# GOVERNMENT COLLEGE OF TECHNOLOGY, COIMBATORE-13

## PROJECT NAME: Air Quality Assessment TN

**Team Members:**

1. Kokulavenkat K

2. Keerthana R

3. Kayalvizhi T

4. Karunasri A

## Phase 4: PROJECT DEVELOPMENT LEVEL 2

**PROJECT DESCRIPTION**:

Gather air quality data from monitoring stations across Tamil Nadu. This data may include measurements of pollutants like PM2.5, PM10, nitrogen dioxide (NO2), sulphur dioxide (SO2), ozone (O3), carbon monoxide (CO), and other relevant parameters.

Ensure the data is structured, accurate, and includes timestamps for each measurement. Aggregate data if multiple stations provide data for the same location. To visualise the data including line charts, bar graphs, heatmaps, and scatter plots and compare different pollutants. Geographic visualizations using tools like GIS software or Python libraries like Folium to map air quality across different regions in Tamil Nadu.

Consider using color-coding to represent air quality levels and station locations on the map. Calculate AQI based on the collected data and the relevant formula for Tamil Nadu. To visualise the air quality find how different pollutants relate to each other and apply time series analysis techniques, autoregressive integrated moving average (ARIMA) or prophet for forecasting air quality.

## LOADING AN DATASET:

- Load a dataset into a pandas DataFrame using `pd.read_csv()` for CSV files or `pd.read_excel()` for Excel files.
- Ensure that 'data' is correctly loaded before applying the provided code for preprocessing.

## DATA CONVERSION:

Convert the 'Sampling_Date' column to a datetime format, handling errors by using NaT.



## DATA FORMATTING:

Create a new 'formatted_date' column with a specific date format (YYYY-MM-DD).

```
[17] data['formatted_date'] = data['Sampling_Date'].dt.strftime('%Y-%m-%d')
```

```
data.head()
```

| | Stn_Code | Sampling_Date | State | City/Town/Village/Area | Location_of_Monitoring_Station | Agency | Type_of_Location | SO2 | NO2 | RSPM/PM10 | PM_2.5 | formatted_date |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 38 | 2014-01-02 | Tamil Nadu | Chennai | Kathivakkam, Municipal Kalyana Mandapam, Chennai | Tamilnadu State Pollution Control Board | Industrial Area | 11.0 | 17.0 | 55.0 | NaN | 2014-01-02 |
| 1 | 38 | 2014-01-07 | Tamil Nadu | Chennai | Kathivakkam, Municipal Kalyana Mandapam, Chennai | Tamilnadu State Pollution Control Board | Industrial Area | 13.0 | 17.0 | 45.0 | NaN | 2014-01-07 |
| 2 | 38 | 2014-01-21 | Tamil Nadu | Chennai | Kathivakkam, Municipal Kalyana Mandapam, Chennai | Tamilnadu State Pollution Control Board | Industrial Area | 12.0 | 18.0 | 50.0 | NaN | 2014-01-21 |
| 3 | 38 | 2014-01-23 | Tamil Nadu | Chennai | Kathivakkam, Municipal Kalyana Mandapam, Chennai | Tamilnadu State Pollution Control Board | Industrial Area | 15.0 | 16.0 | 46.0 | NaN | 2014-01-23 |
| 4 | 38 | 2014-01-28 | Tamil Nadu | Chennai | Kathivakkam, Municipal Kalyana Mandapam, Chennai | Tamilnadu State Pollution Control Board | Industrial Area | 13.0 | 14.0 | 42.0 | NaN | 2014-01-28 |

## DATA AGGREGATION:

Aggregate data by 'City/Town/Village/Area' and 'Sampling_Date', calculating mean values for 'SO2' and 'NO2' for each group.

```
agg_df = data.groupby(['City/Town/Village/Area']).agg({'SO2': 'mean', 'NO2': 'mean'}).reset_index()
agg_df.head()
```

| | City/Town/Village/Area | SO2 | NO2 |
|---|---|---|---|
| 0 | Chennai | 13.014042 | 22.088442 |
| 1 | Coimbatore | 4.541096 | 25.325342 |
| 2 | Cuddalore | 8.965986 | 19.710884 |
| 3 | Madurai | 13.319728 | 25.768707 |
| 4 | Mettur | 8.429268 | 23.185366 |

```
[23] agg_df = data.groupby(['Sampling_Date']).agg({'SO2': 'mean', 'NO2': 'mean'}).reset_index()
agg_df.head()
```

| | Sampling_Date | SO2 | NO2 |
|---|---|---|---|
| 0 | 2014-01-02 | 12.272727 | 19.272727 |
| 1 | 2014-01-03 | 11.461538 | 24.153846 |
| 2 | 2014-01-04 | 9.000000 | 21.250000 |
| 3 | 2014-01-06 | 12.727273 | 22.000000 |
| 4 | 2014-01-07 | 12.230769 | 26.615385 |

## FILTERING UNWANTED DATA:

- Use double square brackets, like `[['NO2', 'SO2', 'City', 'Sampling_Date']]`.
- This creates a new DataFrame with only the selected columns.
- Useful for focusing on specific columns in analysis or data processing tasks.

```
[32] agg_df.columns
     Index(['City/Town/Village/Area', 'Sampling_Date', 'SO2', 'NO2'], dtype='object')

[33] agg_df.head()

     City/Town/Village/Area   Sampling_Date        SO2         NO2
  0               Chennai        2014-01-02    14.000000    18.250000
  1               Chennai        2014-01-03    12.800000    26.000000
  2               Chennai        2014-01-06    15.333333    19.333333
  3               Chennai        2014-01-07    14.666667    32.500000
  4               Chennai        2014-01-08    11.000000    16.000000
```

## ANALYSING AND VISUALISATION:

Creating a heatmap for air quality data can provide a clear visualization of the variations in SO2 and NO2 content across different cities and sampling dates. Since as mentioned that the city column has repeated city names, we may want to aggregate the data before creating the heatmap. Here's a step-by-step guide:
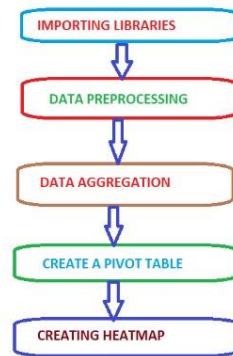
### Aggregate Data:

Calculate the average SO2 and NO2 content for each city on each sampling date. This will help in reducing repeated city names and provide a representative value for each city-date combination.
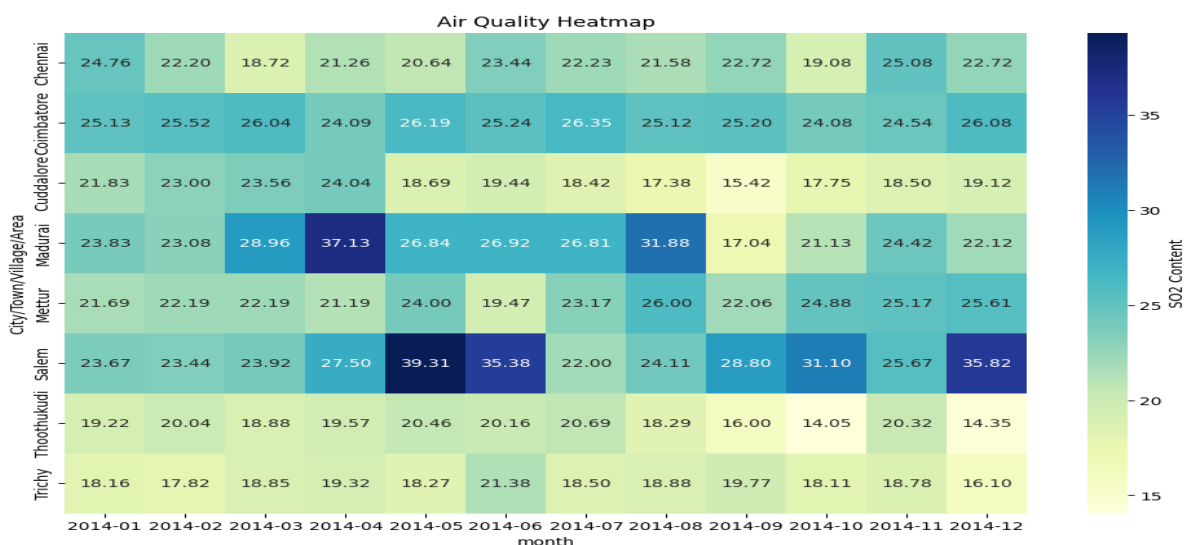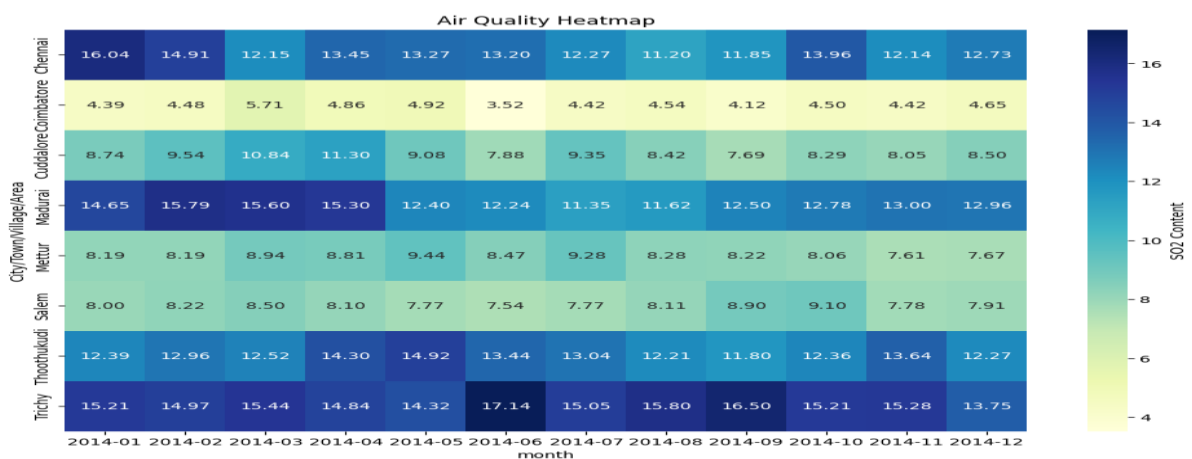
### Pivot Table:

Create a pivot table where cities are along the rows, sampling dates along the columns, and the values are either average SO2 or NO2 content.

### Heatmap:

Use a heatmap to visualize the data. You can use tools like Python with libraries such as Matplotlib and Seaborn or other data visualization tools like Tableau or Excel.

The code preprocesses air quality data by converting dates and creating a formatted date column. It then aggregates the data by calculating mean 'SO2' and 'NO2' values for different locations and dates. Finally, it visualizes the 'SO2' and 'NO2' levels across locations and dates using a heatmap with color-coded values and a color bar, making air quality trends easily interpretable.

```
[ ] import pandas as pd
    import seaborn as sns
    import matplotlib.pyplot as plt

    # Assuming your DataFrame is named 'df'
    # Aggregate data
    data['Sampling_Date'] = pd.to_datetime(data['Sampling_Date'], errors='coerce')

    data['formatted_date'] = data['Sampling_Date'].dt.strftime('%Y-%m-%d')
    data['month'] = data['Sampling_Date'].dt.to_period('M')
    agg_df = data.groupby(['City/Town/Village/Area', 'month']).agg({'SO2': 'mean', 'NO2': 'mean', 'RSPM/PM10': 'mean'}

    # Create a pivot table
    heatmap_data = agg_df.pivot('City/Town/Village/Area', 'month', 'RSPM/PM10')

    # Create heatmap
    plt.figure(figsize=(12, 8))
    sns.heatmap(heatmap_data, cmap='YlGnBu', annot=True, fmt=".2f", cbar_kws={'label': 'RSPM/PM10 Content'})
    plt.title('Air Quality Heatmap')
    plt.show()
```

This code snippet performs the following actions:

- Imports necessary libraries: `pandas`, `seaborn`, and `matplotlib.pyplot`.
- Converts a date column in a DataFrame ('data') to a datetime format.
- Extracts the month and creates a formatted date column.
- Aggregates data by grouping it based on 'City/Town/Village/Area' and 'month', calculating the mean values of 'SO2', 'NO2', and 'RSPM/PM10'.
- Creates a pivot table from the aggregated data with 'City/Town/Village/Area' as rows and 'month' as columns, with 'RSPM/PM10' values in the cells.
- Plots a heatmap using Seaborn to visualize the 'RSPM/PM10' content in the DataFrame. It includes annotations with two decimal places and a colorbar with a label



Air Quality Heatmap