# Database and Analytics Programming

## MSc in Data Analytics, January 2024

1st Alan Babu Manuel
*School of Computing*
*National College of Ireland*
Dublin, Ireland
x23157143@student.ncirl.ie

2nd Melin Mary Lalu
*School of Computing*
*National College of Ireland*
Dublin, Ireland
x23185104@student.ncirl.ie

3rd Vishnunath Nharekatt
*School of Computing*
*National College of Ireland*
Dublin, Ireland
x23157143@student.ncirl.i

*Abstract—* **Sea level rise presents a critical and far-reaching consequence of climate change. As global temperatures climb, glaciers and polar ice sheets melt at an alarming rate, causing ocean expansion and a significant rise in sea levels. Air pollution and CO2 emissions act as key drivers of this process. This project delves into the interconnectedness between these factors and by analyzing it we aim to establish a clearer picture of how human activity contributes to rising sea levels and its potential consequences.**

Sea level rise is a critical and far-reaching consequence of climate change. As global temperatures rise, glaciers and polar ice sheets melt at an alarming rate, causing our oceans to expand. This expansion leads to a significant rise in sea levels, threatening coastlines, communities, and ecosystems worldwide. Air pollution and CO2 emissions play a crucial role in accelerating this process. The warmer climate caused by high levels of greenhouse gases in the atmosphere directly impacts land-based ice, causing it to melt and release the melted water into the oceans. This additional influx of water further exacerbates the problem of rising sea levels.

To gain a deeper understanding of this complex phenomenon, we have chosen three datasets for analysis. One is a structured dataset in CSV format, while the other two are unstructured JSON data. Two datasets are sourced from Kaggle, a renowned platform for data science resources, while the third dataset originates from the official Global Warming website. The structured dataset provides valuable insights into the global trends of sea level rise, while the two unstructured datasets shed light on air pollution and CO2 emission patterns across the globe. By analyzing these datasets, we aim to establish a clearer picture of the interconnectedness between these factors and their impact on rising sea levels.

## I. RELATED WORK

[1] The research paper "Impact Analysis of Sea Level Rising Problems Based on Mathematical Modeling" presents a comprehensive analysis of the impact of sea level rise on various aspects of our lives, including coastal erosion, flooding, and saltwater intrusion into freshwater sources1. The study employs mathematical modeling techniques to simulate the effects of sea level rise, providing valuable insights into the potential consequences of this phenomenon. By applying advanced mathematical models, the authors of this paper contribute to the field of sea level rise research, offering a deeper understanding of the complexities involved in predicting future sea level rise. The mathematical modeling approaches outlined in this paper can be adapted to analyze the sea level rise dataset, enhancing the understanding of the implications of sea level rise on coastal regions and global populations.

[2] This research "Prediction of air pollution through machine learning approaches on the cloud" investigates the potential of machine learning in predicting air pollution, specifically particulate matter (PM2.5). The authors acknowledge the limitations of traditional methods and explore the application of various machine learning models, including linear regression, artificial neural networks (ANNs), and Long Short-Term Memory (LSTM) networks. Their analysis reveals that LSTMs outperform the other models in predicting high PM2.5 values with reasonable accuracy. While ANNs and linear models offer decent overall performance, they struggle with high PM2.5 predictions. The paper acknowledges the need for further exploration, suggesting ensemble methods and techniques for imbalanced datasets as potential areas for improvement. Additionally, the authors highlight the computational demands of LSTMs and the challenges associated with real-time implementation for broader public benefit. They also mention ongoing work on mobile applications for asthma patients that could potentially benefit from these findings, emphasizing the importance of rich and robust datasets for such applications.

[3] The research paper "Application of Grey Prediction Model Based on Python in Carbon Emission Prediction and Low-Carbon Economic Development Analysis" by Tao Peng, Yong Sun, Yushan Zheng, Liang Ge, Feng Jin, Chonghua Wang, Fubo Zhang, Jing Yan, Xin Wang, Jieping Han, and Xiaolong Yang explores the application of the GM (1,1) forecasting model in predicting carbon emissions and analyzing low-carbon economic development.The study focuses on utilizing Python programming to conduct short-term predictions of carbon emissions in a city in Jilin Province, China, from 2017

to 2021, emphasizing the importance of understanding factors affecting carbon emissions and developing forecasting models. By employing the grey prediction model, the authors predict the increase in carbon emissions and conduct correlation analyses to forecast future emissions accurately. The paper suggests strategies to reduce industrial carbon emissions by optimizing energy structures, preventing excessive growth of industrial output, and reducing energy consumption intensity based on the model's outcomes. This research contributes valuable insights into carbon emission prediction and low-carbon economic development, providing a framework for analyzing and forecasting environmental datasets related to carbon emissions.

## II. METHODOLOGY

In this project we are using three dataset such as Sea Level Rise, Air Pollution, and Carbon Emission of the world.

### A. Description of the Datasets

**Sea Level Rise:** In this dataset we are analyzing the climate change in the world based on the sea level rise from the year 2000 to 2023. This data set also includes a wide variety of climate variables, such as information on CO2 emissions, measurements of temperature, and sea level rise observations. This dataset is excellent for deep study and evaluation of climate analysis based on sea level rise.

Dataset Columns:
Temperature: This column indicates the average temperature recorded in degrees Celsius.

CO2 Emissions: This column shows the concentration of carbon dioxide emissions measured in parts per million (ppm).

Sea Level Rise: This column represents the observed rise in sea level, measured in millimeters.

Precipitation: This column displays the amount of rainfall recorded in millimeters.

Humidity: This column indicates the relative humidity level as a percentage.

Wind Speed: This column shows the average wind speed measured in kilometers per hour.

**Air Pollution:** Global Air Pollution Dataset: This dataset provides comprehensive information on global air pollution, covering various aspects related to air quality. It includes data on pollutants like particulate matter (PM2.5), nitrogen dioxide (NO2), carbon monoxide (CO), and ozone (O3) concentrations. The dataset offers insights into air pollution levels across different regions and time periods, facilitating in-depth analysis of air quality trends, pollution sources, and potential impacts on public health and the environment. Researchers can utilize this dataset to study the dynamics of air pollution, assess the effectiveness of

pollution control measures, and develop strategies for improving air quality on a global scale.

Dataset Columns:

Country: This column identifies the specific country or region where the air quality data is collected.

City: This column specifies the name of the city within the country or region for which the air quality data is reported.

AQI Value: This column shows the overall Air Quality Index (AQI) value for the specified city.

AQI Category: This column indicates the overall AQI category for the city, based on the AQI value.

CO AQI Value: This column shows the AQI value specifically for Carbon Monoxide in the city.

CO AQI Category: This column indicates the AQI category for Carbon Monoxide based on its AQI value.

Ozone AQI Value: This column shows the AQI value specifically for Ozone in the city.

Ozone AQI Category: This column indicates the AQI category for Ozone based on its AQI value.

NO2 AQI Category: This column shows the AQI value specifically for Nitrogen Dioxide in the city.

PM2.5 AQI Value: This column shows the AQI value specifically for PM2.5 in the city.

PM2.5 AQI Category: This column indicates the AQI category for PM2.5 based on its AQI value.

**Carbon Emission:** In this dataset we are exploring the historical trends in carbon emissions worldwide from 1740 to 2014. The study aims to identify long-term trends, geographical disparities, and factors driving carbon emissions by analyzing important indicators such as total emissions, fuel-specific contributions (solid, liquid, and gas), cement output, gas flaring, per capita emissions, and bunker fuel usage. The project will use data visualization and statistical analysis to provide insights into historical patterns, the impact of industrialization, energy consumption shifts, and policy interventions on global carbon emissions, thereby contributing to a better understanding of climate change dynamics and informing sustainable mitigation strategies.

Dataset Columns:

Year: This column indicates the year the data was recorded. It specifies the timeframe for which the emissions data is applicable.

Country: This column identifies the country or region for which emissions data are reported. Each entry in the dataset typically represents a specific country or region.

Total_Emission: This column typically shows the total amount of carbon dioxide (CO2) or other greenhouse gases emitted by the country or region in the given year. It represents the total of emissions from the various sources stated in the following columns.

Solid_fuel: This column typically indicates the quantity of CO2 emitted from solid fuel combustion, such as coal, lignite, and biomass.

Liquid_fuel: This column typically depicts the CO2 emissions produced by the combustion of liquid fuels such as gasoline, diesel, and fuel oil.

Gas_fuel: This column typically indicates the amount of CO2 emitted from the combustion of gaseous fuels like natural gas.

Cement: This column typically indicates the amount of CO2 emitted during cement production procedures. The chemical reactions involved in cement production contribute significantly to greenhouse gas emissions.

Gas_flaring: This column typically depicts the CO2 emissions generated by the flaring of related gas during oil extraction and processing processes. Gas flaring is a typical technique in the oil sector, however, it adds to greenhouse gas emissions.

Per_capita: This column usually shows per capita emissions, which are the average emissions per person in a country or region. It is computed by dividing total emissions by population in the country or region.

Bunker_fuels: This column typically indicates the CO2 emissions generated by the combustion of bunker fuels in international shipping and aircraft. Bunker fuels are heavy fuels commonly utilized in marine and aviation engines.

B. *Data Processing*

To analyze a dataset, it is necessary to perform preprocessing processes. This includes reviewing the dataset for quality and integrity. Each dataset undergoes preprocessing processes to get the desired outcome.

**Sea Level Rise:**
- Checked all the data for the null and duplicate values. This dataset has no null values.
- Separated the numerical and categorical data valued columns for better visualization process.

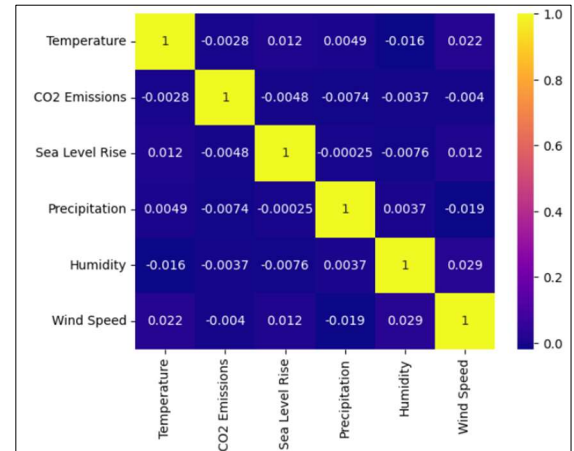- Generated a heatmap for understanding the corelation of numerical data.



Fig no:1 Correlation matrix

- Split the date column to year and date for visualization.
- Checking the sea level rise with each country
- Set the 'country' column as the data frame index for indexing the values in each row.

**Air Pollution:**
- Examined the data for null values in each column and removed them
- Checked the data for duplicate values and eradicated them.
- Dropped the column named id as it is not required for further analysis.
- Divided the columns containing numerical and categorical data values to enhance the visualization process.



- Created a heatmap to explore the correlation among numerical variables.

Fig No. 2: Global Air Pollutants Correlation matrix

**Carbon Emission:**

- Checked for the null values using isna().sum() function. In this dataset, we do not contain any missing values.

- Total was renamed to Total_emission for better representation.

- Checked for duplicate values in the dataset and removed.

- Columns with object values were converted to the proper data type for mathematical operations.

- Converting categorical variables into a numerical format using label encoding suitable for modeling.

- Identified and removed outliers from the Total_emission column.



Fig No.3 Total_emissions Outlier Removal

*C. Tools and Technologies*

- **Docker:** Docker emerges as a pivotal tool in our project's development and deployment workflow. Leveraging containerization technology, Docker provides a consistent environment for our applications across different platforms, eliminating compatibility issues. With Docker, we streamline the deployment process, ensuring that our project remains scalable, portable, and easily maintainable.

- **Python:** Python stands out as a versatile programming language, offering a plethora of libraries and frameworks that streamline development tasks. With libraries like Pandas facilitating the manipulation of data through DataFrames, and tools like Plotly enabling the creation of interactive visualizations, Python proves indispensable for efficient data preprocessing and analysis in our project. Additionally, the Luigi library is employed for ETL (Extract, Transform, Load) tasks, providing a robust framework for orchestrating complex data pipelines with ease and reliability.

- **MongoDB:** In our project, MongoDB emerges as the preferred database solution owing to its NoSQL architecture, which adeptly handles unstructured data such as JSON. Its schema-less design fosters adaptability, allowing seamless adjustments to evolving data structures. MongoDB's scalability and rapid retrieval capabilities are instrumental for accommodating and efficiently managing large and dynamic datasets.

- **PostgreSQL:** Chosen for its robustness and extensibility, PostgreSQL emerges as the backbone of our data management strategy. As an open-source relational database management system, it offers advanced features and strict standards compliance. PostgreSQL excels in scalability, effortlessly managing large datasets and supporting concurrent users. Its compatibility with various platforms ensures seamless integration into our project infrastructure.

*D. Process Flow*

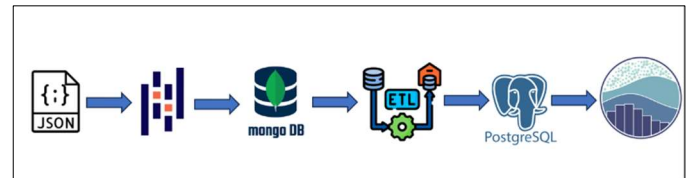The process flow for the analysis of the global air pollution dataset is as follows.



Fig No. 4: The Process Flow

The semi structured dataset, provided in JSON format, was first loaded into a Pandas DataFrame using the pd.read_csv() function.

The data was then imported into a MongoDB database using the PyMongo library. A MongoDB client was established, and the data was inserted into 3 seperate collections named "CO2_Emissions","Air_Pollution_Data" and "Climate_Analysis" within the "DAP_PROJECT" database.
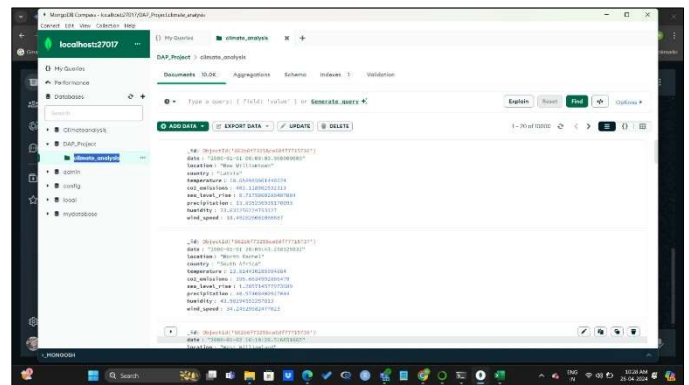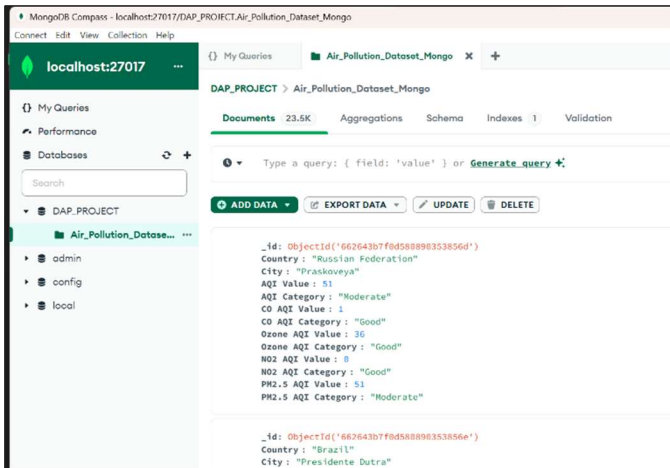


Fig No. 5: Climate_Analysis Mongo Collection

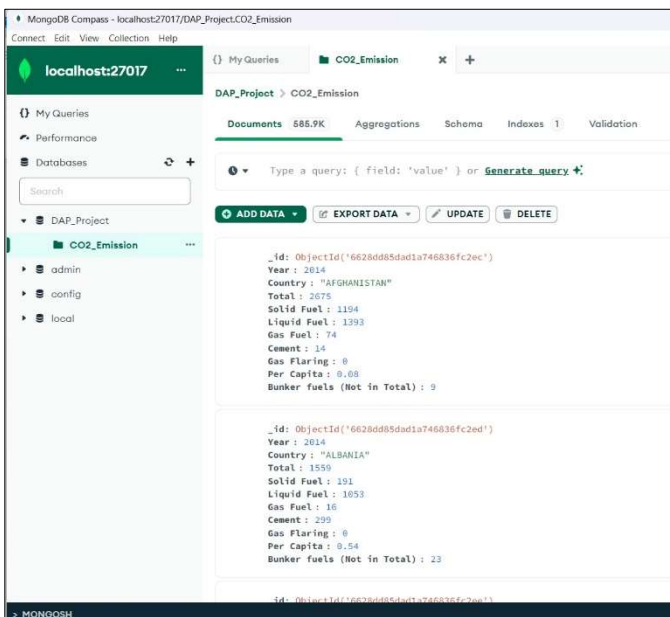Fig No. 6: Air_Pollution_Data Mongo Collection


Fig No. 7: CO2_Emissions Mongo Collection

To automate the data extraction, transformation, and loading process, a Luigi pipeline was implemented.
The pipeline consists of three tasks:

ExtractTask: This task retrieves the data from the MongoDB collection and saves it to a local output CSV/JSON file.

TransformTask: This task reads the extracted data, performs data cleaning (e.g., handling missing values and removing duplicates), and saves the transformed data to a new output CSV/JSON file.

LoadTask: This task reads the transformed data from the CSV file and loads it into a PostgreSQL database table.
The PostgreSQL database connection details, such as the host, port, username, and password, are provided as parameters to the LoadTask
.

Finally, a separate script is provided to retrieve the data from the PostgreSQL database and create a Pandas DataFrame for further analysis and visualization.

## III. RESULT AND EVALUATION

### A. *Data Visualization and Analysis for Sea Level Rise*

#### *1) Analyse the data*


Fig No. 8: Sea Level Rise Dataframe info

The first stage of analysis is to understand the dataset. The sea level dataset is read into a data frame for further analysis. The dataset consists of 8 columns which shows the main factors that affects the climate change.

#### *2) Monthly Sea level rise in each year*


Fig No.9: Bar Plot of Monthly Sea Rise

We are analysing the monthly sea level rise of each year by using the bar plot. Here we can find the month which have the largest sea level rise.

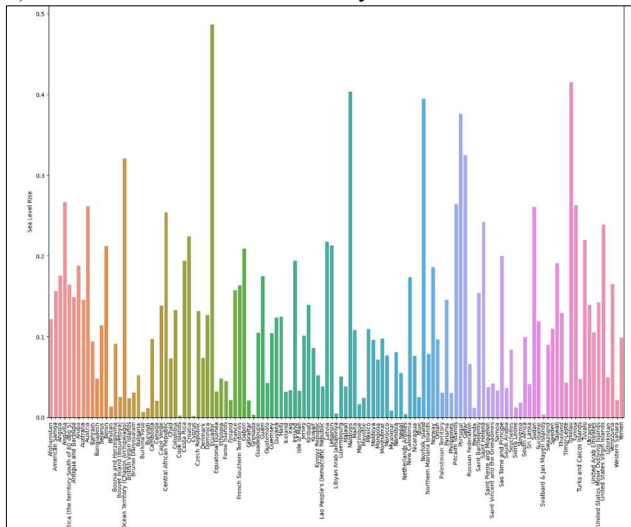## 3) *Sea level rise in each country*



Fig No. 10: Vertical Bar Plot of Sea Level Rise wr.t Countries

We can use histograms to find which country have the largest sea level rise during this period.We found that Ecudor [country in South America] has the largest sea level rise.
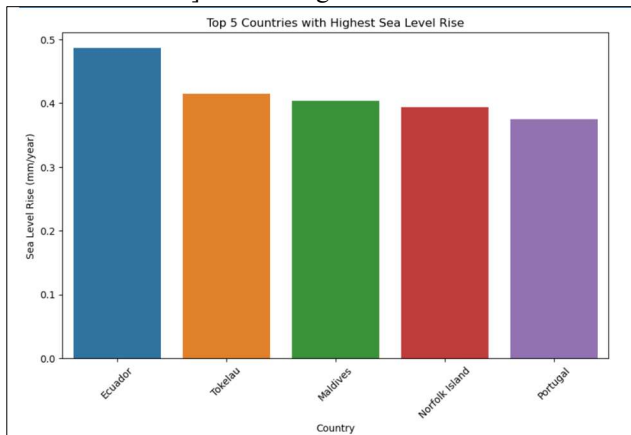


Fig No.11: Bar Plot of Top 5 Countries with Sea level Rise

We can find the top five countries which have sea level rise during this period using bar plot. Ecuador being the top contender.

### 4) *Total sea level rise during the period 2000-2023*

Here we can use the line plot to find the total sea level variation during this time period. We can see from the below figure that the period between 2015 to 2020 was having the largest sea level rise and is dipping thereafter.



Fig No. 12 Line plot of Sea level Rise During 2000-2023

## B) **Data Visualization and Analysis for Air Pollution**

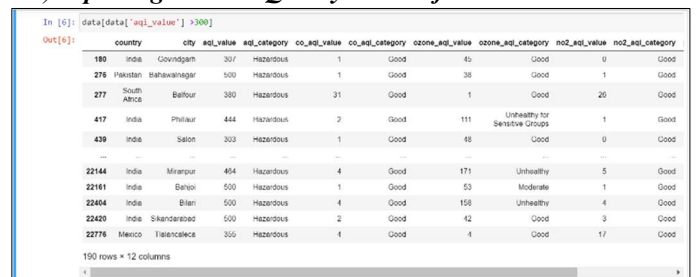### 1)*Exploring The Air Quality Index of Countries*



Fig No.13 Hazardous AQI Value

Studies have shown that an Air Quality Index (AQI) above 300 signifies hazardous air quality, posing severe health risks. According to the data, 190 countries fall into this hazardous category.
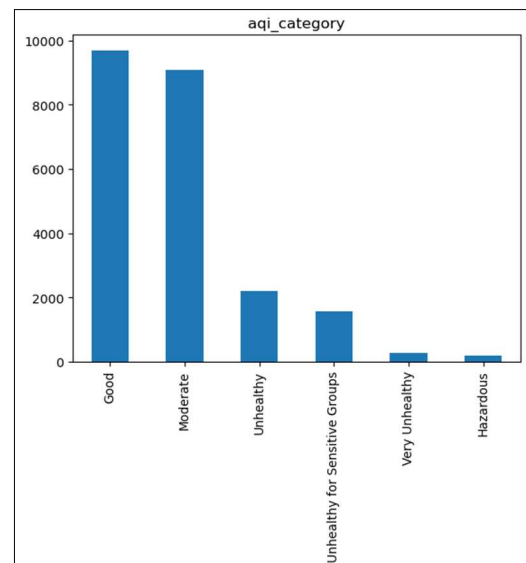


Fig No.14 Bar Plot of AQI Category

Air Quality Index (AQI) is divided into categories based on the level of risk it poses to health. When the AQI value exceeds 100, air quality becomes unhealthy, and the risk to health rises. An AQI below 50 indicates good air quality.

From the data provided, most cities have good air quality (9688). Cities with moderate air quality come second at 9087, followed by 2215 cities with unhealthy air quality. There are 286 cities with very unhealthy air quality and 190 cities with hazardous air quality, posing the greatest health risks.
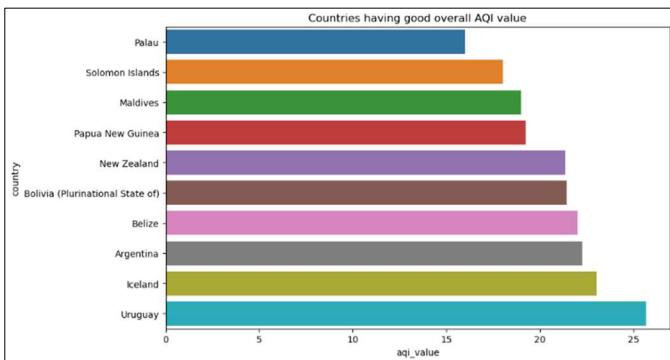


Fig No.15 Worst AQI Countries



Fig No.16 Good AQI Countries

Among the countries, Russia, USA and India have the worst overall AQI, exceeding 450. In contrast, Palau, Solomon Islands, and Maldives boast the cleanest air, with AQI values below 20.

### 2) Correlation between PM2.5 and AQI

The analysis of PM2.5 values reveals a clear and concentrated pattern in its scatter plot, suggesting a substantial contribution of this pollutant to overall air quality. In comparison, CO and Ozone appear to have a lesser influence. Notably, a strong positive correlation exists between PM2.5 AQI and overall AQI, indicating that PM2.5 directly impacts air pollution levels. While a positive correlation is also observed between AQI and both Carbon Monoxide and Ozone, it is significantly weaker compared to PM2.5
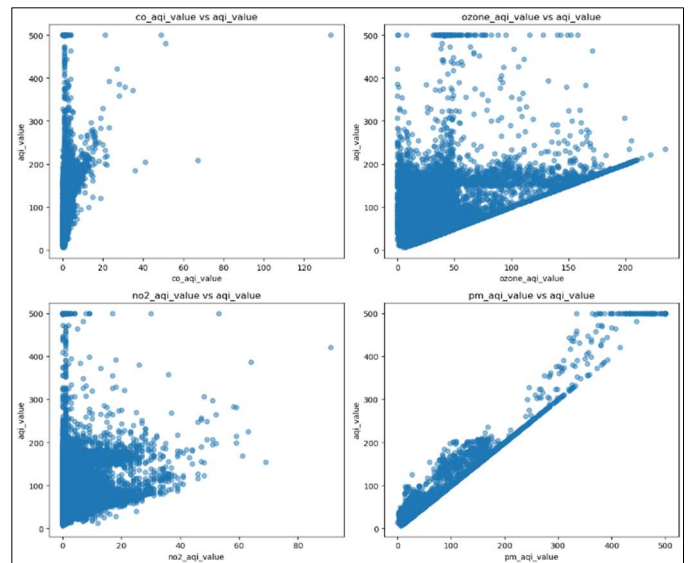


Fig No.17 Scatter plot of pollutants wrt AQI

### C) Data Visualization and Analysis for Carbon Emission

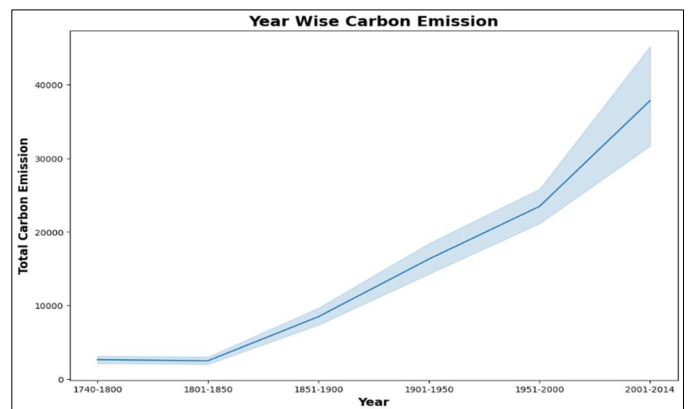#### 1) Year-wise carbon emissions



Fig No.18 Carbon Emission Trend Along the years

The figure above depicts the total carbon emissions of various countries during a certain period. This shows that emissions are increasing year after year.

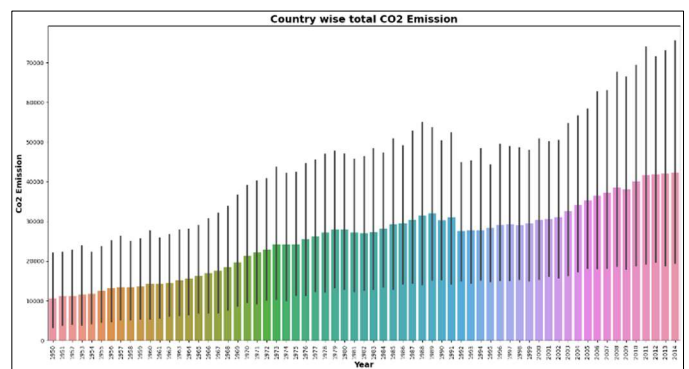#### 2) Total Co2 emissions by country



Fig No.19 Carbon Emission by Each Country

The picture above displays the total carbon emissions of various countries from 1950 to 2014.
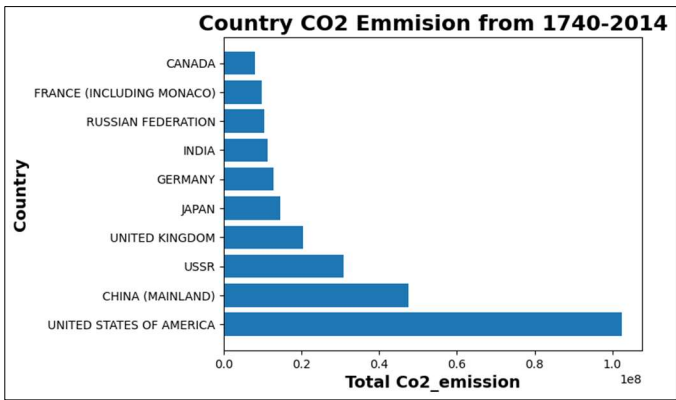
### 3) *Visualization of Top 10*



Fig No.20 Top 10 CO2 emitting Countries

The above figure shows the visualization of the top 10 countries with the highest CO2 emissions. From this, we can understand that the highest carbon emission country is the United States of America.

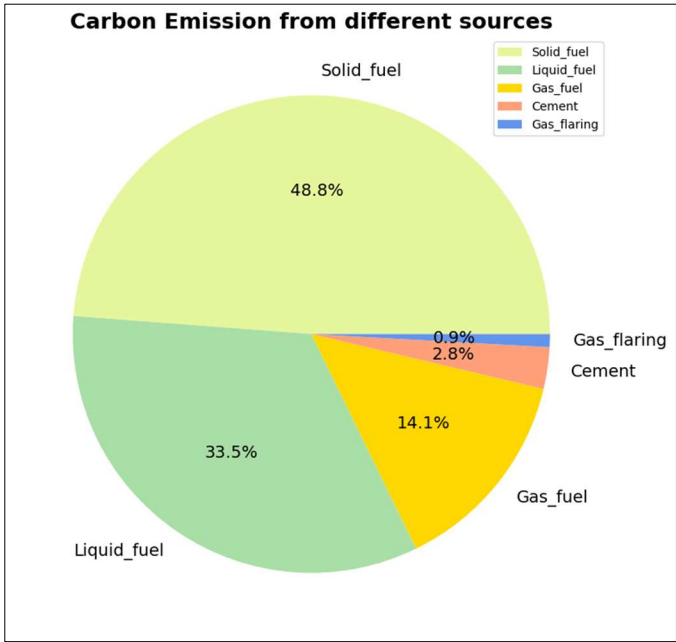### 4) *Carbon Emission from Various Sources*



Fig No.21 Pie Chart of Various Carbon Emission Sources

The figure above exhibits carbon emissions from various sources. Solid fuel emits the most carbon compared to other sources (liquid fuel, gas fuel, cement, gas flaring).

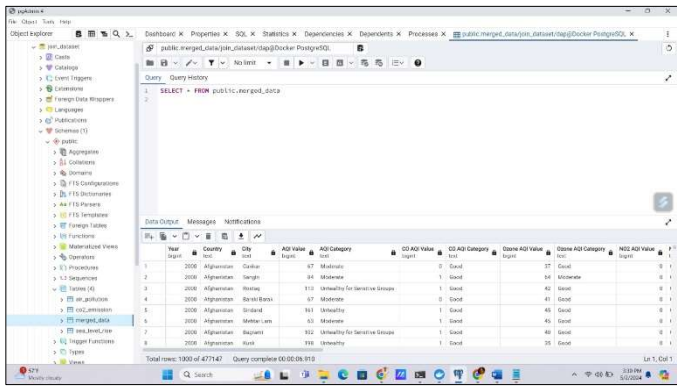## D) Comparitive Visualizations and Analysis



Fig No.22 Joined Datasets

For the purpose of comparing the datasets and draw insights we have joined the three tables namely sea_level_rise , air_pollution and co2_emission in postgre sql.
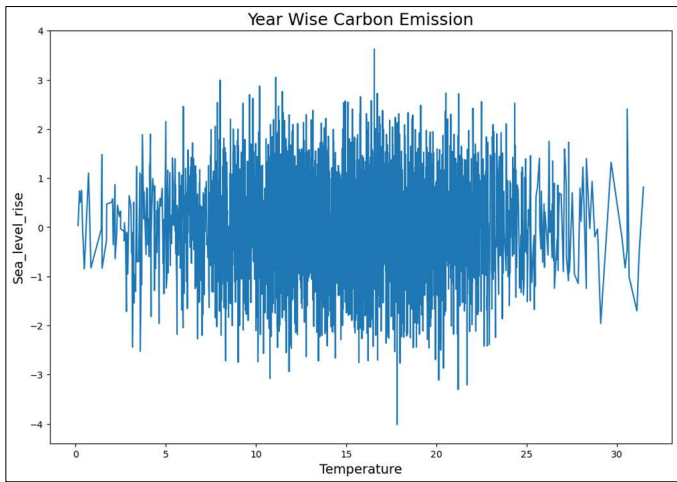


Fig No.23 Comparitive Visualization

The above graph showcases the increase in sea level rise with increase in temperature due to rising carbon emission and air pollution over the years.

## IV. CONCLUSION & FUTURE WORK

Our analysis confirms the alarming trend of rising sea levels, with an observed increase of over 8 inches since 1880, with a significant acceleration in recent years. This phenomenon is primarily driven by air pollution and CO2 emissions, which contribute to global warming. Research suggests that sea levels could potentially rise by one meter by 2050, posing a significant threat to coastal communities and ecosystems worldwide.

While this study highlights the critical link between human activity and rising sea levels, it acknowledges certain limitations. The complex nature of sea level rise is influenced by various natural factors that may not be fully captured by our analysis. Additionally, data availability presents a challenge, as high-quality data is crucial for training accurate models. Furthermore, it's important to acknowledge that any model predictions are inherently susceptible to errors due to their inherent limitations.

Several avenues exist for further investigation. Incorporating higher quality data into our models would enable more precise predictions and a deeper understanding of the contributing factors. Additionally, exploring techniques to improve model accuracy would enhance the reliability of our predictions. Further research could also delve into the potential impacts of sea level rise on specific regions and populations, informing targeted mitigation and adaptation strategies.

## V. REFERENCES

[1] Z. Wu and X. Cheng, "Impact Analysis of Sea Level Rising Problems Based on Mathematical Modeling," 2021 International Conference on Public Management and Intelligent Society (PMIS), Shanghai, China, 2021, pp. 400-403, doi: 10.1109/PMIS52742.2021.00097.

[2] R. O. Sinnott and Z. Guan, "Prediction of air pollution through machine learning approaches on the cloud", *2018 IEEE/ACM 5th International Conference on Big Data Computing Applications and Technologies (BDCAT)*, pp. 51-60, 2018.

[3] T. Peng et al., "Application of Grey Prediction Model Based on Python in Carbon Emission Prediction and Low-Carbon Economic Development Analysis," 2023 International Conference on Network, Multimedia and Information Technology (NMITCON), Bengaluru, India, 2023, pp. 1-6, doi: 10.1109/NMITCON58196.2023.10276259.