

README file for the Challenge Task

Candidate Name: **Ram Kumar Ruppa Surulinathan**

---

I thank you very much for offering this challenge. It was fun and at the same time thought provoking.

I have shared detailed information about different steps of the task along with some intermediate results.

---

## 1. Mandatory Part (A)

### The Sourcer:

I started with the webscraping of the suggested 'tripadvisor' site but unfortunately, I was kept blocking even after I tried with different useragents and different device I have. [Response 403]

Then, I tried webscraping of Yelp website. This time it worked but the response it perfectly encrypted so that the python package was unable to decipher it. [Response 200]

Then, I Googled "berlin restaurants". I came across "Quandoo" webpage. I tried webscraping it and was successful.

For webscraping, I created a script called "quandoo\_webscraper\_app.py" inside the directory "1. quandoo\_restaurants\_scraper"

I tried to handle different edge cases:

- a. Best Case Scenario – Berlin – where there are 33 pages of results.
- b. Best Case Scenario – Frankfurt - where there are 5 pages of results.
- c. Edge Case Scenario – Rostock – where there is just 1 page result.
- d. Worst Case Scenario - Paris – No Result.

As of Dec 16, Quandoo has 647 restaurants enlisted. The webscraping also generated the same number of results.

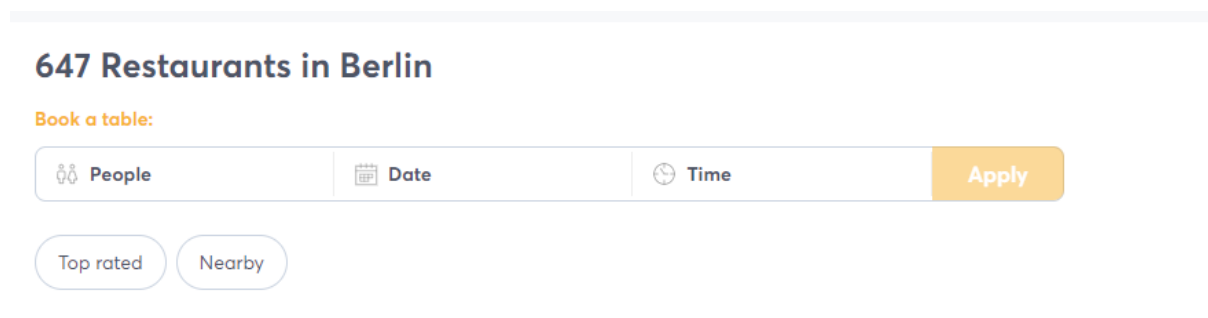


Fig 1: Berlin Restaurant Webpage (<https://www.quandoo.de/en/result?destination=berlin>)  
taken on 16.11.23

I added logging to know the progress of the scraping task and helped me to navigate to the right page in case of failure in the scraping process due to unexpected behavior (missing data).

```
2023-12-16 17:14:51,680 - __main__ - INFO - The city chosen for webscraping is Frankfurt
2023-12-16 17:14:53,147 - __main__ - INFO - There are 5 pages of results for the Frankfurt
2023-12-16 17:14:53,147 - __main__ - INFO - Now Parsing the page number: 2 for Frankfurt .....
2023-12-16 17:14:56,791 - __main__ - INFO - Now Parsing the page number: 3 for Frankfurt .....
2023-12-16 17:14:58,590 - __main__ - INFO - Now Parsing the page number: 4 for Frankfurt .....
2023-12-16 17:15:01,097 - __main__ - INFO - Now Parsing the page number: 5 for Frankfurt .....
2023-12-16 17:15:04,020 - __main__ - INFO - All the pages are successfully parsed!
```

Fig 2: Logging records generated while scraping the restaurant data for Frankfurt

To run the webscraping script:

- a. Please navigate to the directory '1.quandoo\_restaurants\_scraper'
- b. It contains 'quandoo\_webscraper\_app.py' script, Dockerfile, docker-compose file, 'scraped\_data\_results' directory to store the result files generated after scraping.
- c. Docker commands to run:
  1. `docker-compose build`
  2. `docker-compose up -d`
  3. The script can be run independently too, no input required.
- d. The generated files are placed automatically inside the folder called "scraped\_data\_results".

#### Extra improvements suggestions:

- a. To have a schedule for the sourcer to run(daily/hourly/monthly)

I can think of orchestration tools like Airflow or Mage. Here I have created a DAG file called "quandoo\_webscraping\_orchestration.py" It is just a template and didn't execute it.

- b. Think how to fetch only new restaurant data on each run

I will try to compare the new results with the previous results and identify those new restaurants name that are not present and scrape only those and add it to the results. [This could be iterative approach].

---

#### The Database:

- A. Deploy a database. The database is expected to be supported by DBT:
- B. Build the initial database layer to store the sourcer data

I have no prior experience in deploying a real database. My knowledge was limited to pgadmin. For this task, I tried using a docker-compose method referring Youtube tutorials.

To spin up a postgres database:

- a. Please navigate to '2.database\_setup' directory.
- b. It contains docker-compose file, postgres\_connection.py script and quandoo\_berlin\_results.csv file.
- c. Docker commands to run:

1. docker-compose build
  2. docker-compose up -d
  3. Please log on to localhost:8080 to see the database homepage.
- d. The database login page appears like this:

← → ↻ 🏠 ⓘ localhost:8080

Language: English ▼

**Adminer** 4.8.1

(PostgreSQL) postgres@database -

**Login**

|                 |                       |
|-----------------|-----------------------|
| <b>System</b>   | PostgreSQL ▼          |
| <b>Server</b>   | database              |
| <b>Username</b> | postgres              |
| <b>Password</b> | .....                 |
| <b>Database</b> | scraped_data_database |

☒ Permanent login

Fig 3: Loginpage of the database

- e. After login, I created a table called “berlin\_restaurants\_table” in the database.

Language: English ▼

PostgreSQL » database » scraped\_data\_database » public » berlin\_restaurants\_table » Alter table

**Adminer** 4.8.1

DB: scraped\_data\_database ▼  
Schema: public ▼

[SQL command](#) [Import](#)  
[Export](#) [Create table](#)

[select berlin\\_restaurants\\_table](#)

**Alter table: berlin\_restaurants\_table**

Table name: berlin\_restaurants\_table

| Column name         | Type      | Length | Options       | NULL                     | AI?                              | +                                |
|---------------------|-----------|--------|---------------|--------------------------|----------------------------------|----------------------------------|
| Id                  | integer ▼ |        |               | <input type="checkbox"/> | <input checked="" type="radio"/> | <input type="button" value="x"/> |
| Restaurant_name     | text ▼    |        | (collation) ▼ | <input type="checkbox"/> | <input type="radio"/>            | <input type="button" value="x"/> |
| Restaurant_location | text ▼    |        | (collation) ▼ | <input type="checkbox"/> | <input type="radio"/>            | <input type="button" value="x"/> |
| Restaurant_cuisine  | text ▼    |        | (collation) ▼ | <input type="checkbox"/> | <input type="radio"/>            | <input type="button" value="x"/> |
| Restaurant_score    | text ▼    |        | (collation) ▼ | <input type="checkbox"/> | <input type="radio"/>            | <input type="button" value="x"/> |
| Number_of_reviews   | text ▼    |        | (collation) ▼ | <input type="checkbox"/> | <input type="radio"/>            | <input type="button" value="x"/> |

Auto Increment:  ☐ Default values ☐ Comment

Fig 4: berlin\_restaurants\_table in the database

- f. After populating the table with the data from the csv file, it looked like:

Language: English

PostgreSQL > database > scraped\_data\_database > public > SQL command

Adminer 4.8.1

DB: scraped\_data\_database  
Schema: public

SQL command [Import](#)  
[Export](#) [Create table](#)

[select berlin\\_restaurants\\_table](#)

SQL command

```
SELECT *
FROM "berlin_restaurants_table"
WHERE "Restaurant_location" = 'Mitte'
LIMIT 50
```

| ID | Restaurant_name             | Restaurant_location | Restaurant_cuisine      | Restaurant_score | Number_of_reviews |
|----|-----------------------------|---------------------|-------------------------|------------------|-------------------|
| 2  | The Reed                    | Mitte               | Contemporary Restaurant | 5.5              | 397.0             |
| 4  | El Colmado                  | Mitte               | Spanish Restaurant      | 5.5              | 111.0             |
| 8  | Mogg                        | Mitte               | American Restaurant     | 5.5              | 323.0             |
| 12 | SOY Berlin Mitte            | Mitte               | Vietnamese Restaurant   | 5.4              | 230.0             |
| 13 | kopps Restaurant            | Mitte               | Vegan Option            | 5.4              | 943.0             |
| 14 | India Club Berlin           | Mitte               | Indian Restaurant       | 5.4              | 134.0             |
| 16 | Bonfini                     | Mitte               | Italian Restaurant      | 5.4              | 251.0             |
| 18 | Yumcha Heroes               | Mitte               | Asian Restaurant        | 5.3              | 718.0             |
| 22 | gärtnerel. restaurant & bar | Mitte               | German Restaurant       | 5.3              | 199.0             |
| 24 | GaYaYa                      | Mitte               | Asian Fusion Restaurant | 5.3              | 84.0              |
| 28 | Ständige Vertretung (Stäv)  | Mitte               | German Restaurant       | 5.3              | 304.0             |
| 31 | Ristorante Cinque           | Mitte               | Italian Restaurant      | 5.3              | 169.0             |
| 38 | Pho Co Restaurant           | Mitte               | Vietnamese Restaurant   | 5.0              | 48.0              |
| 40 | Restaurant Van-Long         | Mitte               | Asian Restaurant        | 5.2              | 97.0              |
| 43 | Ristorante Roma Pizzeria    | Mitte               | Italian Restaurant      | 5.2              | 82.0              |
| 44 | Nuovo Firenze               | Mitte               | Italian Restaurant      | 5.3              | 66.0              |
| 45 | Golden Rice                 | Mitte               | Asian Restaurant        | 5.2              | 74.0              |
| 51 | Palm Beach Mitte            | Mitte               | Drinks                  | 5.1              | 101.0             |

Fig 5: Simple querying in the database

g. The script called 'postgres\_connection.py' can be used to add data and fetch the data.

### C. Build Datamarts layer to store calculated data

### D. Add more layers if it's needed

### E. Use DBT to process data from the raw to the Datamarts layers

### F. Consider how to update Datamarts on a certain schedule(hourly, daily)

Unfortunately, I have no prior exposure to Datamarts and building it. I tried to refer to some materials but didn't succeed in that. Therefore, I had skipped this segment and downstream segment associated with this.

I started with the data modeling (ER Diagram) that is the starting point for the "Transformation" task of DBT.

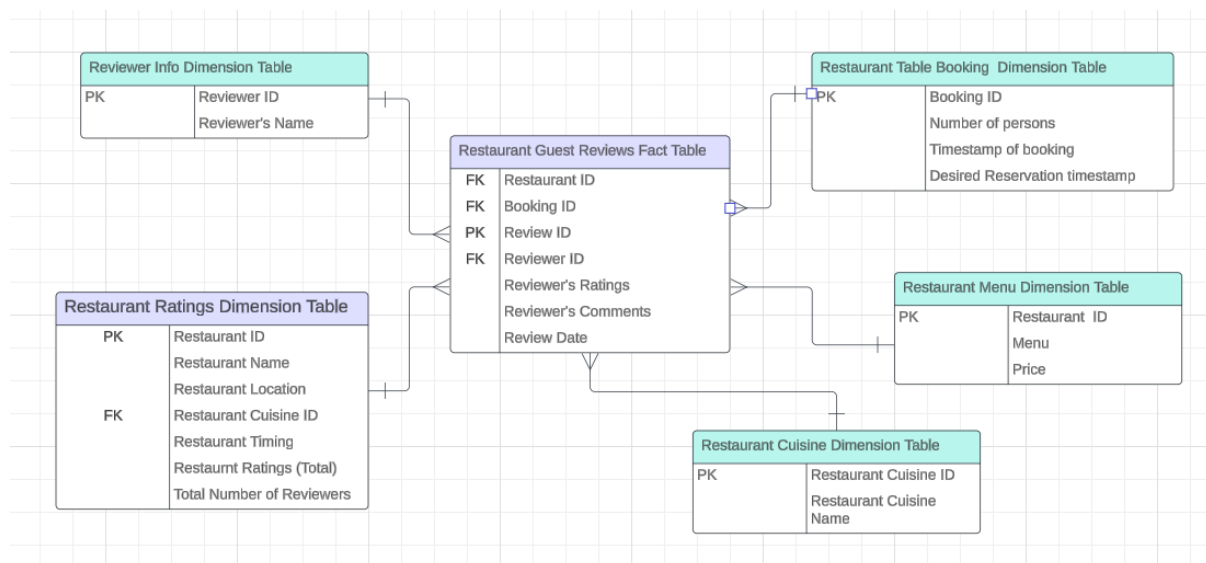


Fig 6: ERD of Quandoo Restaurants Data in Data warehouse (My assumption)

## The Optional Part (B & C)

### Visualising the Data(B)

For this task, I tried to create a simple dashboard using Plotly dash with the scraped data of Berlin. (Berlin city maybe you can relate more)

The dashboard has:

1. Data table to see/filter the scraped results.
2. Scatter plot that represents the Berliner's top 10 highest review for each kind of cuisine.
3. Scatter plot that represents the Berliner's top 10 highest review for each locality.
4. Bar plot that represents the highest number of restaurants in each locality.

To run the visualization script:

- a. Please navigate to '3.quandoo\_dashboard'.
- b. It contains 'berlin\_restaurants\_viz.py' script, docker-compose file, Dockerfile and the directory 'scraped\_data' containing the source file.
- c. Docker commands to run:

1. `docker-compose build`
2. `docker-compose up -d`
3. Please log on to <http://127.0.0.1:9002/> see the dashboard live.
4. The script also can be run independently.

| Quandoo Restaurants in Berlin |                     |                          |                  |                   |
|-------------------------------|---------------------|--------------------------|------------------|-------------------|
| Restaurant_name               | Restaurant_location | Restaurant_cuisine       | Restaurant_score | Number_of_reviews |
| filter data...                |                     |                          |                  |                   |
| Taleh Thai                    | Prenzlauer Berg     | Thai Restaurant          | 5                | 360               |
| The Reed                      | Mitte               | Contemporary Restaurant  | 5.5              | 397               |
| GOLVET                        | Tiergarten          | International Restaurant | 5.9              | 68                |
| El Colmado                    | Mitte               | Spanish Restaurant       | 5.5              | 111               |
| Vineria del Este              | Friedrichshain      | Spanish Restaurant       | 5.7              | 166               |
| Delizie D'Italia              | Prenzlauer Berg     | Italian Restaurant       | 5                | 135               |
| Ha-An                         | Prenzlauer Berg     | Vietnamese Restaurant    | 5.5              | 369               |
| Hogg                          | Mitte               | American Restaurant      | 5.5              | 323               |
| Feel Seoul Good               | Pankow              | Asian Restaurant         | 5.5              | 114               |
| Hitho Chai P-Berg             | Prenzlauer Berg     | Asian Restaurant         | 5.5              | 131               |
| Osteria Fiorello              | Prenzlauer Berg     | Italian Restaurant       | 5                | 75                |
| SOV Berlin Mitte              | Mitte               | Vietnamese Restaurant    | 5.4              | 229               |
| kopps Restaurant              | Mitte               | Vegan Option             | 5.4              | 939               |
| India Club Berlin             | Mitte               | Indian Restaurant        | 5.4              | 134               |
| Jafoy                         | Prenzlauer Berg     | Vietnamese Restaurant    | 5.4              | 167               |

Fig 7: Data table in the dashboard



Fig 8: Plots in the dashboard

**Add Data Lake layer(C)** - Unfortunately skipped this segment.

#### References:

1. <https://python.plainenglish.io/how-to-run-a-python-file-inside-a-docker-container-e17668891c85>
2. <https://python.plainenglish.io/dockerizing-plotly-dash-5c23009fc10b>
3. <https://www.youtube.com/watch?v=nIk0QIPdbcY>
4. <https://docs.getdbt.com/terms/dimensional-modeling>
5. <https://docs.getdbt.com/blog/kimball-dimensional-model>

Thanks for your time in reading my assignment.