

Assignment Documentation

Arbour Education - Data Engineer position technical Assessment

GOAL:

1. Automate the process of uploading data to Postgres db, which includes
 - a. Fetching data from Explore education statistics
 - b. Transforming the data if necessary
 - c. Uploading the data as table to Postgres db
 - d. Validation of data
2. Python code to filter data from existing table according to user selection

DATASET USED :

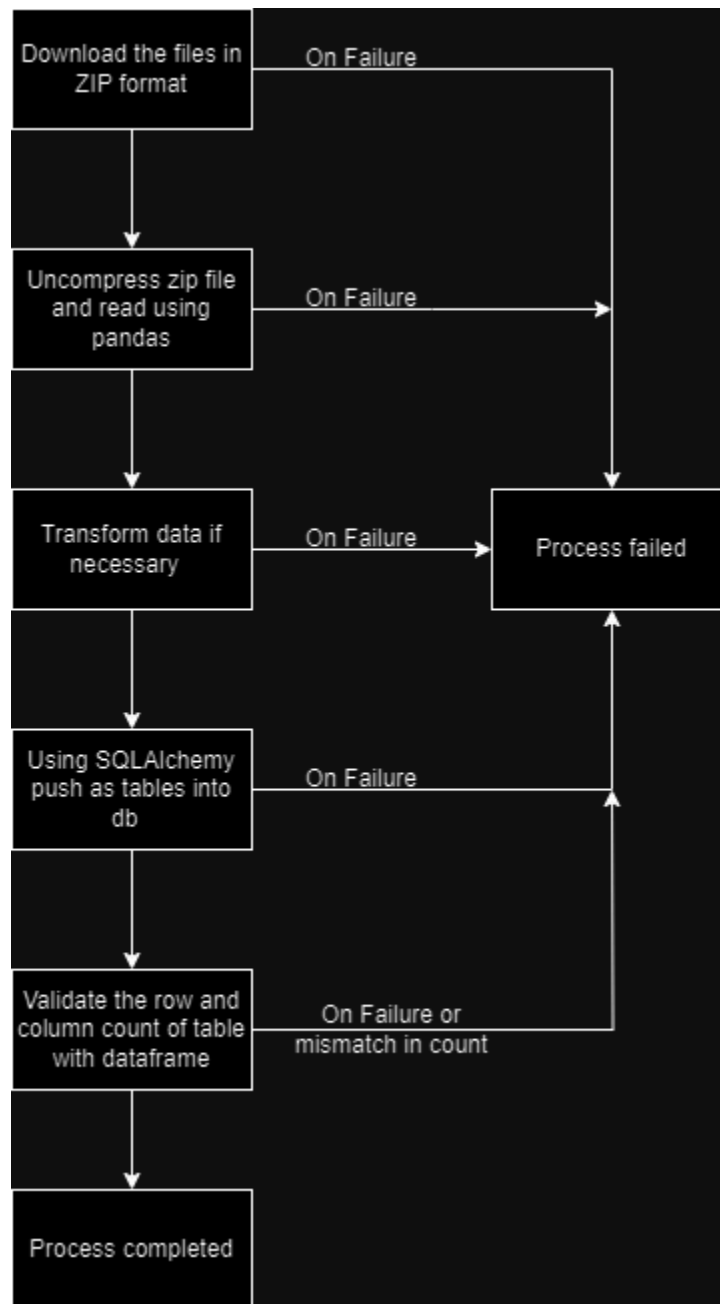
- Academic year 2022/23 School capacity
<https://explore-education-statistics.service.gov.uk/find-statistics/school-capacity>

WORKFLOW:

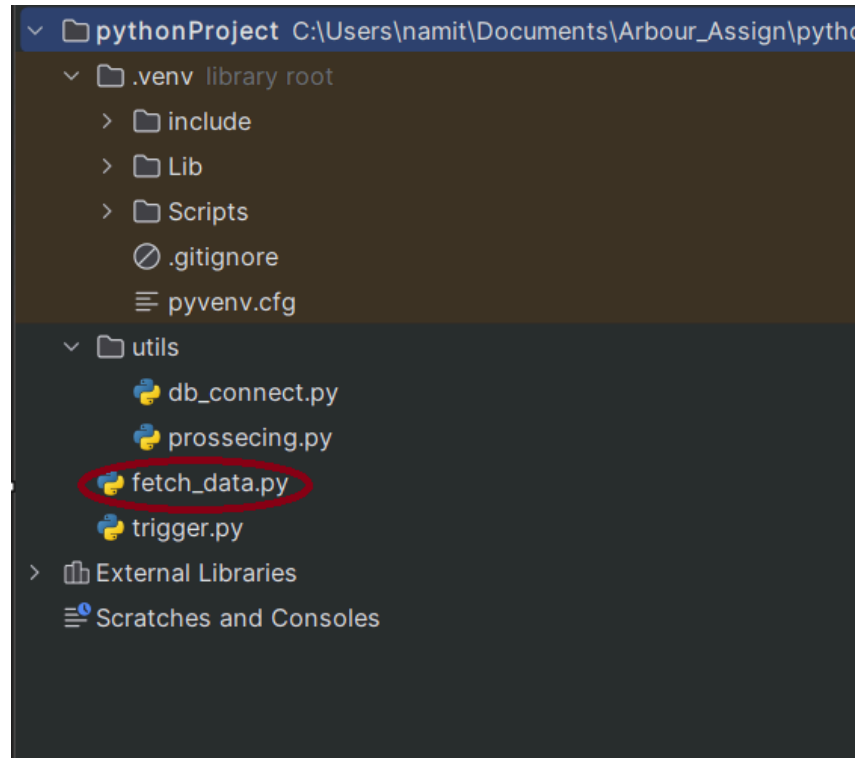
- Automate the process of uploading data to Postgres db
 - **Prerequisite**
 - **URL** from which the zip folder can be downloaded
 - Local directory to store the downloaded files
 - Postgres Db connection details
 - Prior knowledge on column data types

(Scroll down to view code flow)

- **Code Flow**



- **Application Modules**



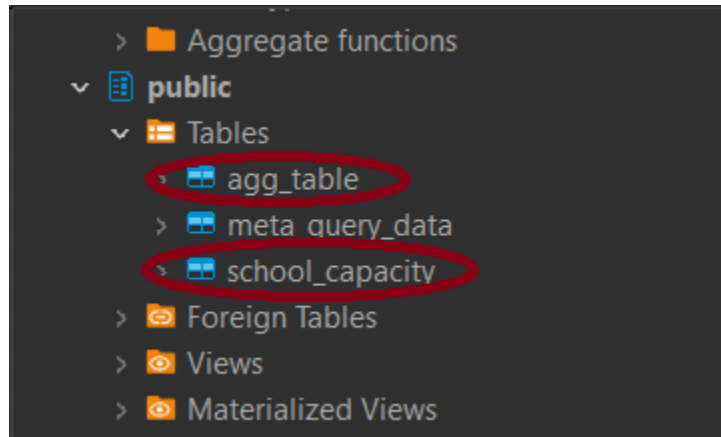
Fetch_data.py is the entry point for the code (no Arguments Required)

- **Code was able to achieve its goal**

```
in  fetch_data (1) x
C:\Users\namit\Documents\Arbour_Assign\pythonProject\.venv\Scripts\python
Begin fetching data
Starting Download
Download Completed Successfully
establishing connection with postgres
connection established
loading school-capacity_200910-202223.csv
C:\Users\namit\Documents\Arbour_Assign\pythonProject\fetch_data.py:65:
    df = pd.read_csv(self.data['downl_loc'] + '\\data\\' + file)
Upload Completed
    validating the table
Validation Started
Validation complted
loading capacity_200910-202223.csv
Upload Completed
    validating the table
Validation Started
Validation complted
Table are created

Process finished with exit code 0
```

- **Table created**



- Python code to filter data from existing table according to user selection

- **Prerequisite**

- Table meta_query_data should be created for holding the selection details. The result from code will depend on the selection criteria

- **Selection**

- **Max_Hier** : The max Hierarchy from the given option

- National
 - Region Code
 - LA code
 - School

Suppose Max_hier is chosen to be Region Code and in choice option you specify a particular region code eg 'E12000001' Then the result will contain all the data of schools of that particular selection only

- **Choice** : Specify the value of Max_hier

- **Agg** : if this is set to be true an additional query will also be created with a result of 39 rows containing total ,primary and secondary school data over the whole time frame.

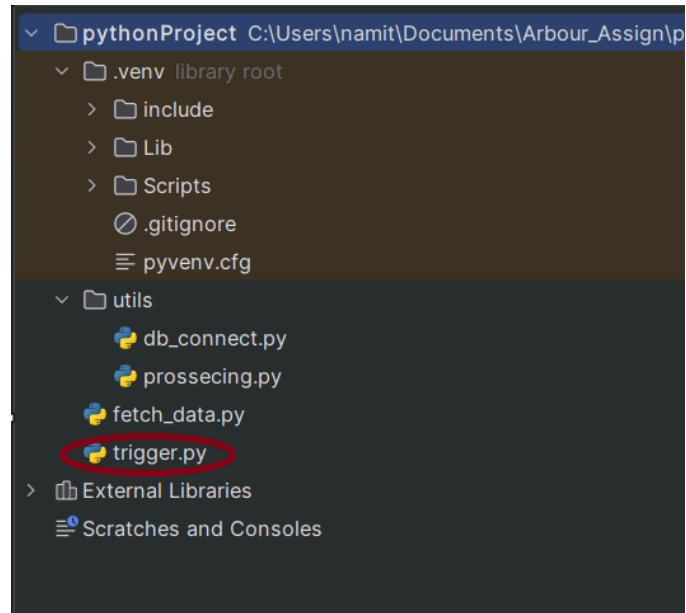
Aggregation will be based on the hier of **Max_Hier** (Agg is not possible when max_hier is School as it is the lowest aggregation)

- **Create table** : if the users want the query result in the form of a table setting this true will fetch the query result in a table

If create table is set as false the query will be save in a txt file and can be used later

- **Table name**: if the above option is set to true then table will be created with the name provided here

- **Application Module**



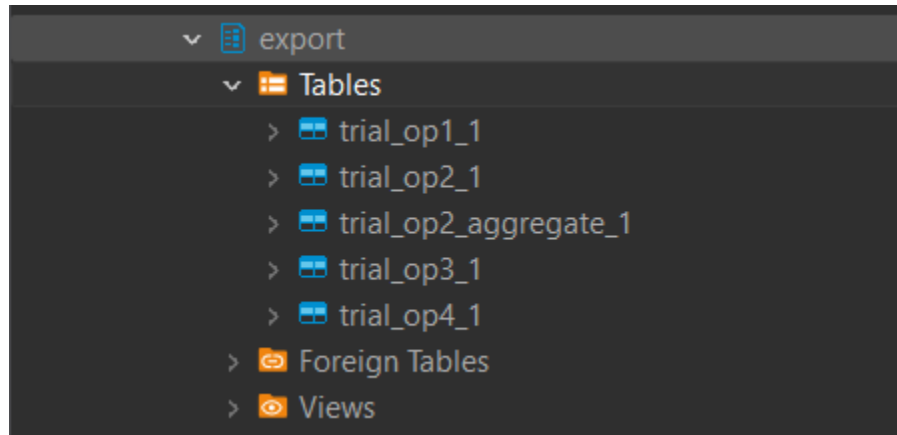
trigger.py is the entrypoint it takes one arg that is task id

- **Metadata table**

A screenshot of a database query result in a tool like DBeaver. The query is `select * from public.meta_query_data where id`. The result is displayed in a table with columns: id, key, and value. The table has 5 rows of data. The 4th row is highlighted.

	id	key	value
1	1	agg	True
2	1	table_name	trial_op2
3	1	create_table	True
4	1	choice	E12000001
5	1	max_hier	region code

- When create_table set to true



- When create_table set to false

