**Politecnico di Torino**
**ICT for smart mobility**

Laboratory 2 - report
Group 2

Marika Attanasio - s311298,
Miriana Passarotto - s307735,
Cristiano Vittori - s316801

Academic Year 2023-2024

# 1 Task 1: Building the time series

The second laboratory's first task requires selecting 30 days and collecting the recorded bookings for each hour after filtering the outliers out.

Considering periods with registered data for both services (*Car2Go* and *Enjoy*), the month having the least missing data is January. However, in the first days of the month, the behavior of the users is strongly influenced by the holidays. Instead in October, even though we have more missing data, the number of rentals has a regular behavior over the whole month. Due to these reasons, October was chosen as the month of interest.

Concerning data filtering, the same filters implemented in the first lab activity are applied in this case. This resulted in excluding bookings lasting less than 5 minutes or more than 2 hours from the analysis. In addition to those, the bookings having the same coordinates for the initial and final positions have been removed.[1]

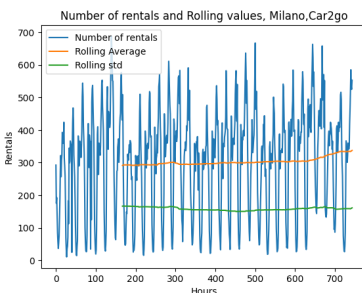The cities which are taken into consideration are Milan (*Car2Go* and *Enjoy*) and Rome (*Car2Go*).

# 2 Task 2: Fitting missing data

For this second task, missing data were identified and substituted. The adopted substitution policy is to check which hours of the day are missing and then substitute those data with the average number of rentals in the same hours calculated on the other days.[1]
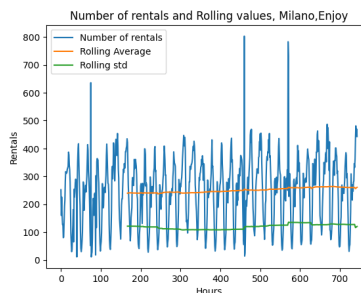
This step is essential for our algorithm to work efficiently since ARIMA models assume a regular time series, with no missing data.
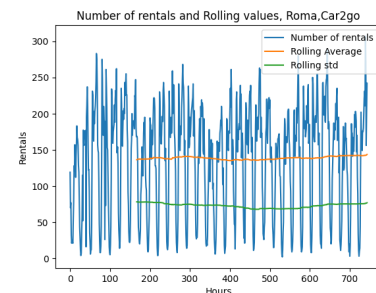
# 3 Task3: Stationarity check

Once missing data were fitted, it was possible to check for stationarity and to decide whether to use differencing or not. To perform this particular analysis, the rolling average and rolling standard deviation were calculated.[2] Figure 1 depicts them.



(a) Milan, Car2Go.          (b) Milan, Enjoy.          (c) Rome, Car2go.

Figure 1: Stationarity check

Their behavior is mostly constant in time, meaning that the time series are stationary. Therefore, there is no need for differentiation. The ARIMA model's hyperparameter $d$ can be set equal to 0, and the model can be represented as ARMA(p, q).

# 4 Task4: ACF and PACF

In this step, the ACF and PACF are computed.[3] These are useful tools for finding adequate values of the $p$ and $q$ parameters if the model is pure AR(p), pure MA(q), or ARMA(p,q). This technique can only be used as a first attempt at creating a good model, further analysis is needed to choose optimal values of $p$ and $q$. Figures 2 and 3 display the ACF and PACF for the first 48 hours (i.e. 2 days) of the number of rentals time series for each city. The PACF provides an initial estimate of the hyper-parameter $p$ by counting the number of significant lags (i.e. those outside the defined blue zone with a significance level less than 5%). For example, for Milano Car2Go (figure 3a), we can count 6 significant peaks; statistically important lags that follow non-significant lags are not taken into account. Consequently, $p=6$ is chosen for Milano Car2Go

*if the process is AR only ...*

*if MA only ...*

and, similarly, $p=1$ for Milano Enjoy and $p=3$ for Roma Car2Go. The ACF graphs (figure 2) clearly show a periodicity: for each city, there is a regular repetition of the pattern in a 24-hour cycle. This result is consistent with the cyclical nature of the time series. A first guess of the parameter $q$ is best done with the correlogram of the ACF; again, it is necessary to check how many significant lags exist, taking care not to overestimate. In this case, $q=4$ can be assumed for all cities.
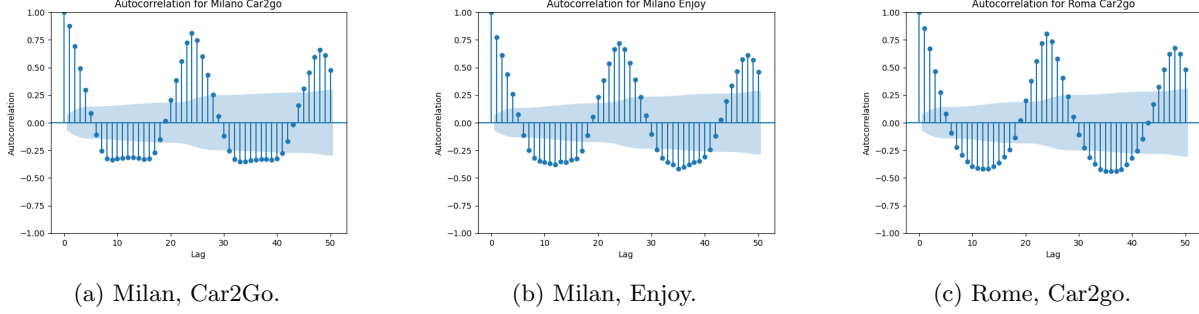


(a) Milan, Car2Go.　　　　(b) Milan, Enjoy.　　　　(c) Rome, Car2go.

Figure 2: ACF plots



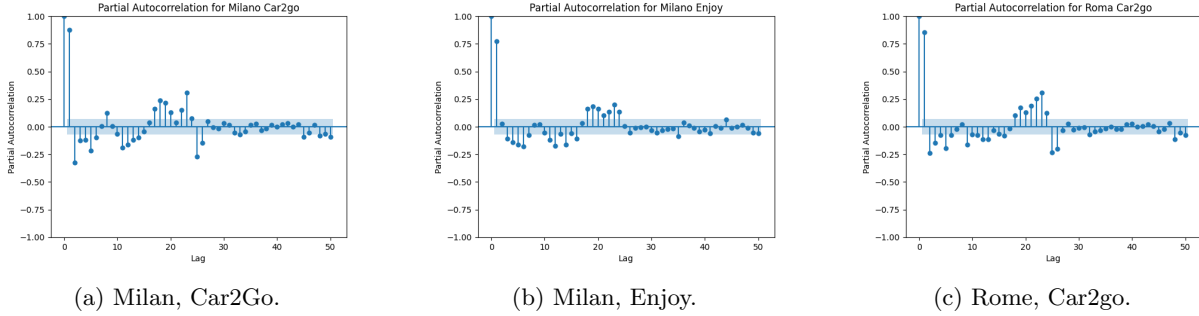(a) Milan, Car2Go.　　　　(b) Milan, Enjoy.　　　　(c) Rome, Car2go.

Figure 3: PACF plots

Since neither the ACF nor the PACF reaches zero after a certain value, all three models will be of ARMA(p,q) type.

# 5  Task5: Training the model

A total of $N=336$ training samples (a period of two weeks) are used, along with 120 hours of testing, to evaluate the performance of a model. To ensure that the model can best generalize from the provided samples, various hyper-parameters will be investigated. In the first set of tasks, an expansion window learning approach will be used to adapt the model to evolving patterns in the data during testing.

# 6  Task6: Model training and testing

To assess the performance of an ARIMA model on the three different cities, we define a single set of parameters (listed in Table 1) and use an expanding window learning strategy. The number of training and test samples coincides with those mentioned in the previous section.

| p | d | q |
|---|---|---|
| 2 | 0 | 4 |

Table 1: ARIMA model parameters: $q=4$ as guessed in the previous section; $p$, instead, is chosen as "average" among the $p$ values previously supposed, in order to determine a single set.

2

Table 2 displays the prediction errors for each city. [4] These errors are calculated by comparing the expected values to the actual ones. From the results, we can deduce that the model works better for Roma Car2Go.

| Case | MAE | MAPE | MSE | R2 |
|------|------|------|------|------|
| Milan, Car2Go | 51.310 | 17.887 | 3818.902 | 0.822 |
| Milan, Enjoy | 36.568 | 15.234 | 2043.024 | 0.820 |
| Rome, Car2Go | 24.223 | 18.314 | 952.847 | 0.788 |

Table 2: Error evaluation.

Moreover, when evaluating an ARIMA model, it is crucial to analyze the residuals, as they provide insight into the goodness of fit and may indicate the need for model adjustments. Figure 4a shows the resulting fitted model for Rome. The expected values appear to be fairly consistent with the actual values; no systematic overestimation or underestimation is visible, but the presence of discrepancies, although not too significant, certainly indicates that the model can be improved. The residual density curve in Figure 4b follows a broadly normal distribution with almost no mean, as ideally should be. [5]
Note that the model described bases the prediction of a value only on the previous two hours (since $p=2$); higher $p$ values could potentially provide more accurate results.
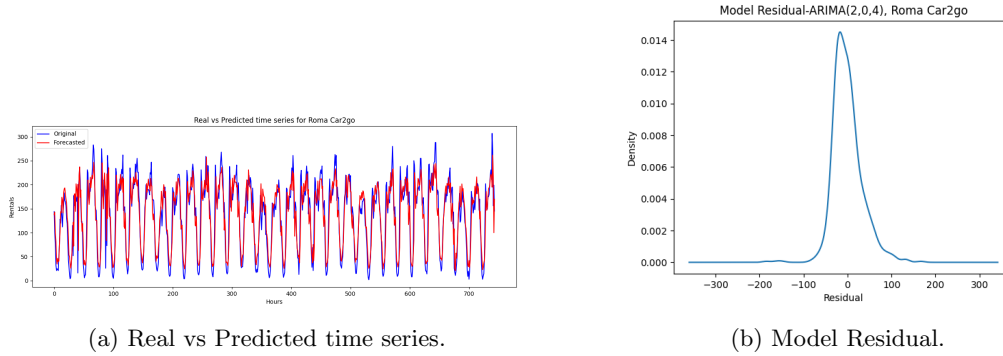


(a) Real vs Predicted time series.

(b) Model Residual.

Figure 4: Fitted model and residuals for Roma Car2Go.

# 7    Task7: Parameters definition
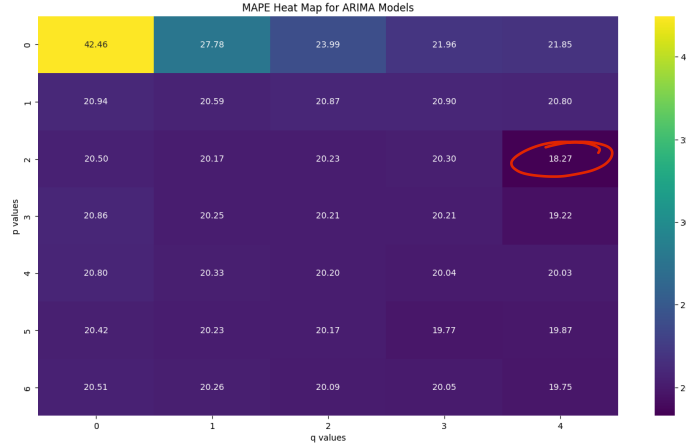
**7a. Grid search: tuning $p,q$.** [6]
The objective of this task is to search for the optimal values of hyperparameters $p$ and $q$ and study the performance of the ARIMA model for different combinations of these. We use an expansion learning strategy and maintain the same training and test dimensions as in the previous sections ($N=336$, test=120). To optimize the performance of the ARIMA model and ensure that it captures the characteristics of the time series at its best, we adopt the "grid search" parameter adjustment approach. The chosen ranges for the parameters are $p$:[0,6] and $q$:[0,4]. This allows us to observe if the values supposed in paragraph 4 can have some kind of match. The best model is identified based on performance evaluation metrics such as MAE and MAPE. Figure 5 shows in an easy-to-read heatmap the mean absolute percentage errors obtained for Roma Car2Go. The best combination for Roma Car2Go is $p = 2$ and $q = 4$, which is given by the lower MAPE. Similarly (see figures 11, 12 in the *References*), the best models for Milano Car2Go and Milano Enjoy are respectively ARIMA(3,0,4) and ARIMA(6,0,2). The results obtained differ from the hypothesis in section 4, which proves that a first intuition is not reliable.

**7b. Tuning $N$ and learning strategy.** [7]
After selecting the optimal configuration of parameters, it's necessary to check how the model performance is affected by the size of the training set and how it changes according to different learning strategies. To
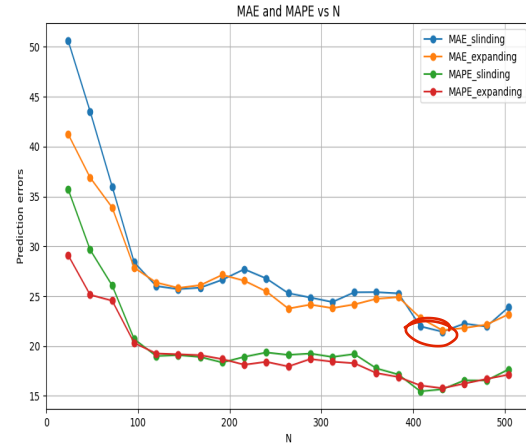
accomplish this, the number of test samples and the parameters $p$ and $q$ remain constant while $N$ varies from 24 to 504 (from 1 to 21 days). We calculate errors using both sliding window and expanding learning approaches. Table 3 displays the results obtained for the city of Rome, and the plot in figure 6 shows the behavior of the errors as $N$ increases. The model achieves the lowest errors with a training set size of 432 samples in combination with the sliding window learning strategy.



Figure 5: Roma (Car2Go), Heat map - MAPE for different combinations of parameters.

| N | MAE_sl | MAE_exp | MAPE_sl | MAPE_exp |
|---|---|---|---|---|
| 24 | 50.624 | 41.253 | 35.725 | 29.112 |
| 48 | 43.519 | 36.889 | 29.650 | 25.133 |
| 72 | 35.992 | 33.890 | 26.070 | 24.547 |
| 96 | 28.392 | 27.859 | 20.699 | 20.310 |
| 120 | 26.012 | 26.368 | 18.994 | 19.254 |
| 144 | 25.706 | 25.827 | 19.084 | 19.173 |
| 168 | 25.851 | 26.096 | 18.903 | 19.082 |
| 192 | 26.641 | 27.136 | 18.351 | 18.692 |
| 216 | 27.704 | 26.550 | 18.915 | 18.127 |
| 240 | 26.773 | 25.474 | 19.355 | 18.416 |
| 264 | 25.295 | 23.756 | 19.113 | 17.951 |
| 288 | 24.863 | 24.150 | 19.247 | 18.696 |
| 312 | 24.412 | 23.805 | 18.901 | 18.431 |
| 336 | 25.389 | 24.170 | 19.196 | 18.273 |
| 360 | 25.405 | 24.734 | 17.758 | 17.290 |
| 384 | 25.280 | 24.908 | 17.146 | 16.893 |
| 408 | 21.970 | 22.808 | 15.454 | 16.043 |
| 432 | 21.435 | 21.573 | 15.665 | 15.766 |
| 456 | 22.251 | 21.808 | 16.562 | 16.232 |
| 480 | 21.946 | 22.145 | 16.544 | 16.694 |
| 504 | 23.884 | 23.174 | 17.653 | 17.128 |

Table 3: Keeping p=2, d=0 and q=4 fixed, varying N for Roma Car2Go.



Figure 6: Errors vs $N$ for Roma Car2Go.

Similar results are observed for the other cities:

- Milano Car2Go: $N = 480$ , learning strategy: sliding window; (Table 5, figure 13)
- Milano Enjoy: $N = 240$, learning strategy: sliding window;(Table 6, figure 14)

### 7c. Comparison between cities.
After defining the best models for all three cities, we compare their performance (see Table 4). Milan Enjoy reveals the best prediction accuracy (lowest MAPE and moderate MAE) with fewer training samples. However, choosing the "best" model may depend on multiple factors such as computational efficiency, number

of training samples, and accuracy metrics. If efficiency is a top priority, since MAE and MAPE performance are acceptable, Rome Car2Go with ARIMA(2,0,4) may be the most suitable model. If accuracy is paramount, Milan Enjoy may be preferred. In the *References* there are the plots of the "fitted" models for the three cities. [8]

| | N | Model | MAPE_Sliding | MAE_Sliding |
|---|---|---|---|---|
| Roma Car2Go | 432 | ARIMA(2,0,4) | 15.662 | 21.435 |
| Milano Car2Go | 480 | ARIMA(3,0,4) | 15.667 | 45.876 |
| Milano Enjoy | 240 | ARIMA(6,0,2) | 13.110 | 33.191 |

Table 4: Results summary.
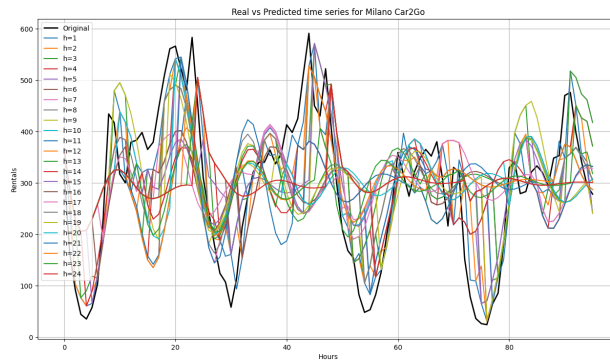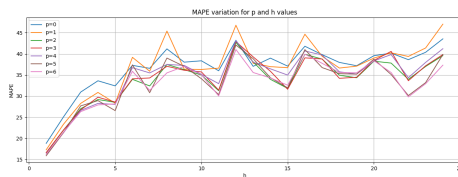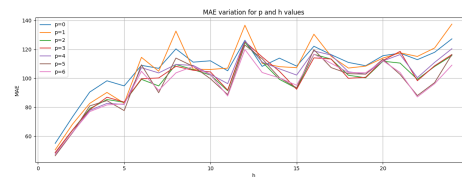
# 8  Task8: Impact of the time horizon



Figure 7: Forecast data for different time horizon values.

For this last task, the impact of the time horizon $h$ on the performance of the ARIMA model was investigated. Therefore, instead of only predicting one single value per time step, we increased the value of $h$, forecasting all the values between $t$ and $t+1$. Indeed, an ARIMA model ($p = 3, d = 0, q = 4$) was trained on a **train set** of 288 samples and tested on a **test set** of 96 samples (i.e. four days), all the used data were referred to **Milan, Car2Go**.[9] Figure 7 shows the results achieved for $h \in [1, 24]$.
Furthermore, both the MAPE and the MAE were evaluated for different values of $h$ and $p$.[10] Figure 8 depicts the results.



(a) MAPE.                                                   (b) MAE.

Figure 8: MAPE and MAE for different values of $p$ and $h$.

In general, the forecast values are closer to the original ones for lower values of $h$, as it can be seen in figure 8. Furthermore, it is clear that the best performance for the model is obtained for $p = 3$, which goes to prove the results stated in task 7.

5

# References

## [1] Task1 and Task2

```python
import pymongo as pm
import matplotlib.pyplot as plt
from datetime import datetime
import pymongo as pm
import pandas as pd
import numpy as np
from statsmodels.graphics.tsaplots import plot_acf, acf, pacf, plot_pacf
from statsmodels.tsa.arima.model import ARIMA
from sklearn.metrics import mean_squared_error
from sklearn.metrics import mean_absolute_error
from sklearn.metrics import r2_score
import seaborn as sns

def bookings_per_hour(db, service, city,time_start, time_end):
    if service == 'Car2go':
        Bookings = db['PermanentBookings']
    elif service == 'Enjoy':
        Bookings = db['enjoy_PermanentBookings']
    else:
        raise Exception('No service found!')

    start_unix = time_start.timestamp()
    end_unix = time_end.timestamp()

    res = Bookings.aggregate([
    # filtro citt  e periodo
    {'$match': {'city': city, 'init_time': {'$gte': start_unix, '$lte': end_unix}}},
    #nuovo format con data e duration
     {"$project": {
       "_id": 0,
       "date": {"$dateToString": {"format": "%Y-%m-%d %H:00:00", "date":"$init_date"}},
       "duration": {"$divide": [{"$subtract": ["$final_time",
       "$init_time"]}, 60]}
       }
    #filtri lab 1
       },
    {"$match": {"duration": {"$gte": 5, "$lte": 120},
                "final_address": {"$ne": "$init_address"}}},
    #nuovo format
    {"$group": {"_id": { "date": "$date"},
    "number_cars": {"$sum": 1}}
  },
  {"$project": {
    "_id": 1,
    "number_cars": 1
}},
{"$sort": {"_id.date": 1}},
])

    df = pd.DataFrame(list(res))
    #plt.figure()
    #df.number_cars.plot(title="Rentals in {}, {} before fitting".format(city,service))
    #t_start = pd.to_datetime("2017-10-01 00:00:00")
    #t_end = pd.to_datetime("2017-10-31 23:00:00")

    return df

def replace_with_hourly_average(df, complete_time_range):
    df["_id.date"] = pd.to_datetime(df["_id.date"])
    df = df.set_index("_id.date")
    #calcolo media per ogni orario del giorno
```

```
62    hourly_avg_df = df.groupby(df.index.hour)["number_cars"].mean().reset_index()
63    for time in complete_time_range:
64        #cerco i missing e sostituisco con la media
65        if time not in df.index:
66            hour_avg = hourly_avg_df[hourly_avg_df["_id.date"] == time.hour]["number_cars
      "].values[0]
67            df.loc[time] = hour_avg
68    df = df.sort_index()
69    return df
70
71 def handle_missing_samples(service, city, df, start_date, end_date, freq="1H"):
72    df["_id.date"] = pd.to_datetime(df["_id"].apply(lambda x: x["date"]))
73    complete_time_range = pd.date_range(start=start_date, end=end_date, freq=freq)
74    days = int(end_date.split("-")[2][:2])
75
76    if len(df)!=days*24:
77        df = replace_with_hourly_average(df, complete_time_range)
78
79    return df
80
81 if __name__ == '__main__':
82    client = pm.MongoClient('bigdatadb.polito.it',
83                            ssl=True,
84                            authSource='carsharing',
85                            username='ictts',
86                            password='Ict4SM22!',
87                            tlsAllowInvalidCertificates=True)
88    db = client['carsharing']
89    data_types = ['bookings']
90    city_names = ['Milano', 'Roma']
91    service_names = ['Car2go', 'Enjoy']
92
93    time_start = datetime(2017, 10, 1, 0, 0, 0)  # ottobre 1
94    time_end = datetime(2017, 10, 31, 23, 00, 00)  # ottobre 30
95    for city in city_names:
96        for service in service_names:
97            df = bookings_per_hour(db, service, city, time_start, time_end)
98            new_df = handle_missing_samples(service, city, df,\
99                                            time_start.strftime("%Y-%m-%d %H:%M:%S"),
      time_end.strftime("%Y-%m-%d %H:%M:%S"))
```

**Task3**

```
[2] def check_stationarity_mean(fitted_df, window_size):
2      fitted_df = fitted_df.set_index("_id")
3      fitted_df['number_cars'] = pd.to_numeric(fitted_df['number_cars'])
4      fitted_df['rolling_avg'] = fitted_df['number_cars'].rolling(window=window_size).mean
      ()
5      f_df = fitted_df.dropna()
6      dropped_values = len(fitted_df) - len(f_df)
7
8      return f_df, dropped_values
9
10 def check_stationarity_std(fitted_df, window_size):
11     fitted_df = fitted_df.set_index("_id")
12     fitted_df['number_cars'] = pd.to_numeric(fitted_df['number_cars'])
13     fitted_df['rolling_std'] = fitted_df['number_cars'].rolling(window=window_size).std()
14     f_df = fitted_df.dropna()
15     dropped_values = len(fitted_df) - len(f_df)
16
17     return f_df,dropped_values
18
19 if __name__ == '__main__':
20     client = pm.MongoClient('bigdatadb.polito.it',
21                             ssl=True,
```

```
22                              authSource='carsharing',
23                              username='ictts',
24                              password='Ict4SM22!',
25                              tlsAllowInvalidCertificates=True)
26      db = client['carsharing']
27      data_types = ['bookings']
28      city_names = ['Milano', 'Roma']
29      service_names = ['Car2go', 'Enjoy']
30
31      time_start = datetime(2017, 10, 1, 0, 0, 0)  # October 1
32      time_end = datetime(2017, 10, 31, 23, 00, 00)  # October 30
33      window_size=168
34      city_df={}
35
36      for city in city_names:
37        for service in service_names:
38          df = bookings_per_hour(db, service, city, time_start, time_end)
39          new_df = handle_missing_samples(service, city, df,
40                                      time_start.strftime("%Y-%m-%d %H:%M:%S"),
41                                      time_end.strftime("%Y-%m-%d %H:%M:%S"))
42          rolling_mean_df, dropped_values_1 = check_stationarity_mean(new_df, window_size)
43          rolling_mean_df['index_column'] = range(dropped_values_1, dropped_values_1 + len(
      rolling_mean_df))
44          rolling_std_df, dropped_values_2 = check_stationarity_std(new_df, window_size)
45          rolling_std_df['index_column'] = range(int(dropped_values_2), int(
      dropped_values_2) + len(rolling_std_df))
46
47
48          # Reset the index for new_df just for plotting (old index was whole date )
49          new_df_reset = new_df.reset_index()
50          new_df_reset['index_column'] = range(len(new_df_reset))
51          city_df[(city, service)] = new_df_reset
52          plt.figure()
53          plt.plot(new_df_reset['index_column'], new_df_reset['number_cars'], label='Number
       of rentals')
54          plt.plot(rolling_mean_df['index_column'], rolling_mean_df['rolling_avg'], label='
      Rolling Average')
55          plt.plot(rolling_std_df['index_column'], rolling_std_df['rolling_std'], label='
      Rolling std')
56          plt.xlabel('Hours')
57          plt.ylabel('Rentals')
58          plt.title('Number of rentals and Rolling values, {},{}'.format(city, service))
59          plt.legend()
60          plt.show()
```

### Task4

```
[3] def acf_plot(df,n_lags,city,service):
 2      plot_acf(df,lags=n_lags)
 3      plt.xlabel('Lag')
 4      plt.ylabel('Autocorrelation')
 5      plt.title(f'Autocorrelation for {city} {service}')
 6      plt.show()
 7
 8  def pacf_plot(df,n_lags,city,service):
 9      plot_pacf(df, lags=n_lags, method='ols')
10      plt.xlabel('Lag')
11      plt.ylabel('Partial Autocorrelation')
12      plt.title(f'Partial Autocorrelation for {city} {service}')
13      plt.show()
14
15  #------------Alternative-------------------
16  def alternative_acf_plot(df,nlags,city,service):
17      acf_lags=acf(df,nlags = nlags) # hp) nlags= len(df)
18      plt.plot(acf_lags)
```

```
19        plt.axis([0,nlags,-1,1])
20        plt.axhline(y= 0,linestyle="--",color = 'gray')
21        plt.axhline(y = -1.96/np.sqrt(len(df)),linestyle = '--',color = 'gray') #95%
          confidentiality
22        plt.axhline(y = 1.96/np.sqrt(len(df)),linestyle = '--',color = 'gray')
23        plt.title(f'Autocorrelation for {city} {service}')
24        plt.xlabel('Lag')
25        plt.ylabel('Autocorrelation')
26        plt.show()
27
28 def alternative_pacf_plot(df, nlags,city,service):
29        pacf_lags=pacf(df,nlags = nlags) # hp) nlags= len(df)
30        plt.plot(pacf_lags)
31        plt.axis([0,nlags,-.5,1])
32        plt.axhline(y= 0,linestyle="--",color = 'gray')
33        plt.axhline(y = -1.96/np.sqrt(len(df)),linestyle = '--',color = 'gray')
34        plt.axhline(y = 1.96/np.sqrt(len(df)),linestyle = '--',color = 'gray')
35        plt.title(f'Partial Autocorrelation for {city} {service}')
36        plt.xlabel('Lag')
37        plt.ylabel('Partial Autocorrelation')
38        plt.show()
39
40 for (city, service), df_reset in city_df.items():
41        if city=='Roma' and service=='Enjoy':
42            pass
43        else:
44            # df should represent a time series with a time-type index
45            df= df_reset[['index_column','number_cars']] #consider just these 2 columns (
       time=index +rentals)
46            df = df.set_index('index_column') # time=index
47            #print(f"DataFrame for {city}, {service}:\n{df}")
48            acf_plot(df,48,city,service)
49            alternative_acf_plot(df,48,city,service)
50            pacf_plot(df,48, city, service)
51            alternative_pacf_plot(df,48,city,service)
```
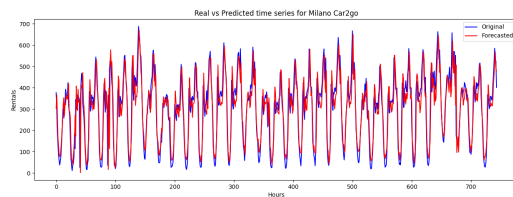
**Task6**

```
[4] p=2
 2 d=0
 3 q=4
 4 for (city, service), df_reset in city_df.items():
 5        key = f'{city}_{service}'
 6        if city=='Roma' and service=='Enjoy':
 7            pass
 8        else:
 9            df= df_reset[['index_column','number_cars']]
10            df = df.set_index('index_column')
11            X = df.values.astype(float)
12            train_len = 216
13            print("Training on ", train_len)
14            test_len = 120
15            predictions = np.zeros( (test_len) )
16            print('Testing ARIMA order (%i,%i,%i) for %s' % (p,d,q,key))
17            train, test = X[0:train_len], X[train_len:(train_len+test_len)] # split the
       train and test - we use an expanding window
18                                                              # approach - where data
       from 0 to now is used to train
19                                                              # and then predict now+1.
        We then shift the time and repeat
20            history = [x for x in train] # this is the past data used for training
21            for t in range(0,test_len): # repeat for all tests we want to do
22                model=ARIMA(history, order=(p,d,q))
23                model_fit = model.fit(method_kwargs={'maxiter':300})#
24                output = model_fit.forecast() # here we get the forecasted data
```

9
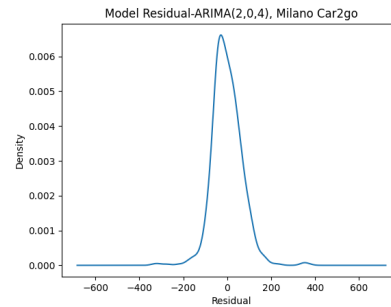
```
25                    yhat = output[0] # get now+1
26                    predictions[t]=yhat #and store it
27                    obs = test[t] # slide over time, by putting now+1 into past
28                    history.append(obs)
29
30            print('(%i,%i,%i) model => MAE: %.3f -- MAPE: %.3f -- MSE: %.3f -- R2: %.3f'
     % (p,d,q, mean_absolute_error(test, predictions),
31                                            mean_absolute_error(test, predictions) / test.
     mean()*100,
32                                            mean_squared_error(test, predictions),
33                                            r2_score(test, predictions)))
```

```
1  for (city, service), df_reset in city_df.items():
2        key = f'{city}_{service}'
3        if city=='Roma' and service=='Enjoy':
4            pass
5        else:
6            df= df_reset[['index_column','number_cars']]
7            df = df.set_index('index_column')
8            model=ARIMA(df, order=(p,d,q))
9            model_fit = model.fit(method_kwargs={'maxiter':300}) # "statespace","
     innovations_mle","hannan_rissanen"
10           #plot results
11           fig = plt.figure(figsize=(15,5))
12           plt.plot(df, color='blue', label='Original')
13           plt.plot(model_fit.fittedvalues, color='red', label='Forecasted')
14           plt.xlabel('Hours')
15           plt.ylabel('Rentals')
16           plt.title('Real vs Predicted time series for {} {}'.format(city,service))
17           plt.legend()
18           plt.show()
19           residuals = pd.DataFrame(model_fit.resid)
20           residuals.plot(kind='kde',  legend=False)
21           plt.xlabel('Residual')
22           plt.ylabel('Density')
23           plt.title('Model Residual-ARIMA({},{},{}), {} {}'. format(p,d,q,city,service)
     )
24           plt.show()
```
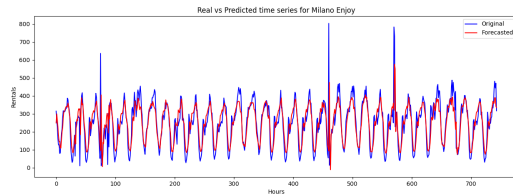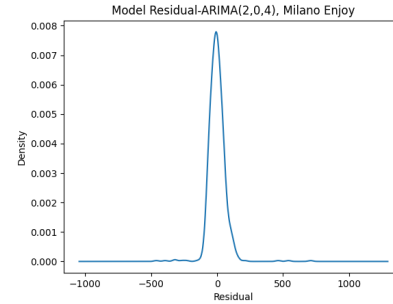


(a) Real vs Predicted time series.



(b) Model Residual.

Figure 9: Fitted model and residuals for Milano Car2Go.

(a) Real vs Predicted time series.
(b) Model Residual.

Figure 10: Fitted model and residuals for Milano Enjoy.

## Task7

```
1  #--------------------Task 7a--------------------------
2  #Heatmap Plot Function
3  def heatmap_forARIMA(matrix,matrix_name):
4      plt.figure(figsize=(12, 8))
5      sns.heatmap(matrix, annot=True, fmt=".2f", cmap="viridis",
6                  xticklabels=diff_degrees, yticklabels=lag_orders)
7      plt.title(f'{matrix_name} Heat Map for ARIMA Models')
8      plt.xlabel('q values')
9      plt.ylabel('p values')
10     plt.show()
11 #----------------------------------------------------
12 train_len = 24*7*2 #2 weeks
13 test_len = 24*5 #5 days
14 d=0
15 lag_orders = (0,1,2,3,4,5,6) #values of p
16 diff_degrees = (0,1,2,3,4) #values of q
17
18 for (city, service), df_reset in city_df.items():
19     key = f'{city}_{service}'
20     if city=='Roma' and service=='Enjoy':
21         pass
22     else:
23         df= df_reset[['index_column','number_cars']]
24         df = df.set_index('index_column')
25         X = df['number_cars'].values.astype(float)
26         predictions = np.zeros((len(lag_orders), len(diff_degrees), test_len))
27         train, test = X[0:train_len], X[train_len:(train_len + test_len)]
28         mape_matrix = np.zeros((len(lag_orders), len(diff_degrees)))
29         mae_matrix = np.zeros((len(lag_orders), len(diff_degrees)))
30         for i, p in enumerate(lag_orders):
31             for j, q in enumerate(diff_degrees):
32                 print(f'Case {key}')
33                 print(f'Testing ARIMA order ({p}, {d}, {q})')
34                 history = [x for x in train]
35                 for t in range(0, test_len):
36                     model = ARIMA(history, order=(p, d, q))
37                     model.initialize_approximate_diffuse()
38                     model_fit = model.fit()
39                     output = model_fit.forecast()
40                     yhat = output[0]
41                     predictions[i][j][t] = yhat
42                     obs = test[t]
43                     history.append(obs)
44
45                 # Calcolo del MAPE
```

11

```
46             mape = mean_absolute_error(test, predictions[i][j]) / test.mean() * 100
47             mape_matrix[i, j] = mape
48             mae = mean_absolute_error(test, predictions[i][j])
49             mae_matrix[i,j]=mae
50
51             print('(%i,%i,%i) model => MAE: %.3f -- MAPE: %.3f' % (p, d, q, mae, mape
    ))
52
53     heatmap_forARIMA(mape_matrix, 'MAPE')
54     heatmap_forARIMA(mae_matrix, 'MAE')
```
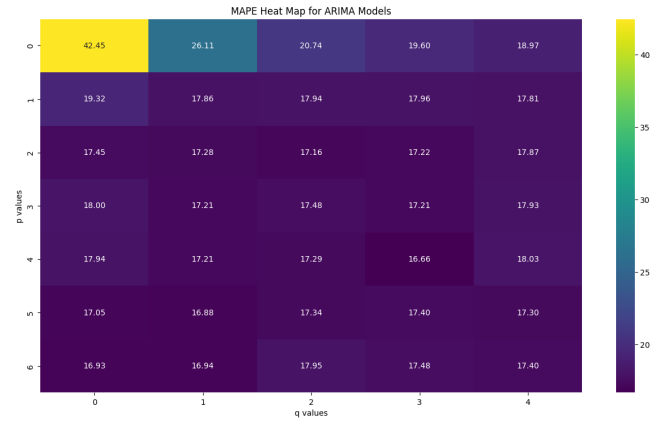


Figure 11: Milano (Car2Go),Heat map - MAPE for different combinations of parameters.



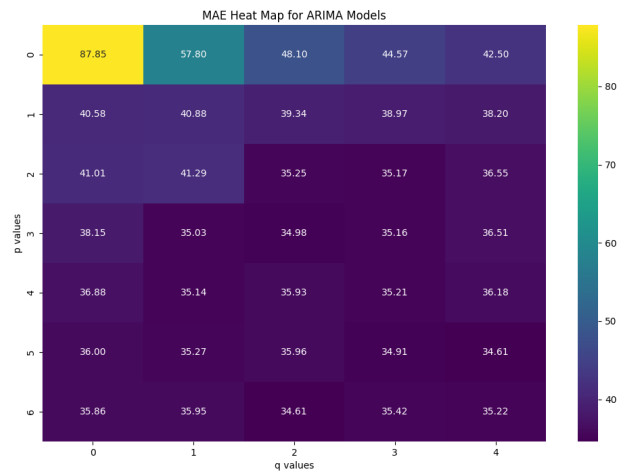Figure 12: Milano (Enjoy),Heat map - MAPE for different combinations of parameters.

```
[7]  #-------------------Task 7b--------------------------------
2   test_len=24*5 #5 days
3   train_len_list= 24 * np.linspace(1, 21, 21) # from 1 to 21 days
```

```python
4  lag_orders=(0,1,2,3,4,5,6)
5  d=0
6  # ###############Roma Car2Go parameters##################
7  p=2
8  q=4
9  df=city_df[("Roma","Car2Go")]
10 # ###############Milano Car2Go parameters##################
11 # p=3                                                    #
12 # q=4                                                    #
13 # df=city_df[("Milano","Car2Go")]                       #
14 # ###########################################################
15
16 # ###############Milano Enjoy parameters##################
17 # p=6                                                    #
18 # q=2                                                    #
19 # df=city_df[("Milano","Enjoy")]                        #
20 # ###########################################################
21
22 predictions_sl=np.zeros((len(lag_orders), test_len))
23 predictions_exp=np.zeros((len(lag_orders), test_len))
24 df = df.set_index('index_column')
25 tot_stats=pd.DataFrame()
26 X = df['number_cars'].values.astype(float) # extract the time series
27 for N in train_len_list:
28     N=int(N)
29     train, test = X[0:N], X[N:(N + test_len)]
30     print(f'Testing N equal to {N}')
31     history_sl = [x for x in train]
32     history_exp = [x for x in train]
33     for t in range(0, test_len):
34         # SLIDING
35         model_sl= ARIMA(history_sl, order=(p, d, q))
36         model_sl.initialize_approximate_diffuse()
37         model_fit_sl = model_sl.fit()
38         output_sl = model_fit_sl.forecast()
39         yhat_sl = output_sl[0]
40         predictions_sl[lag_orders.index(p)][t] = yhat_sl
41         # EXPANDING
42         model_exp = ARIMA(history_exp, order=(p, d, q))
43         model_exp.initialize_approximate_diffuse()
44         model_fit_exp = model_exp.fit()
45         output_exp = model_fit_exp.forecast()
46         yhat_exp = output_exp[0]
47         predictions_exp[lag_orders.index(p)][t] = yhat_exp
48
49         obs = test[t]
50         history_sl.append(obs)
51         history_sl = history_sl[1:]
52         history_exp.append(obs)
53
54     mae_sl=mean_absolute_error(test, predictions_sl[lag_orders.index(p)])
55     mape_sl=mean_absolute_error(test, predictions_sl[lag_orders.index(p)]) / test.mean()
       *100
56     mse_sl=mean_squared_error(test, predictions_sl[lag_orders.index(p)])
57     r2_sl=r2_score(test, predictions_sl[lag_orders.index(p)])
58     print('(%i,%i,%i) model_sliding => MAE: %.3f -- MAPE: %.3f -- MSE: %.3f -- R2: %.3f'
       % (p,d,q, mae_sl,mape_sl,mse_sl,r2_sl))
59
60
61     mae_exp=mean_absolute_error(test, predictions_exp[lag_orders.index(p)])
62     mape_exp=mean_absolute_error(test, predictions_exp[lag_orders.index(p)]) / test.mean
       ()*100
63     mse_exp=mean_squared_error(test, predictions_exp[lag_orders.index(p)])
64     r2_exp=r2_score(test, predictions_exp[lag_orders.index(p)])
65     print('(%i,%i,%i) model_expanding => MAE: %.3f -- MAPE: %.3f -- MSE: %.3f
```

```
           ' % (p,d,q, mae_exp,mape_exp,mse_exp,r2_exp))
66
67      statistics={"N":N, "MAE_sl":mae_sl, "MAE_exp":mae_exp, "MAPE_sl":mape_sl,"MAPE_exp":
        mape_exp,"MSE_sl":mse_sl, "MSE_exp": mse_exp, "R2_sl":r2_sl,"R2_exp":r2_exp }
68      statistics_df=pd.DataFrame(statistics, index=[0])
69      tot_stats=pd.concat([tot_stats, statistics_df], ignore_index=True)
70
71  print(tot_stats)
72  #------------------------------------------------------------
73  #Plot prediction errors vs N
74  N=tot_stats['N'].to_list()
75  MAE_sl=tot_stats['MAE_sl'].to_list()
76  MAE_exp=tot_stats['MAE_exp'].to_list()
77  MAPE_sl=tot_stats['MAPE_sl'].to_list()
78  MAPE_exp=tot_stats['MAPE_exp'].to_list()
79  plt.figure(figsize=(10, 6))
80  plt.plot(N, MAE_sl, marker='o', label='MAE_slinding')
81  plt.plot(N, MAE_exp, marker='o', label='MAE_expanding')
82  plt.plot(N, MAPE_sl, marker='o', label='MAPE_slinding')
83  plt.plot(N, MAPE_exp, marker='o', label='MAPE_expanding')
84  plt.xlabel('N')
85  plt.ylabel('Prediction errors')
86  plt.title('MAE and MAPE vs N')
87  plt.legend()
88  plt.grid(True)
89  plt.show()
```

| N | MAE_sl | MAE_exp | MAPE_sl | MAPE_exp |
|---|---|---|---|---|
| 24 | 91.730 | 76.069 | 31.129 | 25.814 |
| 48 | 71.462 | 67.988 | 22.778 | 21.671 |
| 72 | 61.142 | 63.385 | 20.076 | 20.812 |
| 96 | 52.672 | 54.667 | 17.604 | 18.271 |
| 120 | 52.720 | 54.633 | 17.824 | 18.471 |
| 144 | 50.482 | 49.449 | 17.990 | 17.622 |
| 168 | 48.600 | 48.254 | 17.385 | 17.261 |
| 192 | 51.764 | 52.188 | 17.315 | 17.457 |
| 216 | 51.247 | 52.097 | 16.281 | 16.551 |
| 240 | 50.846 | 50.454 | 16.497 | 16.370 |
| 264 | 48.393 | 48.696 | 16.191 | 16.292 |
| 288 | 47.714 | 48.545 | 16.337 | 16.621 |
| 312 | 47.679 | 48.466 | 16.671 | 16.947 |
| 336 | 48.732 | 51.432 | 16.988 | 17.929 |
| 360 | 52.038 | 54.468 | 17.072 | 17.869 |
| 384 | 53.857 | 56.172 | 16.773 | 17.494 |
| 408 | 53.293 | 54.443 | 16.932 | 17.298 |
| 432 | 50.842 | 51.771 | 16.589 | 16.892 |
| 456 | 47.667 | 48.363 | 16.047 | 16.281 |
| 480 | 45.876 | 47.577 | 15.667 | 16.247 |
| 504 | 47.565 | 47.610 | 16.326 | 16.341 |

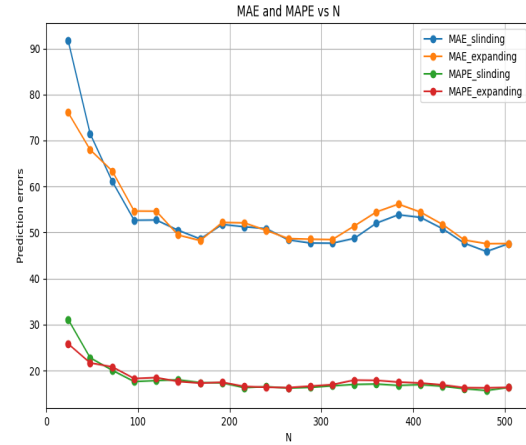Table 5: Keeping p=3, d=0 and q=4 fixed, varying N for Milano Car2Go.



Figure 13: Errors vs $N$ for Milano Car2Go.

| N | MAE_sl | MAE_exp | MAPE_sl | MAPE_exp |
|---|--------|---------|---------|----------|
| 24 | 106.138 | 72.211 | 46.206 | 31.436 |
| 48 | 67.288 | 64.240 | 27.659 | 26.407 |
| 72 | 60.424 | 61.958 | 24.243 | 24.859 |
| 96 | 48.794 | 47.968 | 19.469 | 19.139 |
| 120 | 46.424 | 45.649 | 18.506 | 18.197 |
| 144 | 37.154 | 37.749 | 15.568 | 15.817 |
| 168 | 34.124 | 35.999 | 14.970 | 15.793 |
| 192 | 36.737 | 39.914 | 15.941 | 17.319 |
| 216 | 34.826 | 38.506 | 14.252 | 15.758 |
| 240 | 33.191 | 35.829 | 13.110 | 14.152 |
| 264 | 34.383 | 36.078 | 13.517 | 14.183 |
| 288 | 34.444 | 37.069 | 13.530 | 14.562 |
| 312 | 34.969 | 35.664 | 13.966 | 14.243 |
| 336 | 33.744 | 34.606 | 14.057 | 14.417 |
| 360 | 56.894 | 50.735 | 23.630 | 21.071 |
| 384 | 52.132 | 50.404 | 20.651 | 19.967 |
| 408 | 49.247 | 49.143 | 18.643 | 18.603 |
| 432 | 48.925 | 48.544 | 18.522 | 18.378 |
| 456 | 57.380 | 58.503 | 21.426 | 21.846 |
| 480 | 43.202 | 43.453 | 16.535 | 16.631 |
| 504 | 43.451 | 43.482 | 17.153 | 17.166 |

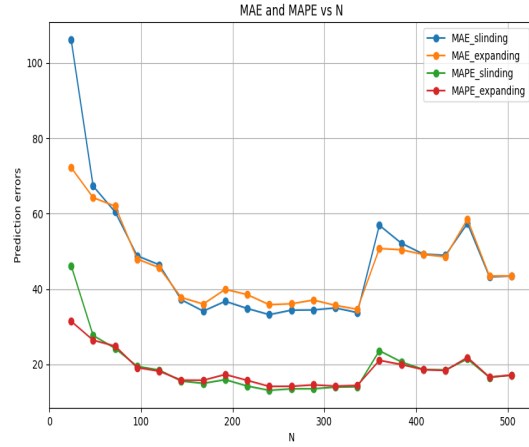Table 6: Keeping p=6, d=0 and q=2 fixed, varying N for Milano Enjoy.
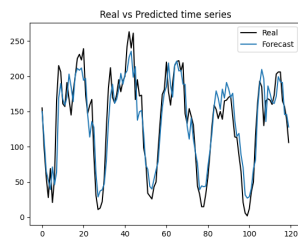


Figure 14: Errors vs $N$ for Milano Enjoy.

```python
# ------------------------ Task 7c ----------------------------
def fit_plot(test,predictions):
    plt.figure()
    plt.plot(test, color='black', label='Real') # plot the real time series
    plt.plot(predictions, label='Forecast')
    plt.title("Real vs Predicted time series")
    plt.legend()
    plt.show()

def buildARIMA_model_sl(p,d,q,N,X,predictions):
    train, test = X[0:N], X[N:(N+ test_len)]
    print(f'Testing N equal to {N}')
    history_sl = [x for x in train]
    for t in range(0, test_len):
        # SLIDING
        model_sl= ARIMA(history_sl, order=(p, d, q))
        model_sl.initialize_approximate_diffuse()
        model_fit_sl = model_sl.fit()
        output_sl = model_fit_sl.forecast()
        yhat_sl = output_sl[0]
        predictions[t] = yhat_sl
        obs = test[t]
        history_sl.append(obs)
        history_sl = history_sl[1:]

    fit_plot(test,predictions)

test_len = 24 * 5
d = 0
#Roma Car2Go
p_Roma = 2
q_Roma = 4
N_Roma = 432
strategy_Roma= "sliding"
df_Roma=city_df[("Roma","Car2go")]
X_Roma = df_Roma['number_cars'].values.astype(float)
predictions_Roma = np.zeros(test_len)
buildARIMA_model_sl(p_Roma,d,q_Roma,N_Roma,X_Roma,predictions_Roma)
#Milano Car2Go
```
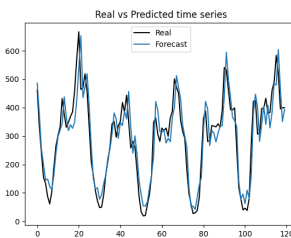
15

```
40  p_MilanoC = 3
41  q_MilanoC = 4
42  N_MilanoC = 480
43  strategy_MilanoC= "sliding"
44  df_MilanoC=city_df[("Milano","Car2go")]
45  X_MilanoC = df_MilanoC['number_cars'].values.astype(float)
46  predictions_MilanoC = np.zeros(test_len)
47  buildARIMA_model_sl(p_MilanoC,d,q_MilanoC,N_MilanoC,X_MilanoC,predictions_MilanoC)
48  #Milano Enjoy
49  p_MilanoE = 6
50  q_MilanoE = 2
51  N_MilanoE = 240
52  strategy_MilanoE= "sliding"
53  df_MilanoE=city_df[("Milano","Enjoy")]
54  X_MilanoE = df_MilanoE['number_cars'].values.astype(float)
55  predictions_MilanoE = np.zeros(test_len)
56  buildARIMA_model_sl(p_MilanoE,d,q_MilanoE,N_MilanoE,X_MilanoE,predictions_MilanoE)
```
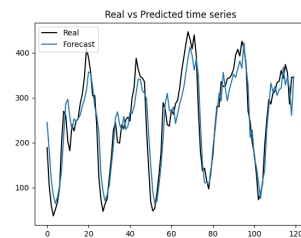


(a) Roma, Car2Go.  (b) Milano, Car2Go.  (c) Milano, Enjoy

Figure 15: Real vs Forecast time series.

**Optional task**

```
[9]  city = 'Milano'
 2   service = 'Car2Go'
 3   # Best parameters for Milano Car2Go -> Found after a grid search
 4   p = 3
 5   q = 4
 6   d = 0
 7
 8   train_len = 288
 9   test_len = 96
10
11   # Forecast step values
12   h_list = list(range(1,25))
13   final_res = []
14   df_MilanoC=city_df[("Milano","Car2go")]
15   df_MilanoC = df_MilanoC.set_index('index_column')
16   X = df_MilanoC['number_cars'].values.astype(float)
17
18   print(f'Testing ARIMA order ({p}, {d}, {q}) on {city}, {service}.')
19   train, test = X[0:train_len], X[train_len:(train_len + test_len)]
20
21   for h in h_list:
22       predictions = np.zeros(test_len)
23       history = [x for x in train]
24       t = 0
25       while t < test_len:
26           model = ARIMA(history, order=(p,d,q))
27           model_fit = model.fit()
28           # Various values for step h
```

```
29          output = model_fit.forecast(steps=h)
30          yhat = output[0:h]
31          predictions[t:t+h] = yhat[:len(predictions[t:t+h])]
32          obs = test[t:t+h]
33          history = history[h:]
34          history = history + list(obs)
35          t += h
36
37      final_res.append(predictions)
38
39 # plot results
40 fig = plt.figure(figsize=(15,5))
41 plt.plot(range(test_len), test, color='black', label='Original', linewidth=2)
42 for h in range(1,25):
43      plt.plot(range(test_len), np.array(final_res[h-1]).flatten(), label=f'h={h}')
44 plt.xlabel('Hours')
45 plt.ylabel('Rentals')
46 plt.title('Real vs Predicted time series for {} {}'.format(city,service))
47 plt.legend()
48 plt.grid()
49 plt.show()
```

```
[10] # Compute MAPE and MAE for different values of p and h
 2 city = 'Milano'
 3 service = 'Car2Go'
 4 # Best parameters for Milano Car2Go -> Found after a grid search
 5 p = 3
 6 q = 4
 7 d = 0
 8 # Range of p values
 9 Np = 7
10
11 train_len = 288
12 test_len = 96
13
14 # Forecast step values
15 h_list = list(range(1,25))
16
17 mape_p = []
18 mae_p = []
19
20 df_MilanoC=city_df[("Milano","Car2go")]
21 df_MilanoC = df_MilanoC.set_index('index_column')
22 X = df_MilanoC['number_cars'].values.astype(float)
23
24 train, test = X[0:train_len], X[train_len:(train_len + test_len)]
25
26 for p in range(Np):
27   mape_h = []
28   mae_h = []
29   print(f'Testing ARIMA order ({p}, {d}, {q}) on {city}, {service}.')
30
31   for h in h_list:
32     predictions = np.zeros(test_len)
33     history = [x for x in train]
34     t = 0
35     while t < test_len:
36       model = ARIMA(history, order=(p,d,q))
37       model_fit = model.fit()
38       # Various values for step h
39       output = model_fit.forecast(steps=h)
40       yhat = output[0:h]
41       predictions[t:t+h] = yhat[:len(predictions[t:t+h])]
42       obs = test[t:t+h]
```

```python
43        history = history[h:]
44        history = history + list(obs)
45        t += h
46
47      mape = mean_absolute_error(test, predictions) / test.mean() * 100
48      mae = mean_absolute_error(test, predictions)
49      mape_h.append(mape)
50      mae_h.append(mae)
51
52    mape_p.append(mape_h)
53    mae_p.append(mae_h)
54
55  # Plot MAPE values
56  fig = plt.figure(figsize=(15,5))
57  for p in range(Np):
58    plt.plot(range(1,25), mape_p[p], label=f'p={p}')
59  plt.xlabel('h')
60  plt.ylabel('MAPE')
61  plt.title('MAPE variation for p and h values')
62  plt.legend()
63  plt.grid()
64  plt.show()
65
66  # Plot MAE values
67  fig = plt.figure(figsize=(15,5))
68  for p in range(Np):
69    plt.plot(range(1,25), mae_p[p], label=f'p={p}')
70  plt.xlabel('h')
71  plt.ylabel('MAE')
72  plt.title('MAE variation for p and h values')
73  plt.legend()
74  plt.grid()
75  plt.show()
```