

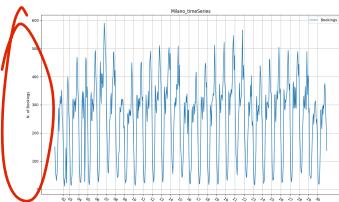
# ICT for smart mobility

Kevin Gumenia s320054, S.Ten. Federico Urso s316888, Alex Umberto Benedetti s319318

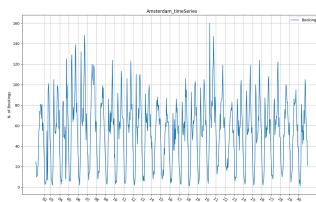
January 6, 2024

## 1 TASK 1: Prediction using ARIMA models

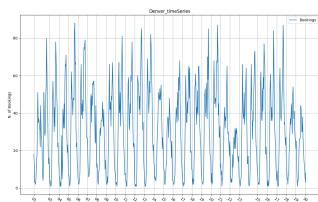
The first task of this step, was to extract the time series of rentals for each city (i.e., Amsterdam, Milano and Denver), considering a period of 30 days. We decided to modify our observation period with respect to the previous analysis because it was not suitable for our purpose. why? ✓  
We decided to consider the time period of 30 days from 1st of October to 30th of October. We choose this period because the period of December was not suitable since the last days of December may be characterized by strange results due to the Christmas holiday and in the period of November there are missing data in a lot of days. ✓



(a) Milano TimeSeries



(b) Amsterdam TimeSeries



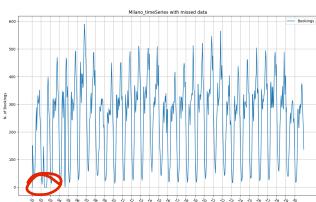
(c) Denver TimeSeries

Figure 1: Original TimeSeries

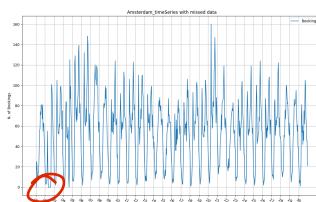
From figures 1a and 1b it is possible to notice that there are some data of some hours missed; i.e. on the 2nd of October and the day before due to a system error. Instead in Denver it is possible to see from the figure 1c are missed some hours between the 1st and the 3rd and between the 23rd and the 25th of October.

## 2 TASK 2: Clean Missing Values

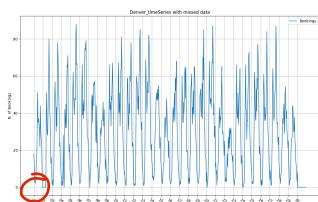
ARIMA models works only with regular time series without missed data , for this reason it will be important to fit the data.. To fit the data firstly we had identified the missing data and after we replaced those missing data with zero. As it is possible to see from figure 2a and figure 2b now the missing data before the 2nd of October have been replaced with zero values. In fact, some lines now are exactly in zero. This cleaning works also for Denver as it is possible to see from the figure 2c.



(a) Milano TimeSeries



(b) Amsterdam TimeSeries



(c) Denver TimeSeries

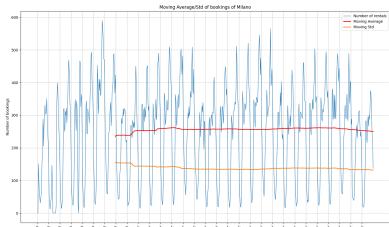
Figure 2: TimeSeries with missed data

From figures 2 it is possible to see that the cleaning process has been effective since it doesn't affect the periodicity.

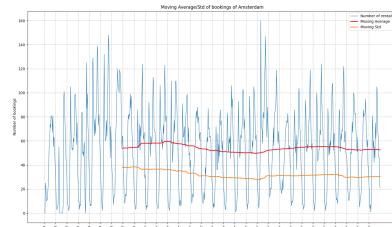
### 3 TASK 3: Stationarity Checking

over which time window ?

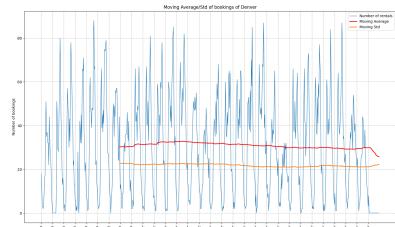
In order to check the stationarity or not of a time series it is possible to plot the rolling statistics: average and standard deviation over the time series that we have plotted in figures 1. A process can be defined stationary if the mean and standard deviation are constant over time.



(a) Milano Stationarity



(b) Amsterdam Stationarity



(c) Denver Stationarity

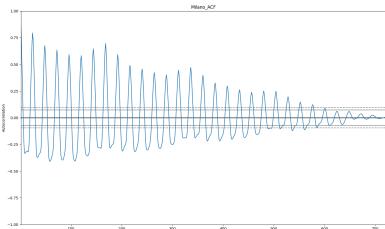
Figure 3: Stationarity

From the figure 3 is possible to understand that both average and standard deviation for all the cities are almost constant over the time. Therefore we deduced that the process is stationary. Due to this stationarity it is possible to set the parameter  $d=0$ . Hence we will have a model  $(p,0,q)$  which represent an ARMA model.

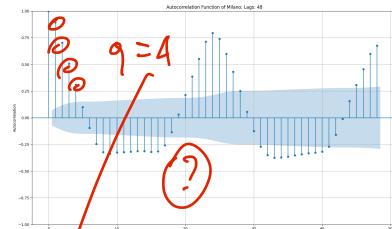
### 4 TASK 4: ACF and PACF

In this section we calculated the Autocorrelation function (**ACF**) and the Partial Autocorrelation Function (**PACF**) which is useful to choose an initial value for the hyperparameters of the model. We have three parameters :

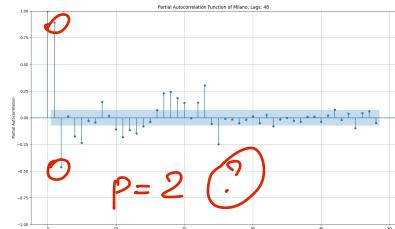
- **p**: The number of lag observations included in the model, also called the lag order.
- **d**: The number of times that the raw observations are differenced, also called the degree of differencing. In our case it is zero as we see before.
- **q**: The size of the moving average window, also called the order of moving average.



(a) Autocorrelation function



(b) ACF - zoom



(c) PACF - zoom

Figure 4: Milano autocorrelation and partial autocorrelation

From the ACF 4a it is possible to notice the 24 hours periodicity and also the weekly periodicity. Therefore as seen before  $d$  is zero, now we have to set the parameter  $p$  and  $q$  by looking at the figures. In fact looking at figure 4b it's possible to determine  $q=3$  Instead looking at figure 4c it's possible to determine  $p=1$  why ?

The chosen method to select the values for  $p$  and  $q$  is based on using the ACF and the PACF's plots. The method to select the best value for  $P$  follows these steps: Consider the PACF's plot, ignore the 0-th one and count the first  $n$ -elements until the closer lag to the dense part is met. Once that lag is met, the  $n-1$  formula is used, where  $n$  represents the counted elements;

The method to choice the best value for  $Q$  is the basically the same but it must be used the ACF plot so the steps are:

Consider the ACF's plot, ignore the 0-th one, count the first  $n$ -elements until the closer lag to the dense part is met. Once that lag is met, the  $n-1$  formula is used, where  $n$  represents the counted elements;

Therefore the choice of parameters  $p$  and  $q$  are:

- **Milano(1,0,3)**
- **Amsterdam (2,0,3)**
- **Denver (2,0,4)**

## 5 TASK 5: Data splitting

After the choice of the initial parameter  $p, d$  and  $q$  we split the dataset in order to consider a time interval of one week out of 30 days as unit of analysis. In fact we use the first week to train model and the second week to test the model. Since our samples are aggregated per hour, we consider as initial guess  $N=168$  samples for training. Moreover we decide to use an expanding window learning strategy.

## 6 TASK 6: Model training

After set the configuration explained in the preceding sections, models have been trained individually for each city. The initial fitted model can be seen in figures 5 and from figures 6 is possible to visualize the error distribution of the residuals. We can observe the existence of residual side-lobes, indicating the persistence of biases. Nevertheless, the observed error exhibits characteristics akin to White Gaussian Noise, aligning with our assumed Autoregressive Moving Average (ARMA) model.

*not really*

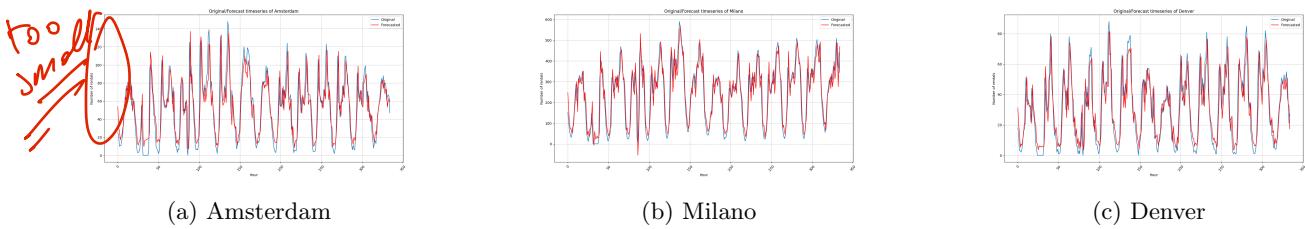


Figure 5: Starting model fitting

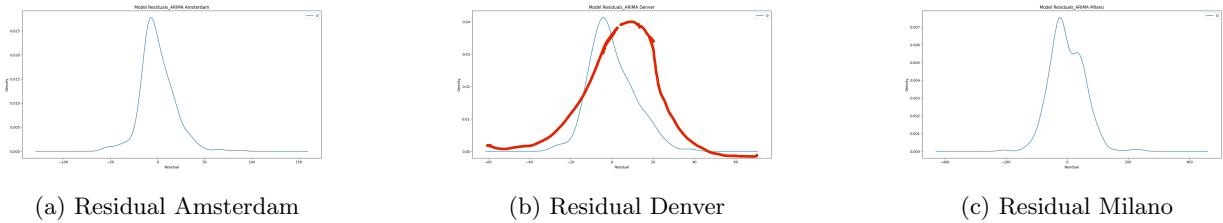


Figure 6: Residuals for each city

In Table 1, it is possible to observe the calculated **Mean Absolute Percentage Error (MAPE)**, which is percentage-based and useful for expressing errors relative to the actual values; the **Mean Square Error (MSE)**, which, by squaring the errors, penalizes larger errors more significantly. These are important metrics that define the accuracy of a forecasting method; and the **Mean Absolute Error (MAE)**, easy to interpret as it is in the same unit as the data but does not account for the direction of errors; and

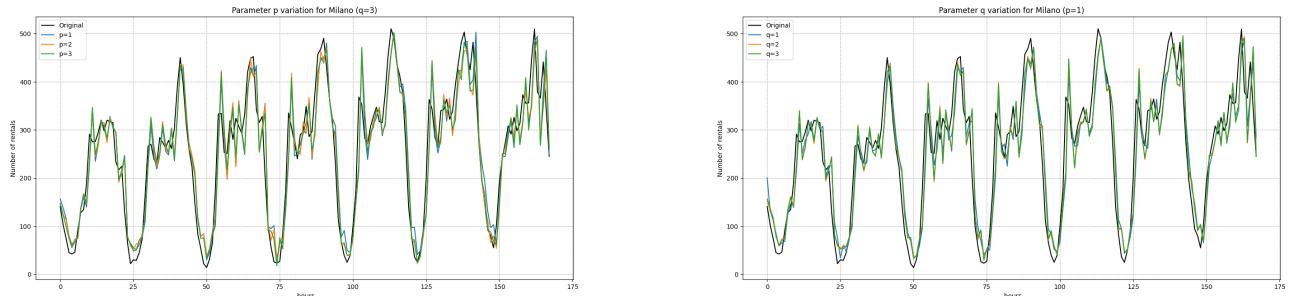
?*direction?* ?

City	Amsterdam	Milano	Denver
<b>MAPE</b>	23.010	16.104	25.522
<b>MSE</b>	400.816	3470.423	134.044
<b>MAE</b>	14.378	46.723	8.965

Table 1: Error metrics for the initial model

## 7 TASK 7: Parameters tuning

Tuning a model involves adjusting a parameter's value to enhance its performance. In this case, a grid search was conducted for parameters  $p$  and  $q$ , keeping the value of  $N$  constant and utilizing only the sliding-window training approach. Grid search is a technique for fine-tuning a model by exploring all possible combinations of parameter values and assessing the model's performance for each combination. For each pair of values  $(p, q)$  within the range of 1-3, various ARIMA models were trained, with  $N$  set to 1 week of training samples and 1 week of test samples.



(a) Variation p for Milano

(b) Variation q for Milano

Figure 7: Comparison of predictions with different p and q values - MILANO

Various types of errors, including mean absolute percentage error (MAPE), and mean square error (MPE), are calculated. However, MAPE is specifically employed for selecting optimal parameters as it enables a direct comparison of results across the three cities. The heatmaps in figures 8 illustrates MAPE and MPE for all parameter combinations. Notably, tuning for these cities involved different ranges for p and q. Some combinations were excluded from the heatmaps due to the occurrence of the error "LinAlgError: LU decomposition error". Heatmaps for Amsterdam and Denver can be found in the appendix (Fig: 14, 15).

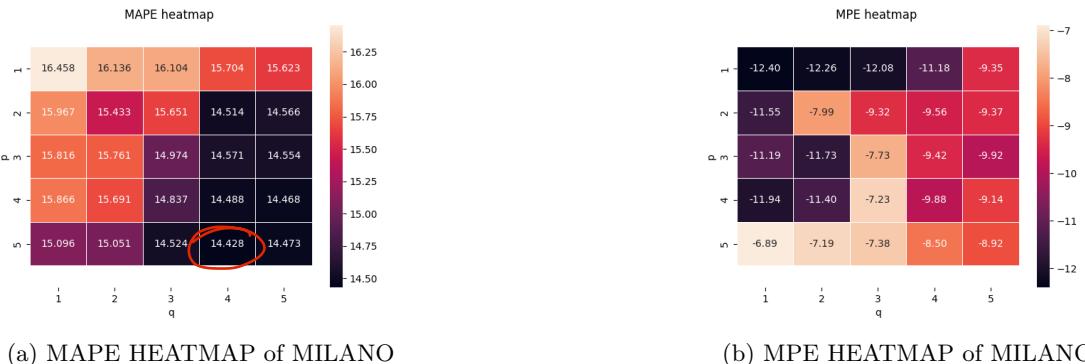


Figure 8: HEATMAP - different combinations of parameters

The most effective parameter combination for Milano is identified as  $p = 5$  and  $q = 4$ , determined by the lower MAPE. The best combinations for Amsterdam and Denver are ARIMA(2,0,5) and ARIMA(3,0,5) respectively.

Once the optimal values for p and q in the ARIMA models were determined, the investigation focuses on assessing the impact of Training Size (N) and the learning strategy by testing both the Sliding Window and Expanding Window approaches. Unlike the sliding window technique, which utilizes overlapping subsets of a fixed size, the expanding window strategy involve training a model on progressively larger subsets of data. To achieve this objective, a grid search is conducted by training several models, each corresponding to a unique combination of Training Size and Learning Strategy. For each of these models, the MAPE is evaluated, and the parameters leading to the best result, i.e., the lowest MAPE, are selected.

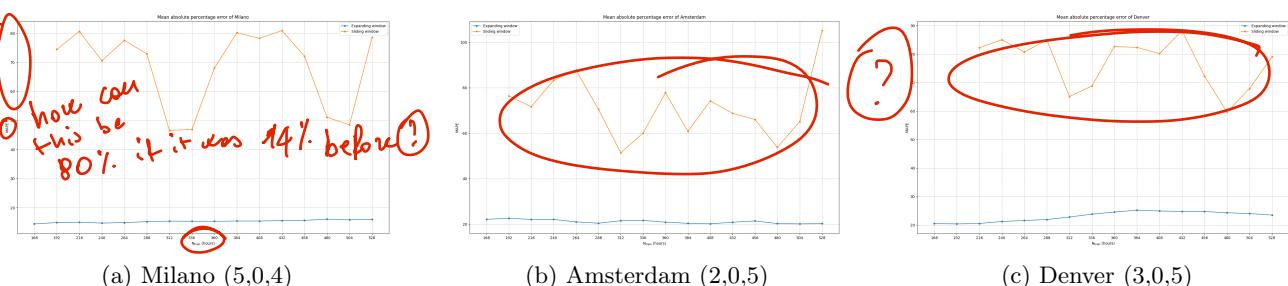


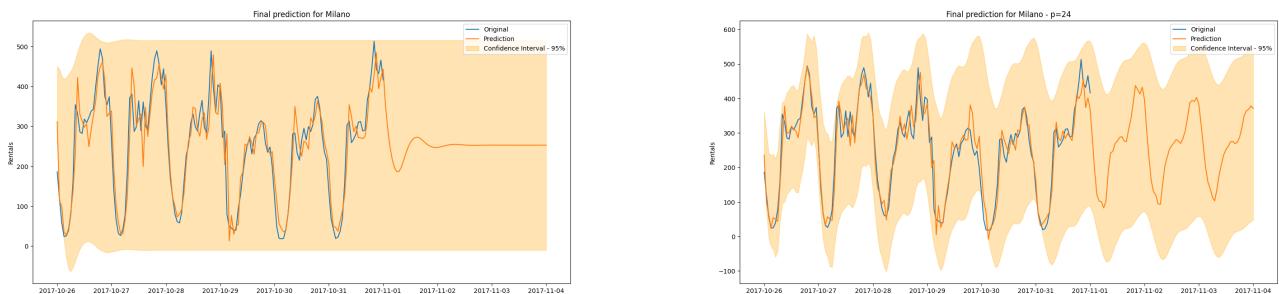
Figure 9: variation of N

From figures 9 it is possible to see that the best N, found with expanding window for all the three cities,

*do you test on the solve test set (?)*

are: for **Milano N=168**, for **Amsterdam N=504** and for **Denver N=192**.

Once the optimal ARIMA model, best N and the chosen windows technique are selected, the final model were implemented and their performance were tested as it is possible to see from figures 10.



(a) Final ARIMA prediction of Milano

(b) Final ARIMA prediction of Milano -  $p=24$

Figure 10: Final ARIMA predictions of MILANO

why  $p=24$  here ?

City	Amsterdam	Milano	Denver
<b>MAPE</b>	20.048	14.428	20.433
<b>MSE</b>	183.465	2213.271	79.004
<b>MAE</b>	10.510	36.906	6.528

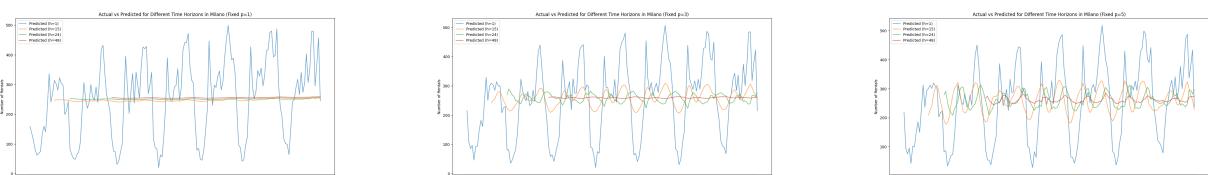
Table 2: Error metrics for final model

As it is possible to see from table 2 the error metrics decrease after the setting of the best parameters get from the task 7.

## 8 TASK 8: Impact of prediction's time horizon

In this task it has been used the time horizon which refers to the lenght of time into the future to consent the model to make predictions. The code implements an ARIMA model for time series forecasting, assessing performance across different time horizons. The series is split into training and test sets, with the model trained on historical data. Predictions are made for each specified time horizon, calculating errors such as MAE and MAPE. The results are visualized in a graph comparing predictions to actual data. This approach allows for evaluating the model's effectiveness across various time horizons, providing a comprehensive view of its performance.

By varying the parameter  $p$  in the ARIMA model, we observe a noticeable impact on the predictions. As  $p$  increases, the lines appear to deviate less from the actual values, suggesting that higher values of  $p$  allow the model to capture and incorporate more historical information. This can lead to predictions that exhibit smoother trends and better alignment with the underlying patterns in the time series data



(a) Horizon prediction with  $p=1$

(b) Horizon prediction with  $p=3$

(c) Horizon prediction with  $p=5$

Figure 11: Final Milano horizon prediction with different  $p$

- complete but some comments are missing
- do not compare against a baseline
- is the test set the same when changing  $N$ ?
- figures are too small
- last part is rushed out

## 9 Appendices

```
1 import pymongo as pm
2 import matplotlib.pyplot as plt
3 import numpy as np
4 import pandas as pd
5 from datetime import datetime, timedelta
6 import calendar
7 import gmplot
8 import math
9 import mne
10 import io
11 import json
12 from statsmodels.tsa.arima.model import ARIMA
13 import seaborn as sns
14 from statsmodels.graphics.tsaplots import plot_pacf, plot_acf
15 from sklearn.metrics import (mean_squared_error, mean_absolute_error, r2_score)
16 import os
17
18
19
20 date_init = datetime.strptime('2017-10-01T00:00:00', '%Y-%m-%dT%H:%M:%S')
21 date_finish = datetime.strptime('2017-10-31T23:59:59', '%Y-%m-%dT%H:%M:%S')
22 #date_init_Denver = datetime.strptime('2017-10-01T06:00:00', '%Y-%m-%dT%H:%M:%S')
23 #date_finish_Denver = datetime.strptime('2017-10-30T07:59:59', '%Y-%m-%dT%H:%M:%S')
24 init_unix = (date_init - datetime(1970, 1, 1)).total_seconds()
25 finish_unix = (date_finish - datetime(1970, 1, 1)).total_seconds()
26
27 timeSeries = []
28 for city in cities:
29     timeSeries['Bookings{}'.format(city)] = Bookings.aggregate([
30         {"$match": {
31             "city": city,
32             "init_time": {"$gte": init_unix, "$lte": finish_unix}
33         },
34         {"$project": {
35             "duration": {
36                 "$subtract": ["$final_time", "$init_time"]
37             },
38             "or_de": {
39                 "$ne": [
40                     {"$arrayElemAt": ["$origin_destination.coordinates", 0],
41                     },
42                     {
43                         "$arrayElemAt": ["$origin_destination.coordinates", 1]
44                     }
45                 ]
46             },
47         },
48     },
49
50     {"dow": {"$dayOfYear": "$init_date"},
51      "hour": {"$hour": "$init_date"}
52    },
53    {"$match": {
54        "duration": {
55            "$gte": 5 * 60, "$lte": 150 * 60
56        },
57        "or_de": True
58    },
59  },
60
61    {"$group": {
62        "_id": {"dow": "$dow", "hour": "$hour"},
63        "totOFbookings": {"$sum": 1}
64    }
65  },
66  ],
67  {"$sort": {
68      "_id": 1
69  }
70 }
71 ])
72 })])
```

## 9.1 TASK 1

```
1 def timeSeries_day(title, Bookings):
2     listBookings = list(Bookings)
3     DFBBookings = pd.DataFrame(listBookings)
4     labels = []
5     ticks = []
6
7     for i in range(DFBookings.shape[0] - 1):
8         if (DFBookings["_id"][i]["hour"] == 0) or (
9             (DFBookings["_id"][i + 1]["dow"] - DFBBookings["_id"][i]["dow"]) != 0) and (
10                (DFBookings["_id"][i + 1]["dow"] - DFBBookings["_id"][i]["dow"]) != 1):
11            ticks.append(i)
12            formatted_date = day_of_year_to_date(DFBookings["_id"][i + 1]["dow"], 2017)
13            labels.append(formatted_date)
14
15     plt.figure(figsize=(10, 6))
16     plt.xlabel("Days of October")
17     plt.ylabel("N. of Bookings")
18     plt.title(title)
19     plt.plot(DFBookings["totOfbookings"], label="Bookings")
20     plt.xticks(ticks,
21                labels=labels,
22                rotation=-30)
23     plt.legend(loc='best')
24     plt.grid(True, which="both")
25     plt.show()
26
27     DFBBookings.columns = ['_id', 'totOfbookings']
28     df_result = DFBBookings.to_csv(index=False) # Save to CSV without index
29
30     return df_result
31
32
33 df = time_series_df = timeSeries_day(city + "_timeSeries", timeSeries['Bookings' + city])
```

## 9.2 TASK 2

```
1 def add_miss(df, title):
2     dow_hour_combinations = [(dow, hour) for dow in range(275, 305) for hour in range(24)]
3     new_rows = []
4
5     #print(df)
6
7     for dow, hour in dow_hour_combinations:
8         if not any(((eval(entry['_id'])['dow'] == dow) and (eval(entry['_id'])['hour'] == hour))
9                    for _, entry in
10                     df.iterrows()):
11             new_rows.append({'_id': {'dow': dow, 'hour': hour}, 'totOfbookings': 0})
12
13     new_df = pd.DataFrame(new_rows)
14     df = pd.concat([df, new_df], ignore_index=True)
15     #print(df)
16
17     df[['dow', 'hour']] = df['_id'].apply(
18         lambda x: pd.Series([x['dow'], x['hour']]) if isinstance(x, dict) else x.split(','))
19
20     df['dow'] = df['dow'].apply(lambda x: int(x.split(':')[1].strip('{}')) if isinstance(x, str)
21                                 else int(x))
22     df['hour'] = df['hour'].apply(lambda x: int(x.split(':')[1].strip('{}')) if isinstance(x, str)
23                                  else int(x))
24
25     df = df.sort_values(by=['dow', 'hour']).reset_index(drop=True)
26     #print(df)
27
28     labels = []
29     ticks = []
30
31     for i in range(df.shape[0] - 1):
32         if (df["hour"][i] == 0) or (
33             (df["dow"][i + 1] - df["dow"][i]) != 0) and (
34                (df["dow"][i + 1] - df["dow"][i]) != 1):
35             ticks.append(i)
36             print(df["dow"][i + 1])
37             formatted_date = day_of_year_to_date(int(df["dow"][i + 1]), 2017)
```

```

35     labels.append(formatted_date)
36     plt.figure(figsize=(10, 6))
37     plt.xlabel("Days of October")
38     plt.ylabel("N. of Bookings")
39     plt.title(title)
40     plt.plot(df["totOFbookings"], label="Bookings")
41     plt.xticks(ticks=ticks,
42                 labels=labels,
43                 rotation=-30)
44     plt.legend(loc='best')
45     plt.grid(True, which="both")
46     plt.show()
47
48     return df
49
50 df_miss = add_miss(pd.read_csv(io.StringIO(df)), city + "_timeSeries with missed data")

```

### 9.3 TASK 3

```

1 def stationarity(df, title):
2     labels = []
3     ticks = []
4     for i in range(df.shape[0] - 1):
5         if (df["hour"][i] == 0) or (
6             (df["dow"][i + 1] - df["dow"][i]) != 0) and (
7                 (df["dow"][i + 1] - df["dow"][i]) != 1):
8             ticks.append(i)
9             print(df["dow"][i + 1])
10            formatted_date = day_of_year_to_date(int(df["dow"][i + 1]), 2017)
11            labels.append(formatted_date)
12
13 df['MA'] = df['totOFbookings'].rolling(24 * 7).mean() # Moving average
14 df['MS'] = df['totOFbookings'].rolling(24 * 7).std() # Moving std
15 plt.figure(constrained_layout=True)
16 plt.plot(df['totOFbookings'], linewidth=1, label='Number of rentals')
17 plt.plot(df['MA'], linewidth=2, color='r', label='Moving Average')
18 plt.plot(df['MS'], linewidth=2, label='Moving Std')
19 plt.title('Moving Average & Standard deviation of bookings of ' + title)
20 plt.xlabel('Date')
21 plt.ylabel('Number of bookings')
22 plt.xticks(ticks=ticks,
23             labels=labels,
24             rotation=-30)
25 plt.grid(linestyle='--', linewidth=0.8)
26 plt.legend()
27 plt.savefig(title + 'MEAN_&_STD.png', format='png')
28 plt.show()
29
30     return df
31
32 df_stat = stationarity(df_miss, city)

```

### 9.4 TASK 4

```

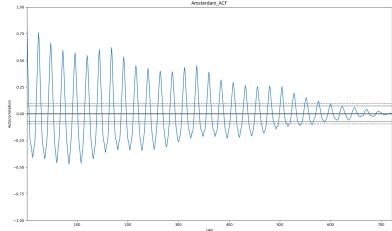
1 def autocorrelation(df, title):
2     plt.figure(constrained_layout=True)
3     pd.plotting.autocorrelation_plot(df["totOFbookings"])
4     plt.title(title + '_ACF')
5     plt.grid()
6     plt.show()
7
8 n_lags = 48
9 fig, ax = plt.subplots(constrained_layout=True)
10 plot_acf(df["totOFbookings"], ax=ax, lags=n_lags)
11 plt.title(title + '_ACF - Lags: %d' % n_lags)
12 plt.grid(which='both')
13 plt.xlabel("Lag")
14 plt.ylabel("Autocorrelation")
15 plt.show()
16
17 n_lags = 48
18 fig, ax = plt.subplots(constrained_layout=True)
19 plot_pacf(df["totOFbookings"], ax=ax, lags=n_lags)

```

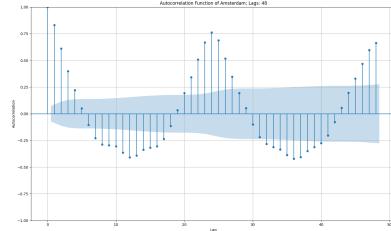
```

20 plt.title(title + '_PACF - Lags: %d' % n_lags)
21 plt.grid(which='both')
22 plt.xlabel("Lag")
23 plt.ylabel("Partial Autocorrelation")
24 plt.show()
25
26 autocorrelation(df_miss, city)

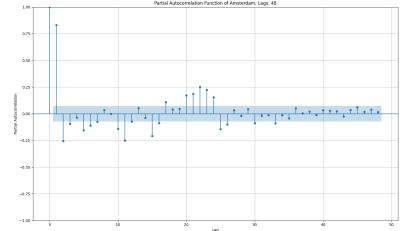
```



(a) Autocorrelation function

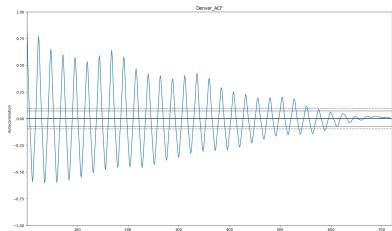


(b) ACF - zoom

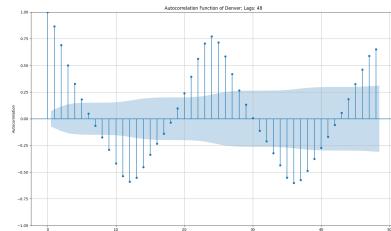


(c) PACF - zoom

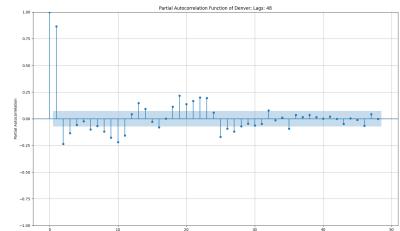
Figure 12: AMSTERDAM autocorrelation and partial autocorrelation



(a) Autocorrelation function



(b) ACF - zoom



(c) PACF - zoom

Figure 13: DENVER autocorrelation and partial autocorrelation

## 9.5 TASK 5

```

1 def split(df):
2     data = df["totOfbookings"].values.astype(float)
3     N = 7 * 24
4     train_set, test_set = data[0:N], data[N:(2 * N)]
5
6     return train_set, test_set, data
7
8 train, test, data = split(df_miss)

```

## 9.6 TASK 6

```

1 def model_training(train, test, data, title):
2     p=5
3     q=4
4     order = (p, 0, q)
5     model = ARIMA(train.astype(float), order=order)
6     model_fit = model.fit(method='statespace')
7
8     print(model_fit.summary())
9
10    plt.plot(train, label='Original')
11    plt.plot(model_fit.fittedvalues, color="red", label='Forecasted')
12    plt.title('Original/Forecast timeseries of ' + title + ' - Training phase')
13    plt.xlabel("Hour")
14    plt.ylabel("Number of rentals")
15    plt.xticks(rotation=50)
16    plt.legend(loc='upper right')
17    plt.grid(linestyle='--', linewidth=0.8)
18    plt.show()
19
20    mae = mean_absolute_error(train[0:len(model_fit.fittedvalues)], model_fit.fittedvalues)

```

```

21 mape = mae / np.mean(train[0:len(model_fit.fittedvalues)]) * 100
22 print("TRAIN DATASET : (%i,0,%i) model => MAE: %.3f -- MSE: %.3f -- R2: %.3f -- MAPE: %.3f
23 " % (1, 3, mean_absolute_error(train[0:len(model_fit.fittedvalues)], model_fit.
24 fittedvalues), mean_squared_error(train[0:len(model_fit.fittedvalues)], model_fit.
25 fittedvalues), r2_score(train[0:len(model_fit.fittedvalues)], model_fit.fittedvalues),mape
26 ))
27
28 history = train.astype(float)
29 predictions = []
30 for t in range(0, len(test)):
31     model = ARIMA(history, order=order)
32     model_fit = model.fit(method='statespace')
33     output = model_fit.forecast()
34     yhat = output[0]
35     predictions.append(yhat)
36     obs = test[t]
37     history = np.append(history, obs) # expanding window
38
39 plt.plot(test, label='Original')
40 plt.plot(predictions, color="red", label='Forecasted')
41 plt.title('Original/Forecast timeseries of ' + title + ' - Testing phase')
42 plt.xlabel("Hour")
43 plt.ylabel("Number of rentals")
44 plt.xticks(rotation=50)
45 plt.legend(loc='best')
46 plt.grid(linestyle='--', linewidth=0.8)
47 plt.show()
48
49 mae = mean_absolute_error(test, predictions)
50 mape = mae / np.mean(test[0:len(model_fit.fittedvalues)]) * 100
51 print(str(mae) + " " + str(mape))
52
53 plt.plot(data[0:len(model_fit.fittedvalues)], label='Original')
54 plt.plot(model_fit.fittedvalues, color="red", label='Forecasted')
55 plt.title('Original/Forecast timeseries of ' + title)
56 plt.xlabel("Hour")
57 plt.ylabel("Number of rentals")
58 plt.xticks(rotation=50)
59 plt.legend(loc='upper right')
60 plt.grid(linestyle='--', linewidth=0.8)
61
62 residuals = pd.DataFrame(model_fit.resid)
63 residuals.plot()
64 residuals.plot(kind='kde')
65 plt.title('Model Residuals_ARIMA ' + title)
66 plt.xlabel("Residual")
67 plt.ylabel("Density")
68 plt.show()
69
70 model_training(train, test, data, city)

```

## 9.7 Task 7

```

1 def variation(df,title,q,d):
2     data=df['totOFbookings'].values.astype(float)
3     train,test= data[0:7*24], data [7*24:(2*7*24)]
4     len_test=len(test)
5     p_var=(1,2,3)
6
7     predictions = np.zeros((len(p_var),len_test))
8     results = {"p": [], "d": [], "q": [], "mse": [], "mae": [], "mape": []}
9     try:
10         for p in p_var:
11             print('Testing ARIMA order (%i, 0, %i)' % (p,q))
12             train, test = data[0:7*24], data[7*24:(7*24+len_test)]
13             history = [x for x in train]
14             for t in range(0, len_test):
15                 model = ARIMA(history, order= (p, d, q))
16                 model_fit = model.fit( method='statespace')
17                 output = model_fit.forecast()
18                 yhat = output[0]
19                 predictions[p_var.index(p)][t] = yhat
20                 obs = test[t]
21                 history.append(obs) #expanding window
22                 history=history[1:]

```

```

23     except Exception as e:
24         print(f"Si è verificata un'eccezione di tipo {type(e).__name__}{str(e)}")
25         pass
26
27     plt.plot(test,color = 'black', label = "Original")
28     for p in p_var:
29         print("%i,0,2) model => MAE: %.3f -- MSE: %.3f -- R2: %.3f" %(p, mean_absolute_error(
30         test,predictions[p_var.index(p)]), mean_squared_error(test,predictions[p_var.index(p)]),
31         r2_score(test,predictions[p_var.index(p)])))
32
33     mae = mean_absolute_error(test,predictions[p_var.index(p)])
34     mape = mae/np.mean(test)*100
35     results["p"].append(p)
36     results["d"].append(0)
37     results["q"].append(q)    results["mse"].append(mean_squared_error(test,predictions[
38         p_var.index(p)]))    results["mae"].append(mean_absolute_error(test,predictions[p_var.
39         index(p)]))
40     results["mape"].append(mape)
41     plt.plot(predictions[p_var.index(p)],label='p=%i' %p)
42
43     plt.title('Parameter p variation for '+ city + ' (q=3)')
44     plt.xlabel(" hours")
45     plt.ylabel("Number of rentals")
46     plt.legend(loc='best')
47     plt.grid(linestyle = '--', linewidth=0.8)
48     plt.show()
49
50     p = 2
51     MA_orders = (1, 2, 3)
52     predictions = np.zeros((len(p_var), len_test))
53     for q in MA_orders:
54         print('Testing ARIMA order (%i, 0, %i)' % (p, q))
55         train, test = data[0:7*24], data[7*24:(7*24 + len_test)]
56         history = [w for w in train]
57         for t in range(0, len_test):
58             model = ARIMA(history, order=(p, d, q))
59             model_fit = model.fit()
60             output = model_fit.forecast()
61             yhat = output[0]
62             predictions[p_var.index(q)][t] = yhat
63             obs = test[t]
64             history.append(obs)
65
66     plt.plot(test, color="black", label="Original")
67     for q in MA_orders:
68         print("%i,0,%i) model => MAE: %.3f -- MSE: %.3f -- R2: %.3f" % (p, q,
69         mean_absolute_error(test,predictions[p_var.index(q)]), mean_squared_error(test,
70         predictions[p_var.index(q)]), r2_score(test, predictions[p_var.index(q)])))
71
72     mae = mean_absolute_error(test, predictions[p_var.index(p)])
73     mape = mae / np.mean(test) * 100
74     results["p"].append(p)
75     results["d"].append(0)
76     results["q"].append(q)
77     results["mse"].append(mean_squared_error(test, predictions[p_var.index(q)]))
78     results["mae"].append(mean_absolute_error(test, predictions[p_var.index(q)]))
79     results["mape"].append(mape)
80     plt.plot(predictions[p_var.index(q)], label='q=%i' % q)
81
82     plt.title('Parameter q variation for '+ city + " (p=1)")
83     plt.xlabel(" hours")
84     plt.ylabel("Number of rentals")
85     plt.legend(loc='best')
86     plt.grid(linestyle = '--', linewidth=0.8)
87     plt.show()
88
89 def variation_p_d(df):
90     data = df['totOfbookings'].values.astype(float)
91     train, test = data[0:7 * 24], data[7 * 24:(2 * 7 * 24)]
92     test_len = len(test)
93
94     N = 7 * 24

```

```

93 train, test = data[0:N], data[N:(2*N)] test_len = len(test)
94 lag_orders = (1, 2, 3, 4, 5)
95 MA_orders = (1, 2, 3, 4, 5)
96 # predictions = np.zeros((len(lag_orders), test_len))
97 predictions = np.zeros(((len(lag_orders) * len(MA_orders)), test_len))
98 results = {"p": [], "d": [], "q": [], "mse": [], "mae": [], "mape": [], "mpe": []}
99 combinations = range(0, (len(lag_orders) * len(MA_orders)))
100
101 comb = 0
102 for p in lag_orders:
103     for q in MA_orders:
104         print('Testing ARIMA order (%i, %i, %i)' % (p, 0, q))
105         train, test = data[0:N], data[N:(N + test_len)]
106         history = [x for x in train]
107         try:
108             for t in range(0, test_len):
109                 model = ARIMA(history, order=(p, 0, q))
110                 model_fit = model.fit( method='statespace')
111                 output = model_fit.forecast()
112                 yhat = output[0]
113                 predictions[comb][t] = yhat
114                 obs = test[t]
115                 history.append(obs) make sliding window
116                 # history = history[1:]
117
118             print("(%.i,%i,%.i) model => MAE: %.3f -- MSE: %.3f -- R2: %.3f" % (p, 0, q,
119             mean_absolute_error(test, predictions[comb]), mean_squared_error(test, predictions[comb]),
120             r2_score(test, predictions[comb])))
121
122             mae = mean_absolute_error(test, predictions[comb])
123             mape = mae / np.mean(test) * 100
124             adder = [(a - b) / a for a, b in zip(test, predictions[comb])]
125             mpe = (100 / test_len) * np.sum(add)
126             results["p"].append(p)
127             results["d"].append(0)
128             results["q"].append(q)
129             results["mse"].append(mean_squared_error(test, predictions[comb]))
130             results["mae"].append(mean_absolute_error(test, predictions[comb]))
131             results["mape"].append(mape)
132             results["mpe"].append(mpe)
133             comb += 1
134         except:
135             pass
136
137     return results
138
139 def heat_map (results):
140     results_df = pd.DataFrame(results)
141     print(results_df)
142
143     plt.figure()
144     heat_df_mpe = results_df.pivot(index='p', columns='q', values='mape')
145     ax = sns.heatmap(heat_df_mpe, annot=True, linewidths=.5, fmt='%.3f')
146     bottom, top = ax.get_ylim()
147     ax.set_ylim(bottom + 0.5, top - 0.5)
148     plt.title('MAPE heatmap ')
149     plt.show()
150
151     # MPE
152     plt.figure()
153     heat_df_mpe = results_df.pivot(index='p', columns='q', values='mpe')
154     ax = sns.heatmap(heat_df_mpe, annot=True, linewidths=.5, fmt='%.2f')
155     bottom, top = ax.get_ylim()
156     ax.set_ylim(bottom + 0.5, top - 0.5)
157     plt.title('MPE heatmap ')
158     plt.show()
159
160     best = results_df["mape"].idxmin()
161     p = results_df.loc[best]['p'].astype(int)
162     d = 0
163     q = results_df.loc[best]['q'].astype(int)
164     order = (p, d, q)
165     print("BEST: ", order)
166

```

```

167     return order, p,q
168
169
170 def N_variation (test,p,q):
171     d = 0
172     order = (p,d,q)
173     results = {"N": [], "window": [], "mse": [], "mae": [], "mape": [], "mpe": []}
174     comb = 0
175     train_size = [24 * x for x in range(7, 23)]
176     test_len = len(test)
177     predictions = np.zeros((len(train_size) * 2), test_len))
178     window = [0, 1]
179     for N in train_size:
180         for w in window:
181             try:
182                 if w == 0:
183                     a = 'Expanding Window'
184                 else:
185                     a = 'Sliding Window'
186                 print(f'Testing ARIMA best order, size: {N} hours and {a}')
187                 train, test = data[0:N], data[N:(N + test_len)]
188
189                 history = [x for x in train]
190
191                 for t in range(0, test_len):
192                     model = ARIMA(history, order=order)
193                     model_fit = model.fit( method='statespace')
194                     output = model_fit.forecast()
195
196                     yhat = output [0]
197
198                     predictions[comb][t] = yhat
199                     obs = test[t]
200                     if w == 0:
201                         history.append(obs) # expanding window
202                     else:
203                         history = history[1:] # to make sliding window
204                     print("(%,%,%) model => MAE: %.3f -- MSE: %.3f -- R2: %.3f" % (p, d, q,
205 mean_absolute_error(test, predictions[comb]), mean_squared_error(test, predictions[comb]),
206 r2_score(test, predictions[comb])))
207
208                     mae = mean_absolute_error(test, predictions[comb])
209                     mape = mae / np.mean(test) * 100
210                     adder = [(a - b) / a for a, b in zip(test, predictions[comb])]
211                     mpe = (100 / test_len) * np.sum(add)
212                     results["N"].append(N)
213                     results["window"].append(w)
214                     results["mse"].append(mean_squared_error(test, predictions[comb]))
215                     results["mae"].append(mean_absolute_error(test, predictions[comb]))
216                     results["mape"].append(mape)
217                     results["mpe"].append(mpe)
218                     comb += 1
219
220             except:
221                 pass
222
223
224             results = pd.DataFrame(results)
225             mape_expanding = results.pivot(index='N', columns='window', values='mape')
226             print(mape_expanding)
227             plt.figure(constrained_layout=True)
228             plt.plot(mape_expanding, linestyle='-', marker='o', markersize=4)
229             plt.title('Mean absolute percentage error of ')
230             plt.xticks(train_size)
231             plt.xlabel(r'N$\_{\mathit{train}}$' + ' (hours)')
232             plt.ylabel("MAPE")
233             plt.legend(["Expanding window", "Sliding window"])
234             plt.grid(linestyle = '--', linewidth=0.8)
235             plt.show()
236
237             best = results["mape"].idxmin()
238             N = results.loc[best]['N'].astype(int)
239             window = results.loc[best]['window'].astype(int)
240             if window == 0:
241                 window = 'Expanding Window'
242             else:
243                 window = 'Sliding Window'

```

```

241 print("BEST: N: %d, window: %s " % (N, window))
242 return N
243
244
245 def testing_model(N, test, order, p, q, df, city):
246     train, test = df['totOFbookings'][0:N], df['totOFbookings'][N:(N + len(test))]

247     history = train.astype(float)
248     predictions = []
249     for t in range(0, len(test)):
250         model = ARIMA(history, order=order)
251         model_fit = model.fit(method='statespace')
252         output = model_fit.forecast()
253         yhat = output[0]
254         predictions.append(yhat)
255         obs = test[t]
256         history = np.append(history, obs) # expanding window
257         plt.plot(predictions, color = "red", label='Forecasted')

258
259     plt.plot(test, label='Original')
260     plt.title('Original/Forecast timeseries of '+city+' - Testing phase')
261     plt.xlabel("Date")
262     plt.ylabel("Number of rentals")
263     plt.xticks(rotation=50)
264     plt.legend(loc='best')
265     plt.grid(linestyle = '--', linewidth=0.8)
266     plt.show()

267
268     mae = mean_absolute_error(test, predictions)
269     mape = mae / np.mean(test) * 100
270     print("TEST DATASET : (%i,0,%i) model => MAE: %.3f -- MSE: %.3f -- R2: %.3f -- MAPE: %.3f"
271           % (p, q, mean_absolute_error(test, predictions), mean_squared_error(test, predictions),
272               r2_score(test, predictions), mape))

272
273     print(df.columns)
274     train, test = df['totOFbookings'][-N:], df['totOFbookings'][-len(test):]
275     history = [x for x in train]
276     predictions = np.zeros((len(train), len(test)))
277     model = ARIMA(train.astype(float), order=order)
278     model_fit = model.fit(method='statespace')

279
280     t_start = pd.to_datetime("2017-11-15 00:00:00")
281     t_end = pd.to_datetime("2017-10-31 00:00:00")
282     past = pd.Timedelta(days=5)
283     future = pd.Timedelta(days=4)

284
285     predict = model_fit.predict(t_end - past, t_end + future)

286
287     df.index = pd.to_datetime(df.index)
288     fig, ax = plt.subplots(1, figsize=(15, 5))
289     ax.plot(df['totOFbookings'].loc[t_end - past:], label='Original')
290     ax.plot(predict, label='Prediction')
291     confidence_interval = model_fit.get_forecast(steps=len(predict)).conf_int(alpha=0.05)
292     ax.fill_between(predict.index, confidence_interval.iloc[:, 0], confidence_interval.iloc[:, 1], color='orange', alpha=0.3, label='Confidence Interval - 95%')
293     plt.title('Final prediction for ' + city)
294     plt.ylabel("Rentals")
295     plt.legend(loc='best')
296     plt.show()

297
298 variation(df_miss,city,3,0)
299
300 results = variation_p_d(df_miss)
301
302 order, p, q = heat_map(results)
303
304 N=N_variation(test,p,q)
305
306 colonne_da_rimuovere = ['dow', 'hour', 'MA', 'MS']
307 df_miss = df_miss.drop(colonne_da_rimuovere, axis=1)
308 df_miss['_id'] = df_miss['_id'].apply(lambda x: eval(x) if isinstance(x, str) else x)
309 df_miss['_id'] = df_miss['_id'].apply(lambda x: day_of_year_to_date_(int(x['dow']), 2017) + (f
310     " {str(x['hour'])}.zfill(2)}:00:00" if 'hour' in x else ''))
311 df_miss = df_miss.set_index('_id')
312 testing_model(N, test, (p, 0, q), p, q, df_miss, city)

```

```

313
314 testing_model(N, test, (p, 0, q), p, q, df_miss, city + ' - p=24')

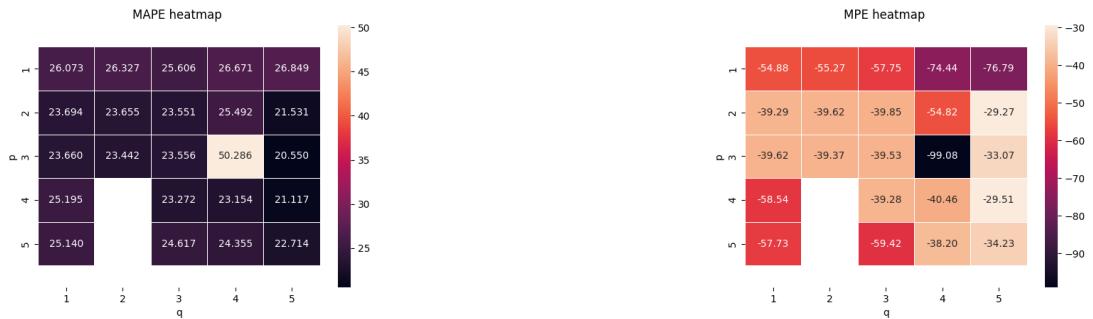
```



(a) MAPE HEATMAP of AMSTERDAM

(b) MPE HEATMAP of AMSTERDAM

Figure 14: HEATMAP - different combinations of parameters



(a) MAPE HEATMAP of DENVER

(b) MPE HEATMAP of DENVER

Figure 15: HEATMAP - different combinations of parameters

## 9.8 TASK 8

```

1 def testing_model_with_time_horizon(N, test, order, varying_p_values, q, df, city,
2   time_horizons):
3
4   train, test = df['totOFbookings'][0:N], df['totOFbookings'][N:(N + len(test))]
5
6   history = train.astype(float)
7
8   for fixed_p in varying_p_values:
9     print(f"Fixed parameter p = {fixed_p}:")
10
11   plt.figure(figsize=(10, 6))
12   for h in time_horizons:
13     predictions = []
14
15     for t in range(0, len(test) - h + 1):
16       model = ARIMA(history, order=(fixed_p, 0, q))
17       model_fit = model.fit(method='statespace')
18       output = model_fit.get_forecast(steps=h)
19       yhat = output.predicted_mean[-1]
20       predictions.append(yhat)
21       obs = test[t + h - 1]
22       history = np.append(history, obs)
23
24     mae = mean_absolute_error(test[:len(predictions)], predictions)
25     mape = mae / np.mean(test[:len(predictions)]) * 100
26     print(f"Time Horizon (h) = {h}: MAE: {mae:.3f} -- MAPE: {mape:.3f}")
27
28     plt.plot(df.index[-len(predictions):], predictions, label=f'Predicted (h={h})',
29           alpha=0.7)

```

```
29     plt.title(f'Actual vs Predicted for Different Time Horizons in {city} (Fixed p={fixed_p})')
30     plt.xlabel('Date')
31     plt.ylabel('Number of Rentals')
32     plt.legend()
33     plt.show()
34
35 testing_model_with_time_horizon(168, test, (5,0,4), [1,3,5], 4, df_miss, city, [1,15,24,48])
```