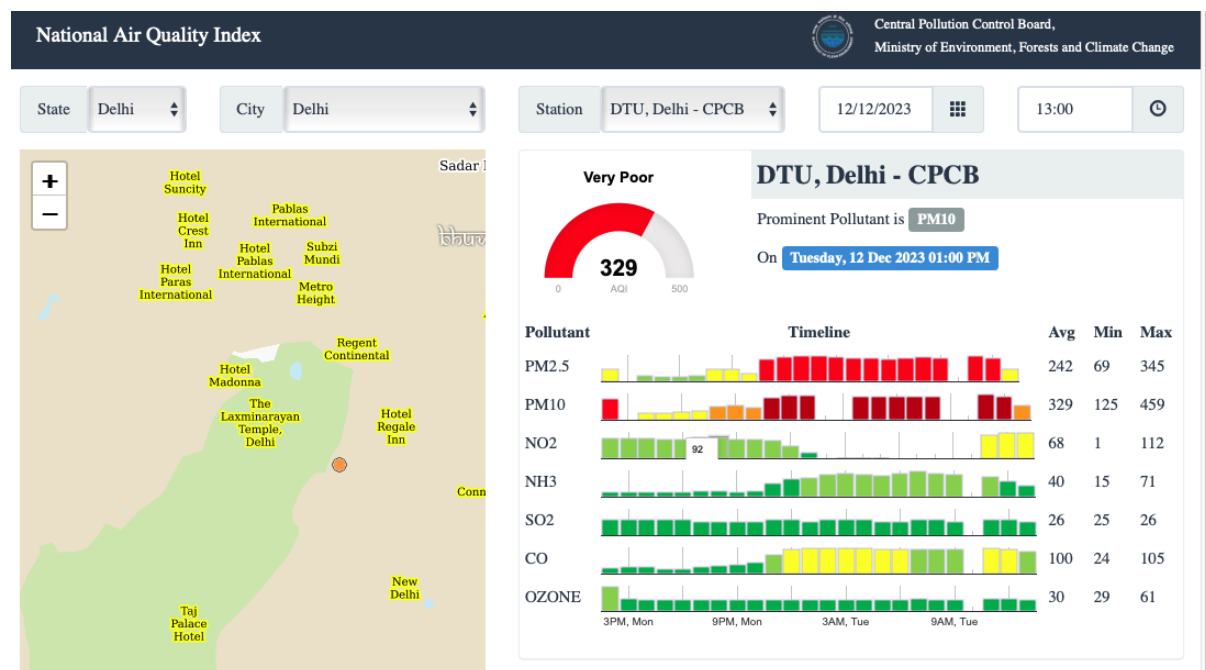# Database & Analytics Programming Documentation of Project Work Done

Thapelo Khantsi
*School of Computing*
*National College of Ireland*
*Dublin, Ireland*
*x23131535@student.ncirl.ie*

The following document entails details of the steps done in the group project.

## 1. Data Selection:

The original data from Central Pollution Control (CPCB ) which is an organisation for air and water quality monitoring in India. Since the objective was to work with semi-structured data, Kaggle was used to get data that was either XML or JSON. And the  data source had both, therefore the JSON file was selected instead. The data acquired was India City Air Quality that had different elements present in air. The goal of the project was to programmatically process and analyse the data to predict or determine pollutants that are most significant to poor air quality in India Cities.



## 2. MongoDB and PostgreSQL Set-up

These two databases were set up using Docker, MongodDB was mainly to take the JSON file dataset and put it into MongodDB, then perform ETL processes in order to analyse the data. PostgreSQL was to store structured dataset from MongoDB.

## 3. Connection to MongoDB

Python was used to achieve the connection to MongoDB. With this connection; a database and collection were created. Then the JSON file was loaded and inserted into the collection. The challenge here was that first the data was inserted not considering the structure of data as a python

dictionary. This later created problems in data processing and transformation as the data normalization efforts did not provide accurate structure for a data frame that could be used to do analysis. As we can see from the screenshot below. The first step had created just a single document which required rigorous steps in normalization of data.

```
[188]:  # Print the first few rows of the further normalized DataFrame
        print("Further Normalized DataFrame:")
        print(df_normalized_states.head())

        Further Normalized DataFrame:
                          id                                   City  \
        0      Andhra_Pradesh  [{'id': 'Amaravati', 'Station': {'id': 'Secret...
        1  Arunachal_Pradesh                                      NaN
        2              Assam                                       NaN
        3              Bihar  [{'id': 'Bettiah', 'Station': {'id': 'Kamalnat...
        4          Chandigarh                                      NaN

               City_id              City_Station_id City_Station_lastupdate  \
        0          NaN                          NaN                     NaN
        1    Naharlagun  Naharlagun, Naharlagun – APSPCB     07-12-2021 10:00:00
        2     Guwahati                           NaN                     NaN
        3          NaN                          NaN                     NaN
        4    Chandigarh                          NaN                     NaN

                          City_Station_Pollutant_Index  \
        0                                          NaN
        1  [{'Avg': 'NA', 'Max': 'NA', 'Min': 'NA', 'id':...
        2                                          NaN
        3                                          NaN
        4                                          NaN
```
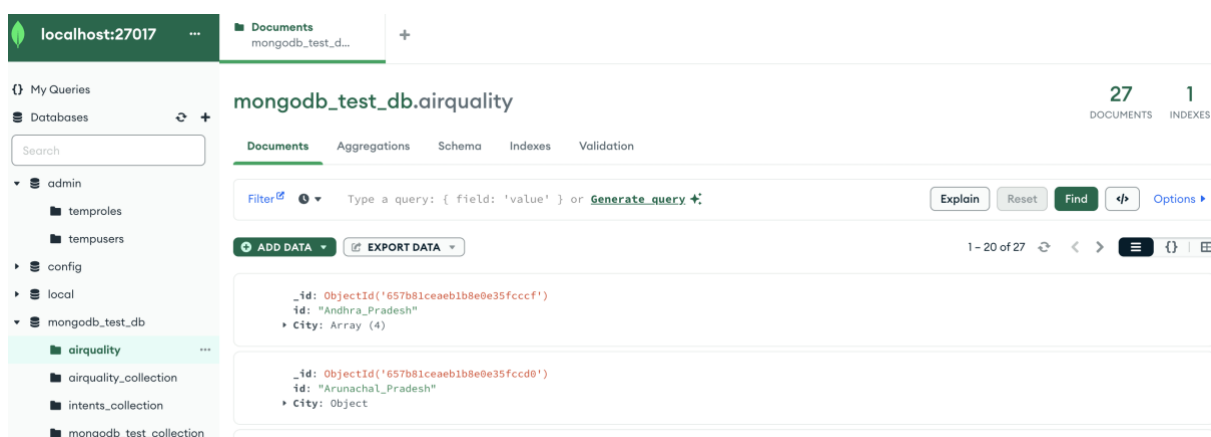
To manage this, JSON file was inserted into mongoDB with an idea of data being a python dictionary and indexing it to get a specific structure within the data.

```
# Insert the JSON data into the collection
insert_result = collection.insert_many(data['AqIndex']['Country']['State'])
print(f"Inserted {len(insert_result.inserted_ids)} documents.")
```

This solution resulted to 27 documents being created which made is more manageable to do further transformation and create a data frame that can be used to analyse the data.



4.  **Data Retrieval From MondoDB and Data Cleaning**

Python was also used to retrieve data from mongoDB in order to create a data frame that will be used for analysis. The first data frame had a shape of  (27,4) . And this data frame was not suitable

enough to work with. The data was further transformed into a second data frame that had a shape of 829 rows and 8 columns, which was much better to use for analysis.
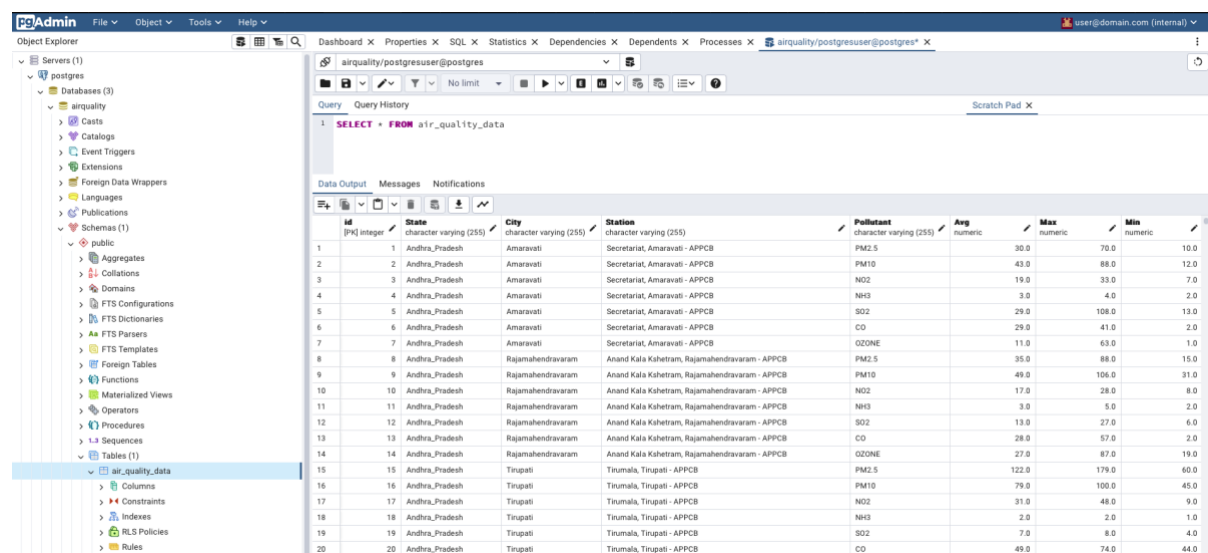
## 5. Data Visualisation

Some data visualisation were done to understand the data further and get insides.

## 6. Storing Data frame created in MongoDB into PostgreSQL database

In this process:

- Connection to SQL server was performed
- Database for air quality was created as well which would be used to process the data.
- Provided PostgreSQL details in order to create a table in PostgreSQL that will store the data based on the MongoDB data frame.
- Lastly the data from PostgreSQL was read into a data frame to use for further visualizations.



## 7. Further Predictive Analysis Using Machine Learning Models.

Linear Regression and Random Forest Regression models were used to support the analysis done from data visualisations about the factors that influence poor air quality in India.

## 8. GitHub Code Storage

Code was lastly committed to GitHub.

## 9. Reporting

We collectively written a report for the work done. Dividing each section into Dataset 1, Dataset 2, and Dataset 3; so that every member can document work done for their dataset.

## 10. Video Presentation

This was the last requirement in the project and we used teams to create a video presentation.