

LUMEN DATA SCIENCE 2024.

TECHNICAL DOCUMENTATION

# PREDICTION OF HOTEL OCCUPANCY

*Matija Kukić, Ivan Džanija, Valentina Vidović, Dora Baričević*

Zagreb, April 2024

# *Table of Contents*

|                                     |           |
|-------------------------------------|-----------|
| <b>1. Model implementation.....</b> | <b>3</b>  |
| 1.1. Cleaning.....                  | 3         |
| 1.2. Model training and saving..... | 4         |
| 1.3. Prediction making.....         | 7         |
| <b>2. App.....</b>                  | <b>8</b>  |
| 2.1. FastAPI.....                   | 8         |
| 2.2. Streamlit.....                 | 8         |
| 2.3. Docker.....                    | 9         |
| 2.4. Instructions.....              | 9         |
| <b>3. References.....</b>           | <b>10</b> |

# 1. Model implementation

First we opened the given train dataset and started cleaning it.

## 1.1. Cleaning

```
train_df = pd.read_parquet(path)
train_df["date_from"] = pd.to_datetime(train_df["date_from"])
train_df["reservation_date"] = pd.to_datetime(train_df["reservation_date"])
train_df["date_to"] = pd.to_datetime(train_df["date_to"])
train_df["stay_date"] = pd.to_datetime(train_df["stay_date"])
train_df["cancel_date"] = pd.to_datetime(train_df["cancel_date"])
#print(df.head())
#df.info()
#df.describe()
b = list((train_df["date_to"] - train_df["date_from"]) / np.timedelta64(1, "D"))
del train_df["night_number"]
for i in range(len(b)):
    if b[i] == 0:
        b[i] = 1
train_df["stay_nights"] = b
train_df["price_per_night"] = train_df["price"] / train_df["stay_nights"]
#f = open("data_u_txt.txt", "w+")
#f.write(df.head().to_string())
#f.close()
```

Figure 1.1.

We opened the parquet with the given train dataset and calculated the price per night and added it as a column in the data frame.

```
df = data1
df = df[df["date_from"] >= df["reservation_date"]]
df = df[(df["reservation_date"] <= df["cancel_date"]) | (df["cancel_date"].isna())]
df = df[df["adult_cnt"] > 0]
df = df[(df["cancel_date"] < df["date_to"]) | (df["cancel_date"].isna())]
df.drop_duplicates(subset = "reservation_id", inplace = True, keep = "first")
df = df[df["reservation_status"] == "Checked-out"]
df["guest_count"] = df["adult_cnt"] + df["children_cnt"]
df["nightly_price_per_guest"] = df["price_per_night"] / df["adult_cnt"]
prices_df = df[df["price_per_night"] > 150]
prices_df = prices_df[prices_df["nightly_price_per_guest"] < 6500]
#prices_df.to_parquet("train_data_cleaned.parquet")
return prices_df
```

Figure 1.2.

Then we cleaned the data from incorrect data (eg. reservation date must be before cancel date) and deleted some outliers that were statistically insignificant (guests who paid less than 150 of the given currency and those who paid more than 6500).

## 1.2. Model training and saving

```
path = Path.cwd().parent.parent
filePath = str(path) + "/data_cleanup/second_dataset/train_data_price_corrected.parquet"
df = pd.read_parquet(filePath)
df.info()
```

*Figure 1.3.*

Here we open the parquet file where data cleaning was finished and store it as a data frame.

When executed this code gives us a data frame where for each date in the set we have the amount of rooms filled on that given day (we delete the date column and turn it into a list for easier plotting).

```
# Apply logarithm function to guest_count data
room_occupancy["room_count_log"] = np.log(room_occupancy["room_cnt"])

# Seasonal differencing
seasonal_difference = room_occupancy['room_count_log'].diff(periods=365)

# Remove NaN values resulted from differencing
seasonal_difference = seasonal_difference.dropna()
```

*Figure 1.4.*

Here we take the log of the number of rooms occupied and differentiate the data with the number of occupied rooms from a year ago (365 periods).

```
test = room_occupancy.iloc[2:368]
train = room_occupancy.iloc[368:-1]
print(test.shape, train.shape)
```

*Figure 1.5.*

Here we take a test set and a train set on which we will be training and testing the model on, train set is the data from 2009 and the test set is the data from 2008.

```
test_list = test["room_cnt"].tolist()
train_list = train["room_cnt"].tolist()
mid_point_test = list()
date_list = test.index.tolist()
for i in range(len(test)):
    mid_point_test.append((test_list[i] + train_list[i]) / 2)
```

*Figure 1.6.*

Now we create the mid\_point\_test list that allows us to take both years into account for the quantile that will be used for the conformal prediction.

```
mod = ARIMA(seasonal_difference_df["room_cnt"], order = (0,0,2), seasonal_order = (1,0,1,7))
mod = mod.fit()
mod.summary()
```

*Figure 1.7.*

Here is the code that trains our ARIMA model with the parameters that we calculated to give us the best results.

```
start = len(seasonal_difference_df)
end = len(test) + len(seasonal_difference_df) - 1

pred = mod.predict(start = start, end = end)

original_train = room_occupancy.iloc[369:]
print(len(pred), len(test))
prediction_list = list()
for i, data in enumerate(train["room_count_log"]):
    prediction_list.append(data)
for i, prediction in enumerate(pred):
    prediction_list[i] += prediction
for i in range(len(prediction_list)):
    prediction_list[i] = math.exp(prediction_list[i])
print(prediction_list)
# Now, 'reverted_values' contains the predicted values in their original scale
```

*Figure 1.8.*

Here we tested our prediction on the test set. Manual re-differentiation and returning of the predictions value from its logarithmic value is visible in this code.

```
#residual calculation
residuals = list()
for i in range(len(mid_point_test)):
    residuals.append(mid_point_test[i] - predic
#residuals = mid_point_test - prediction_list
#residuals_list = residuals.values.tolist()
```

Figure 1.9.

```
residuals_list = residuals
residuals_sorted = np.sort(residuals_list)

p_values = [(np.sum(residuals_list >= r) + 1) / (len(residuals_list) + 1) for r in residuals_sorted]
alpha = 0.05

upper_quantile = np.quantile(residuals_sorted, 1 - alpha / 2)
lower_quantile = np.quantile(residuals_sorted, alpha / 2)

def predict_new_point(prediction_point):
    upper_bound = prediction_point + upper_quantile
    lower_bound = prediction_point + lower_quantile
    return lower_bound, upper_bound

lower_list = list()
upper_list = list()
for i in prediction_list:
    lo, hi = predict_new_point(i)
    lower_list.append(lo)
    upper_list.append(hi)
```

Figure 1.10.

Here we calculated the residuals of our prediction and the test set, then calculated the quantiles that we used for the conformal prediction.

```
import joblib
joblib.dump(mod, 'arima_trained.pkl')
```

Figure 1.11.

We saved the trained model as a pickle file that can be opened and used for prediction in predict.py that is used by our docker deployment implementation.

### 1.3. Prediction making

The beginning of predict.py looks something like this. First we open our model and the parquet where the test data is contained. The part left out in the commented line is where we take the data frame and group the data by date and get room\_cnt or the number of rooms occupied on a given date from the test set the exact same as we did in model training, the room\_cnt is only used to compare the data we predict later. Then we extract the starting date from the test data and the end date and give our models prediction. The rest of the code is exactly the same as the model training code, and a plot comparing our conformal prediction and the test data is given at the end.

```
model = jl.load('arima_trained.pkl')
df = pd.read_parquet('test.parquet')

#... grouping by day...

df['date_from'] = pd.to_datetime(df['date_from'])
df['date_to'] = pd.to_datetime(df['date_to'])
date1 = df["date_from"].min()
date2 = df["date_to"].max()
start = (date1 - datetime(2010,1,2)).days
end = int((date2-date1).days) + start - 1
#print(start,end,date1,date2)
pred = model.predict(start+366,end+366)
```

Figure 1.12.

## 2. App

We put together FastAPI, Streamlit and Docker Compose to make a web app. FastAPI handles the behind-the-scenes stuff, making our APIs strong and adaptable, which is super important for machine learning. Streamlit makes it easy for users to interact with and see data, so everyone can understand it better. Docker Compose helps us manage everything so our app runs smoothly.

### 2.1. FastAPI

FastAPI is a modern, fast (as the name suggests), web framework for building APIs with Python 3.7+ based on standard Python type hints. It is built on top of Starlette for the web parts and Pydantic for the data parts. FastAPI is known for its high performance due to its asynchronous capabilities and automatic generation of interactive API documentation using Swagger UI and ReDoc. It's often chosen for its ease of use, typing support and performance, making it a popular choice for building APIs, especially in the data science community.

We made two endpoints:

- one with GET method that sends you a welcome message
- and one with POST method that analyzes a parquet file and makes predictions

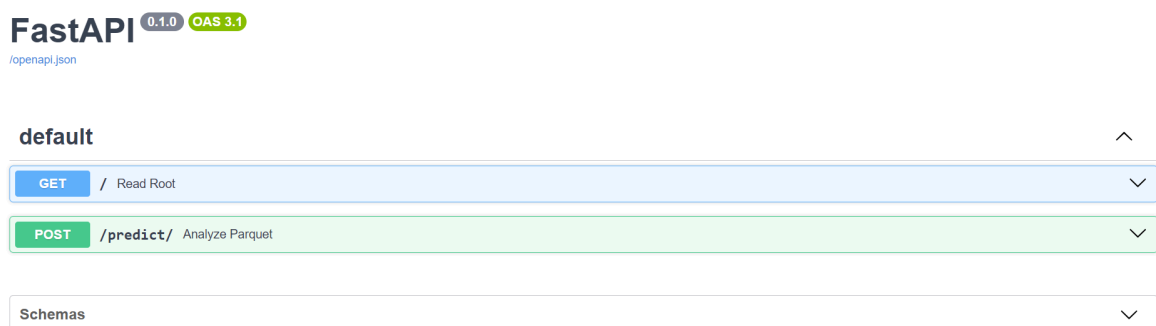


Figure 2.1.



## 2.2. Streamlit

Streamlit is an open-source Python library that makes it easy to create interactive web apps for data science and machine learning projects. It allows you to build interactive web applications directly from Python scripts, without the need for HTML, CSS or JavaScript. Streamlit is designed to be simple and intuitive, enabling data scientists to quickly prototype and share their work with others. It provides various widgets and components for building interactive elements such as sliders, buttons and plots and it automatically updates the app in real-time as the underlying Python code changes.

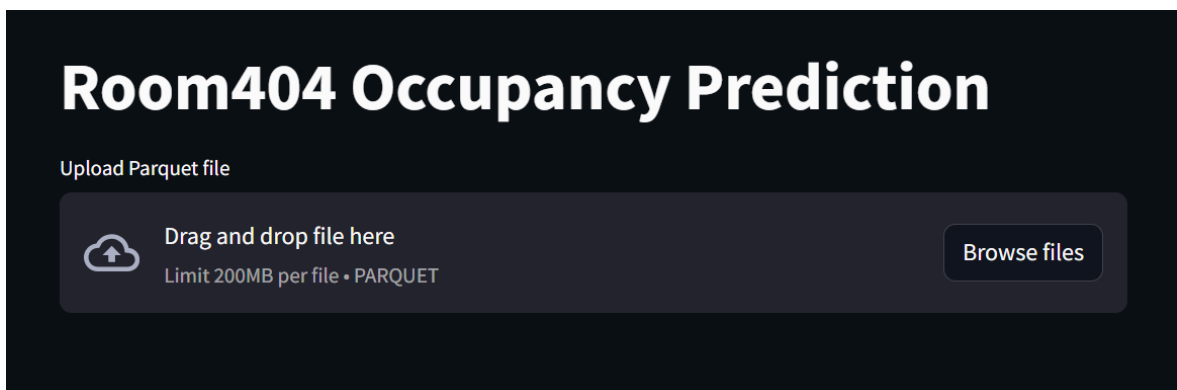


Figure 2.2.

## 2.3. Docker

Docker is a platform that allows you to package, distribute and run applications in lightweight, isolated containers. Containers are similar to virtual machines but more lightweight and portable, as they share the host operating system's kernel and only contain the application and its dependencies. Docker uses a declarative syntax called Dockerfile to define the environment and dependencies required for running an application. Docker Compose is a tool that allows you to define and run multi-container Docker applications using a YAML file. It simplifies the process of managing complex applications with multiple services by defining them in a single configuration file and enabling easy orchestration of containers.


## 2.4. Instructions

To use the app, follow these steps:


- Use Docker Compose to build and run the Docker containers. Open your terminal and enter the command: `docker-compose up --build`
- This command will launch two services: one for the frontend and one for the backend
- After the containers are up and running, open your web browser and visit <http://localhost:8501> to see the app
- For testing the model, upload the Parquet file that has same data columns as the train dataset
- The output on the screen will show the plot that has predictions for given time
- You can also check out the API service at <http://localhost:8000>

### Room404 Occupancy Prediction

Upload Parquet file

 Drag and drop file here  
Limit 200MB per file • PARQUET

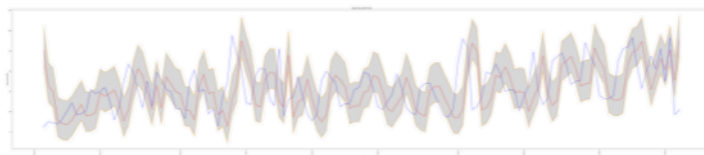
Browse files

 test.parquet 118.0KB ×

Uploading file...

Blue line - given data

Red line - prediction data



2011-01-03

True value: 5 Prediction: 81 Upper: 104 Lower: 65

2011-01-04

True value: 10 Prediction: 45 Upper: 69 Lower: 30

Figure 2.3.

### 3. References

- [1] Python - <https://docs.python.org/3/>
- [2] FastAPI - <https://fastapi.tiangolo.com/>
- [3] Streamlit - <https://docs.streamlit.io/>
- [4] Docker - <https://docs.docker.com/>