349	-18.429	34.025	-17.439
384	34.004	-18.020	33.478	-16.502
408	33.821	-17.554	33.112	-16.148
432	33.767	-17.298	33.152	-16.347
456	33.785	-17.324	33.380	-16.557
480	33.959	-17.380	33.647	-17.075

Table 3: Parameter N tuning for Torino Enjoy with $d = 0$, $p = 2$ and $q = 2$.