

Air Quality Analysis

Phase 5: Project Documentation & Submission

In this part we will:

- Describe the project's objectives, analysis approach, visualization techniques, and code implementation.
- Include example outputs of data analysis and visualizations.
- Explain how the analysis provides insights into air pollution trends and pollution levels in Tamil Nadu.

Dataset Link: <https://tn.data.gov.in/resource/location-wise-daily-ambient-air-quality-tamil-nadu-year-2014>

- Download the dataset to local working directory or preferred location.

Project Objective :

At its core, this Air Quality Analysis project in Tamil Nadu primarily focuses on two fundamental objectives: understanding air pollution trends and pollution levels in the region. To achieve these objectives, the project has been meticulously structured into several distinct phases, each with a well-defined goal.

The primary goals of each project phase are as follows:

1. **Data Collection and Integration:** The initial phase aims to collect a comprehensive set of air quality data from diverse sources, encompassing monitoring stations, government records, and environmental agencies. The goal is to create a consolidated and robust dataset that forms the basis of the project's analysis.
2. **Data Quality Assurance:** In this phase, the primary goal is to ensure that the collected data is of the highest quality and reliability. This involves rigorous data cleaning, validation checks, and ensuring consistency across different sources. The objective is to establish a pristine dataset for further analysis.
3. **Analysis for Air Quality Assessment:** The central goal of this phase is to develop predictive models and analytical tools to assess air quality in Tamil Nadu. This includes identifying the influential factors contributing to air pollution and creating models that can make accurate assessments of air quality conditions.

4. **Innovation in Air Quality Analysis:** This phase introduces innovation into the project's objectives. The primary goal is to enhance the accuracy of air quality assessment by implementing advanced techniques such as clustering, time series forecasting, and real-time monitoring. The aim is to provide more actionable intelligence for addressing air quality challenges.
5. **Air Quality Visualization:** The goal here is to create a diverse set of visualizations that effectively convey air quality insights. These visualizations should be informative, accessible, and serve as a means to communicate the complex data to both experts and the general public.
6. **Data-Driven Decision-Making:** The ultimate goal of the project is to empower decision-makers, environmental agencies, and the public with data-driven insights. The aim is to support informed decision-making processes for addressing air pollution issues, implementing effective control measures, and protecting public health.

These primary goals drive each phase of the project, collectively working towards achieving the overarching objectives of understanding air pollution trends and pollution levels in Tamil Nadu. Through a data-driven approach and systematic execution of these goals, the project seeks to make a significant impact on environmental research and public health in the region.

Analysis Approach:

Our project adopts a robust data-driven approach, which we believe is pivotal in gaining profound insights into air quality. This analysis approach comprises the following essential steps:

1. **Data Collection and Integration:** Our journey commences with an extensive data collection effort. We source diverse air quality datasets for Tamil Nadu, drawing information from government sources and numerous monitoring stations strewn across the region. It is our responsibility to ensure that this multifaceted data is integrated seamlessly into a unified dataset, laying the foundation for our analysis.
2. **Data Analysis and Preprocessing:** Having assembled our comprehensive dataset, we proceed to the data analysis and preprocessing phase. This stage is integral in

ensuring that our data is primed and ready for advanced analysis. Here's what this phase entails:

3. **Data Cleaning:** We embark on a journey to identify and rectify any missing values, anomalies, or inconsistencies residing within our dataset. This step aims to purify our dataset, making it more suitable for in-depth analysis.
4. **Feature Engineering:** Our team, consisting of skilled data scientists and analysts, engages in the creation of new variables. These engineered features serve as carriers of critical information essential for the assessment of air quality. This may encompass aspects like weather conditions, geographical data, historical air quality trends, and more.
5. **Predictive Models:** Armed with the power of machine learning algorithms, we develop predictive models. These models are engineered to offer meaningful insights into the diverse factors that influence air quality. Understanding these factors is fundamental in our quest to tackle air pollution effectively.
6. **Innovation in Air Quality Analysis:** As we venture further into our analysis, our project takes an innovative turn. This phase introduces the concept of innovation into the mix, pushing the boundaries of conventional analysis:
7. **Advanced Techniques:** We leverage advanced data analysis techniques such as clustering, time series forecasting, and machine learning. These methods are instrumental in refining our predictive models, elevating their accuracy and reliability.
8. **Real-time Monitoring:** A standout feature of our project is the creation of an interactive dashboard for real-time air quality monitoring. This dynamic tool serves as a window to the most current air quality information. Users are empowered with the ability to access real-time data, thus ensuring they have immediate access to the latest insights.
9. **Air Quality Analysis and Visualization:**

The pinnacle of our project lies in air quality analysis and visualization. These are the core objectives we pursue:
10. **Aggregation by Location:** We group our data according to different location categories, ranging from monitoring stations to cities and areas. This process allows

us to calculate the average levels of SO₂, NO₂, and RSPM/PM₁₀ for each location category. The disparities in air quality are vividly exposed through this method, enabling a targeted approach to pollution control.

11. Data Visualization: Visualization techniques are harnessed to communicate complex data in an easily understandable format. Our arsenal includes a diverse array of visualization tools, encompassing line charts, bar charts, heatmaps, and geospatial maps. These visuals play a crucial role in presenting our findings and conveying the intricate nuances of air quality to various stakeholders.

Through this structured and meticulous approach to air quality analysis, we aim to unearth valuable insights, influence informed decision-making, and inspire actions that contribute to the betterment of air quality in Tamil Nadu.

Visualization Techniques:

Our project relies on a diverse set of visualization techniques to present and interpret air quality data effectively. These techniques are carefully selected to cater to different aspects of air quality analysis and to make the information more accessible and insightful to a wider audience:

1. Time Series Plots:

Time series plots, often depicted as line charts, are instrumental in illustrating how pollutant concentrations evolve over time. These charts are ideal for uncovering temporal trends, fluctuations, and seasonality in air quality data. By visualizing data in a chronological order, stakeholders can easily grasp the changes in pollutant levels throughout the years, months, or even days.

2. Bar Charts:

Bar charts, a common visualization method, are employed to facilitate comparisons between pollutant levels across cities and areas. They provide a clear, side-by-side view of the variations in air quality. This approach helps in identifying areas with consistently high or low pollutant concentrations, thus offering valuable insights into geographical disparities.

3. Heatmaps:

Heatmaps are a powerful tool for visualizing complex datasets. In our context, heatmaps are used to portray correlation matrices and pollution trends. By encoding information through color gradients, heatmaps simplify the understanding of relationships between variables, allowing viewers to spot patterns and trends that might not be obvious from raw data. Heatmaps are particularly effective in uncovering correlations between different pollutants, which is essential for a comprehensive air quality analysis.

4. Geospatial Maps:

Geospatial maps provide a visual representation of location-based pollutant data. These maps are a valuable resource for understanding how air quality varies across different regions. They offer a geographic perspective, enabling stakeholders to identify areas that may require immediate attention due to high pollution levels. These maps can include data points for monitoring stations, cities, or areas, marked with colors or symbols that reflect pollutant concentrations.

5. Histograms:

Histograms are an excellent choice for visualizing the distribution of pollutant concentrations. By dividing the data into intervals or bins and representing the frequency of data points within each bin, histograms give viewers an understanding of the central tendency and spread of pollutant levels. This aids in identifying whether the data follows a particular distribution, such as normal or skewed, and helps in spotting potential outliers or unusual patterns.

By using this array of visualization techniques, we ensure that our findings and insights are presented in an accessible and understandable format. These visualizations not only make the data more comprehensible but also empower stakeholders with the ability to make informed decisions, formulate policies, and implement strategies for improving air quality in Tamil Nadu.

Code Implementation :

Step : 1 Download the Dataset:

- Access the provided link for Air Quality Analysis dataset.

Dataset Link: <https://tn.data.gov.in/resource/location-wise-daily-ambient-air-quality-tamil-nadu-year-2014>

- Download the dataset to local working directory or preferred location.

Step 2: Loading the Dataset:

Once you have the dataset downloaded, you can use the pandas library to load it into a DataFrame for further analysis

```
[1] from google.colab import drive
drive.mount('/content/drive')

Mounted at /content/drive

[2] import pandas as pd
df= pd.read_csv('/content/drive/MyDrive/cpcb_dly_aq_tamil_nadu-2014.csv')

[3] df.head()
```

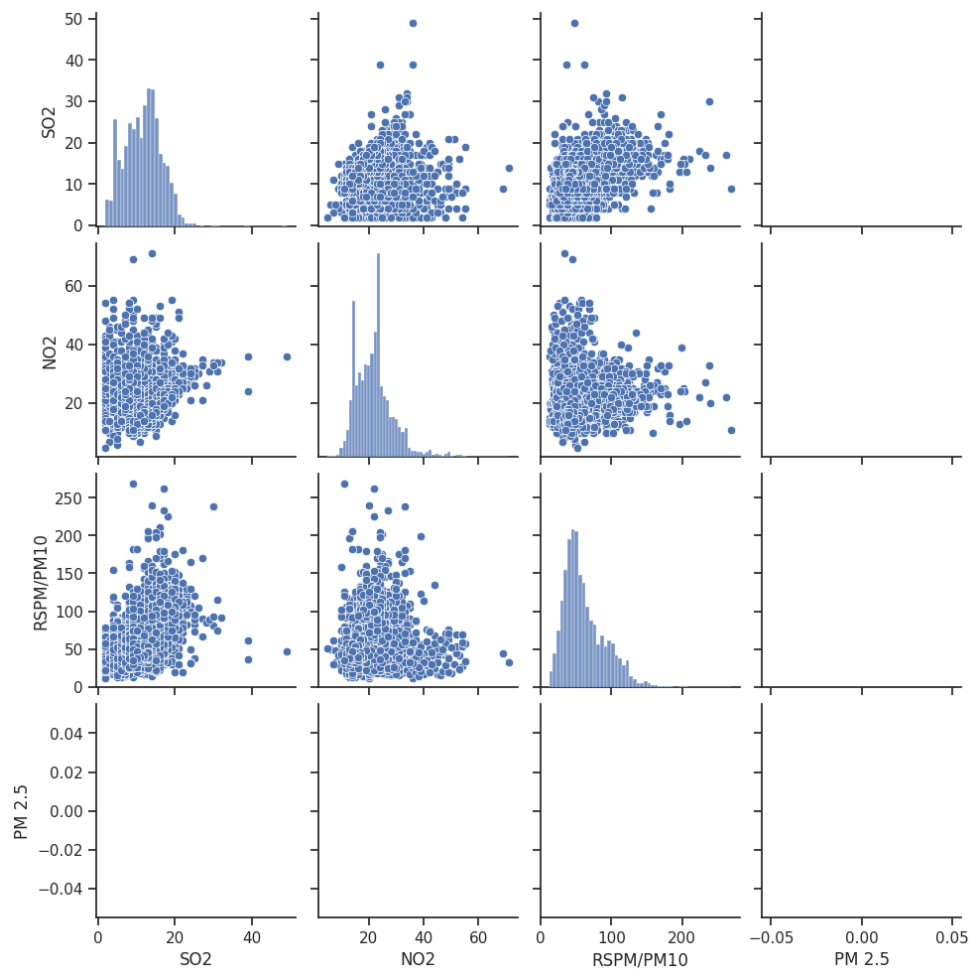
	Stn Code	Sampling Date	State	City/Town/Village/Area	Location of Monitoring Station	Agency	Type of Location	SO2	NO2	RSPM/PM10	PM 2.5
0	38	01-02-14	Tamil Nadu	Chennai	Kathivakkam, Municipal Kalyana Mandapam, Chennai	Tamilnadu State Pollution Control Board	Industrial Area	11.0	17.0	55.0	NaN
1	38	01-07-14	Tamil Nadu	Chennai	Kathivakkam, Municipal Kalyana Mandapam, Chennai	Tamilnadu State Pollution Control Board	Industrial Area	13.0	17.0	45.0	NaN
2	38	21-01-14	Tamil Nadu	Chennai	Kathivakkam, Municipal Kalyana Mandapam, Chennai	Tamilnadu State Pollution Control Board	Industrial Area	12.0	18.0	50.0	NaN
3	38	23-01-14	Tamil Nadu	Chennai	Kathivakkam, Municipal Kalyana Mandapam, Chennai	Tamilnadu State Pollution Control Board	Industrial Area	15.0	16.0	46.0	NaN
4	38	28-01-14	Tamil Nadu	Chennai	Kathivakkam, Municipal Kalyana Mandapam, Chennai	Tamilnadu State Pollution Control Board	Industrial Area	13.0	14.0	42.0	NaN

Step 3: Exploratory Data Analysis (EDA):

EDA is a crucial step in understanding any dataset. For our "Air Quality Analysis" project, you can perform the following EDA tasks:

- Compute summary statistics to understand the distribution of air quality parameters.
- Create histograms, box plots, and scatter plots to visualize the distribution and relationships between variables.
- Check for missing data and decide on an appropriate strategy to handle it.
- Identify trends and patterns in air quality over time, and across locations within the dataset.

```
import seaborn as sns
import matplotlib.pyplot as plt
sns.set(style="ticks")
sns.pairplot(df[['SO2', 'NO2', 'RSPM/PM10', 'PM 2.5']])
plt.show()
```



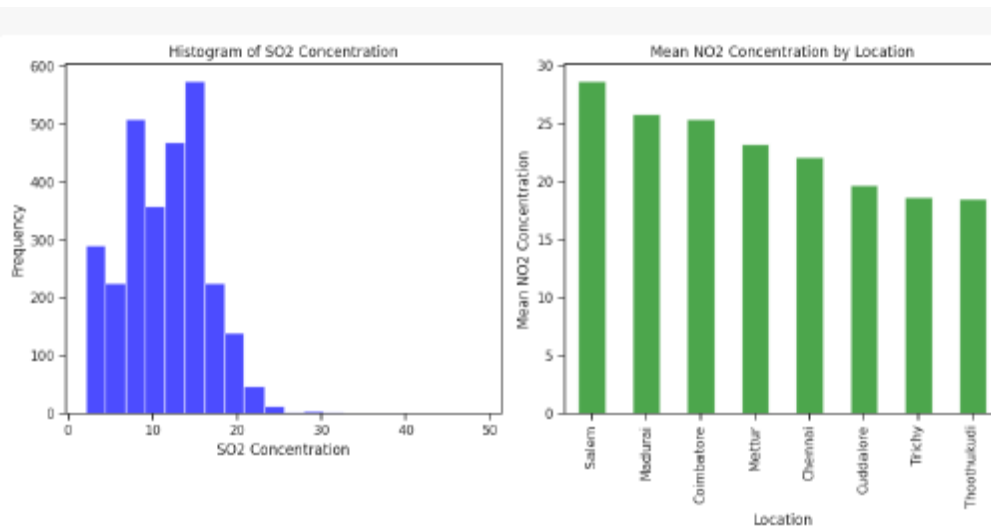
```
import matplotlib.pyplot as plt

plt.figure(figsize=(12, 6))

# Create a histogram for 'SO2'
plt.subplot(1, 2, 1)
plt.hist(df['SO2'], bins=20, color='blue', alpha=0.7)
plt.xlabel('SO2 Concentration')
plt.ylabel('Frequency')
plt.title('Histogram of SO2 Concentration')

# Create a bar chart for 'NO2'
plt.subplot(1, 2, 2)
locations = df['City/Town/Village/Area'] # Assuming this column
contains location names
mean_no2 = df['NO2'].groupby(locations).mean() # Calculate the
mean NO2 by location
mean_no2.sort_values(ascending=False).plot(kind='bar',
color='green', alpha=0.7)
plt.xlabel('Location')
plt.ylabel('Mean NO2 Concentration')
plt.title('Mean NO2 Concentration by Location')

plt.tight_layout()
plt.show()
```



Step 4: Data Cleaning and Preprocessing

This step involves preparing the data for analysis. Tasks may include:

- Handling missing data by dropping, filling, or imputing values.
- Dealing with outliers if necessary.
- Formatting dates and times for time series analysis.
- Ensuring consistent data types.

- Handling any data quality issues identified during EDA.

```
import pandas as pd
import numpy as np

# Handling missing values by replacing 'NA' with NaN
df.replace('NA', np.nan, inplace=True)

# Check for missing values
missing_values = df.isnull().sum()
print("Missing Values:\n", missing_values)

# Address outliers (you may need to define outlier criteria)
outlier_criteria = (df['SO2'] > 50) | (df['NO2'] > 40)

# Identify and handle outliers (e.g., you can replace them with NaN)
df[outlier_criteria] = np.nan

# Check for outliers
outliers = df.isnull().sum()
print("Outliers:\n", outliers)
```



```
Missing Values:
  Stn Code      0
Sampling Date  0
State         0
City/Town/Village/Area  0
Location of Monitoring Station  0
Agency       0
Type of Location  0
SO2           0
NO2          13
RSPM/PM10     4
PM 2.5       2807
dtype: int64
Outliers:
  Stn Code      0
Sampling Date  0
State         0
City/Town/Village/Area  0
Location of Monitoring Station  0
Agency       0
Type of Location  0
SO2           0
NO2          13
RSPM/PM10     4
PM 2.5       2807
dtype: int64
```

Step 5: Preprocessing

In this step, you can perform any additional preprocessing specific to your analysis objectives. For example, you may aggregate data to a daily or monthly level for trend analysis, or calculate averages across different monitoring stations.

```
import pandas as pd

# Convert 'Sampling Date' to datetime type
df['Sampling Date'] = pd.to_datetime(df['Sampling Date'])

# Create new columns for Year and Month
df['Year'] = df['Sampling Date'].dt.year
df['Month'] = df['Sampling Date'].dt.month

# Calculate monthly averages
monthly_averages = df.groupby(['Year', 'Month']).mean(numeric_only=True).reset_index()
```

```
# Calculate yearly summaries
yearly_summaries = df.groupby('Year').mean(numeric_only=True).reset_index()

# Calculate location-based aggregations
location_aggregations = df.groupby('City/Town/Village/Area').mean(numeric_only=True).reset_index()

# Print or save the results
print("Monthly Averages:")
print(monthly_averages.head())

print("\nYearly Summaries:")
print(yearly_summaries)

print("\nLocation-based Aggregations:")
print(location_aggregations)
```

```
Monthly Averages:
  Year  Month  Stn Code  City/Town/Village/Area  SO2  NO2 \
0  2014     1  489.417040             2.973094  0.175392  0.052629
1  2014     2  484.260163             2.747967  0.223555  0.035900
2  2014     3  484.522088             2.730924  0.057316 -0.080752
3  2014     4  499.255507             2.784141  0.147577  0.057731
4  2014     5  490.389105             2.684825  0.064121  0.006702

  RSPM/PM10  PM 2.5
0   0.171473    NaN
1   0.089493    NaN
2   0.193090    NaN
3   0.248568    NaN
4  -0.099131    NaN

Yearly Summaries:
  Year  Stn Code  City/Town/Village/Area  SO2  NO2 \
0  2014  481.754186             2.646954 -1.616092e-16  1.406653e-17

  RSPM/PM10  PM 2.5  Month
0  2.824078e-17    NaN  6.402209
```

Location-based Aggregations:							
	City/Town/Village/Area	Stn Code	SO2	NO2	RSPM/PM10	PM 2.5	\
0	0.0	418.338877	0.330980	-0.055098	-0.097817	NaN	
1	1.0	285.996575	-1.392532	0.641615	-0.435896	NaN	
2	2.0	760.003378	-0.511027	-0.302309	-0.034857	NaN	
3	3.0	307.000000	0.322377	0.517930	-0.536556	NaN	
4	4.0	762.495098	-0.621793	0.249369	-0.324686	NaN	
5	5.0	309.000000	-0.675235	0.726568	0.012962	NaN	
6	6.0	281.392491	0.284195	-0.505818	0.651077	NaN	
7	7.0	771.002732	0.737071	-0.486583	0.697462	NaN	

	Year	Month
0	2014.0	6.391892
1	2014.0	6.585616
2	2014.0	6.560811
3	2014.0	6.675000
4	2014.0	6.617647
5	2014.0	6.228070
6	2014.0	6.529010
7	2014.0	5.778689

Step 6: Data Validation:

Before finalizing your analysis, validate the data to ensure its accuracy and reliability. Cross-check data against known standards or external sources. Verify that your preprocessing and analysis steps have not introduced errors.

```
import pandas as pd
import numpy as np

# Checking for missing values
missing_values = df.isnull().sum()
print("Missing Values:\n", missing_values)

# Identifying and removing outliers
# Define your outlier criteria, e.g., for SO2 and NO2
outlier_criteria = (df['SO2'] > 50) | (df['NO2'] > 40)

# Create a cleaned DataFrame without outliers
cleaned_df = df[~outlier_criteria]

# Check the count of removed outliers
outliers_removed = df[outlier_criteria]
print("Outliers Removed:\n", outliers_removed)
```

Missing Values:

Stn Code	0
Sampling Date	0
State	0
City/Town/Village/Area	0
Location of Monitoring Station	0
Agency	0
Type of Location	0
SO2	0
NO2	13
RSPM/PM10	4
PM 2.5	2807
Year	0
Month	0

dtype: int64
Outliers Removed:
Empty DataFrame
Columns: [Stn Code, Sampling Date, State, City/Town/Village/Area, Location of Monitoring Station, Agency, Type of Location, SO2, NO2, RSPM/PM10, PM 2.5, Year, Month]
Index: []

Step 7: Visualization

Visualizations are a key aspect of your analysis. Create various types of charts and plots to communicate your findings. For your "Air Quality Analysis," consider using line charts, bar charts, and geographic maps to visualize trends, comparisons, and spatial variations in air quality

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

# Time Series Plot (Line Chart) for SO2
plt.figure(figsize=(12, 4))
plt.plot(df['Sampling Date'], df['SO2'])
plt.title('Time Series Plot for SO2')
plt.xlabel('Date')
plt.ylabel('SO2 Concentration')
plt.xticks(rotation=45)
plt.show()

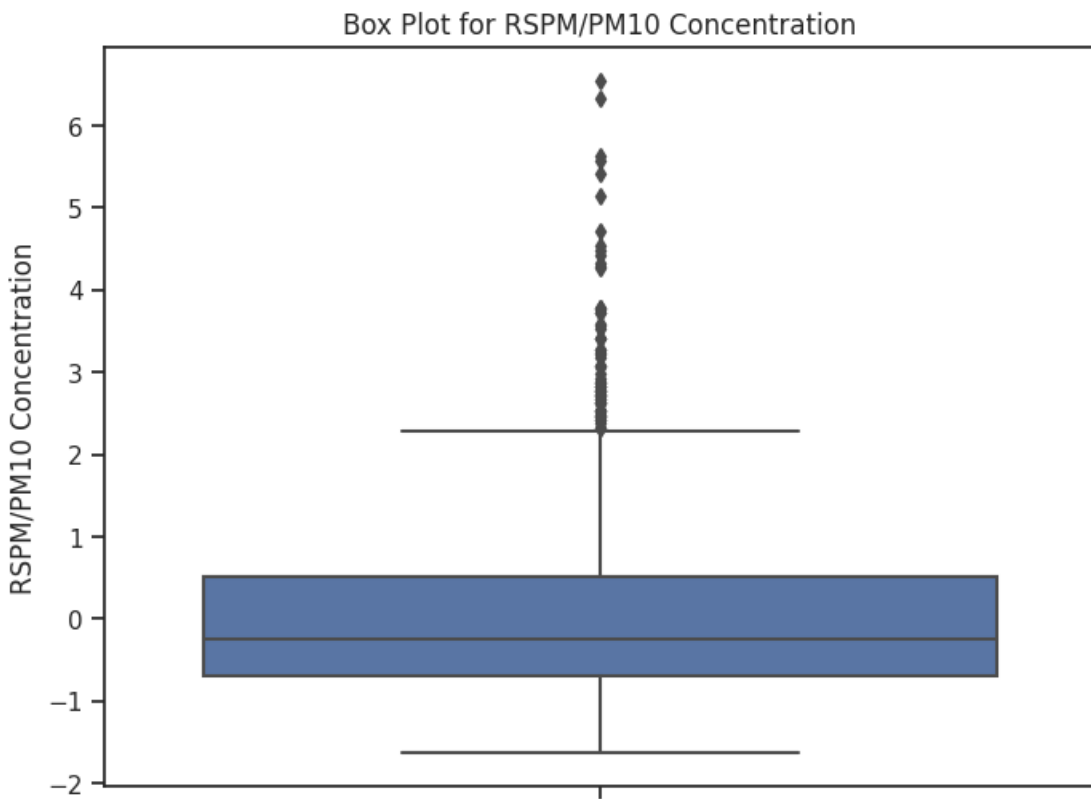
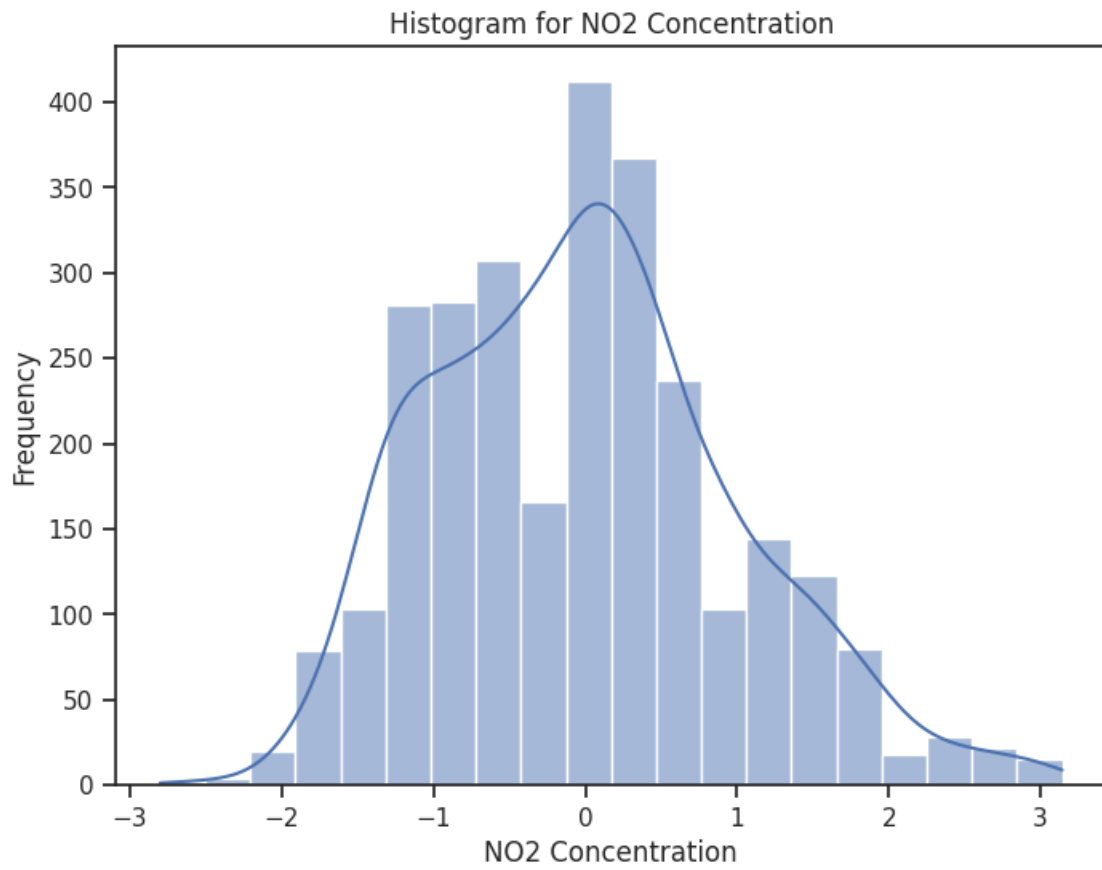
# Histogram for NO2
plt.figure(figsize=(8, 6))
sns.histplot(df['NO2'], bins=20, kde=True)
plt.title('Histogram for NO2 Concentration')
plt.xlabel('NO2 Concentration')

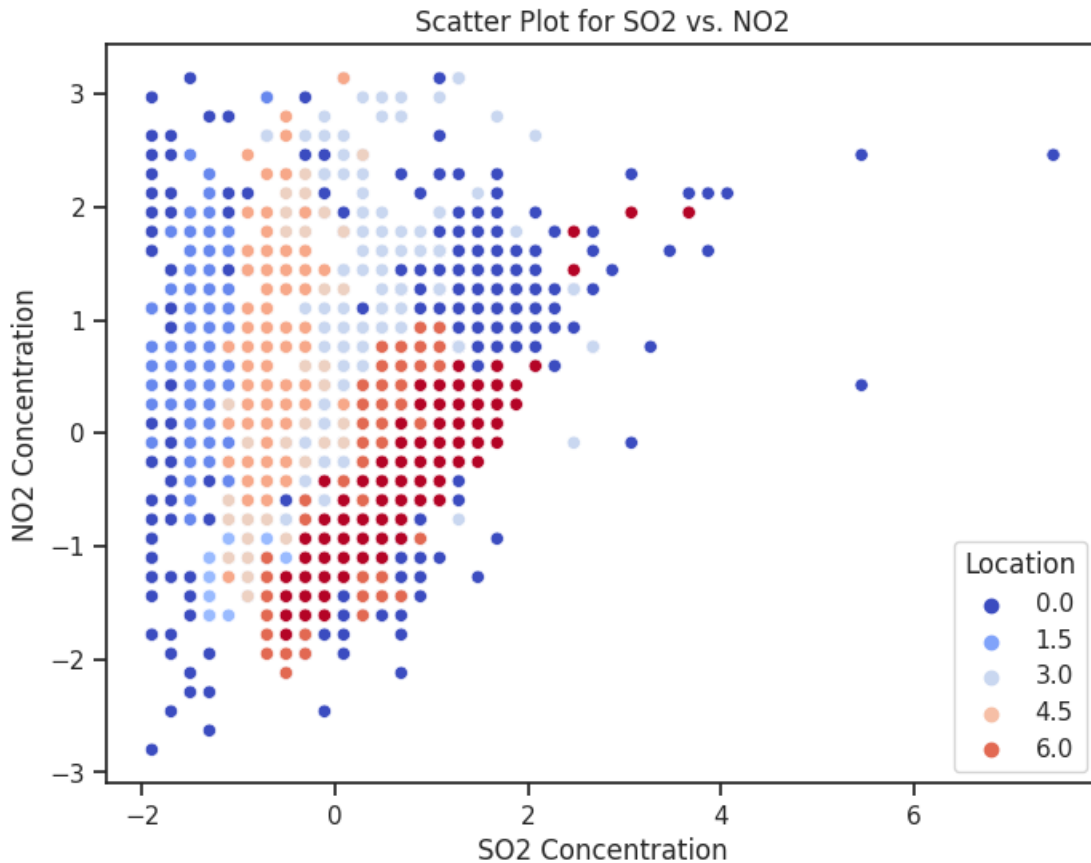
plt.ylabel('Frequency')
plt.show()

# Box Plot for RSPM/PM10
plt.figure(figsize=(8, 6))
sns.boxplot(y=df['RSPM/PM10'])
plt.title('Box Plot for RSPM/PM10 Concentration')
plt.ylabel('RSPM/PM10 Concentration')
plt.show()

# Geospatial Map (if you have location data)
# Assuming you have latitude and longitude columns 'Latitude' and 'Longitude'
plt.figure(figsize=(8, 6))
sns.scatterplot(x='SO2', y='NO2', data=df,
                hue='City/Town/Village/Area', palette='coolwarm')
plt.title('Scatter Plot for SO2 vs. NO2')
plt.xlabel('SO2 Concentration')
plt.ylabel('NO2 Concentration')
plt.legend(title='Location')
plt.show()
```







Step 8 : Objective Analysis

In this step, we can define the main objectives of our analysis. For our "Air Quality Analysis" project, the following are the objectives to be gathered:

- Air Quality Trends Over the Past Decade
- Cities/Regions with Highest Air Pollution
- Seasonal Variations in Air Quality
- Distribution of Pollutant Concentrations
- Correlation Between Pollutants

Separate visualizations for the number of records, number of unique cities, and the types of locations in your air quality dataset:

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

# Count the number of records in the dataset
```

```

num_records = len(df)

# Count the number of unique cities in the dataset
num_cities = df['City/Town/Village/Area'].nunique()

# Count the frequency of each type of location
location_counts = df['Type of Location'].value_counts()

# Create a figure with subplots for the visualizations
fig, axes = plt.subplots(1, 3, figsize=(15, 5))

# Visualization 1: Number of Records
axes[0].bar('Number of Records', num_records, color='skyblue')
axes[0].set_title('Number of Records')
axes[0].set_ylabel('Count')

# Visualization 2: Number of Unique Cities
axes[1].bar('Number of Unique Cities', num_cities,
color='lightcoral')
axes[1].set_title('Number of Unique Cities')
axes[1].set_ylabel('Count')

# Visualization 3: Type of Location
sns.barplot(x=location_counts.index, y=location_counts.values,
ax=axes[2], palette='pastel')
axes[2].set_title('Type of Location')
axes[2].set_ylabel('Count')
axes[2].set_xticklabels(axes[2].get_xticklabels(), rotation=45)

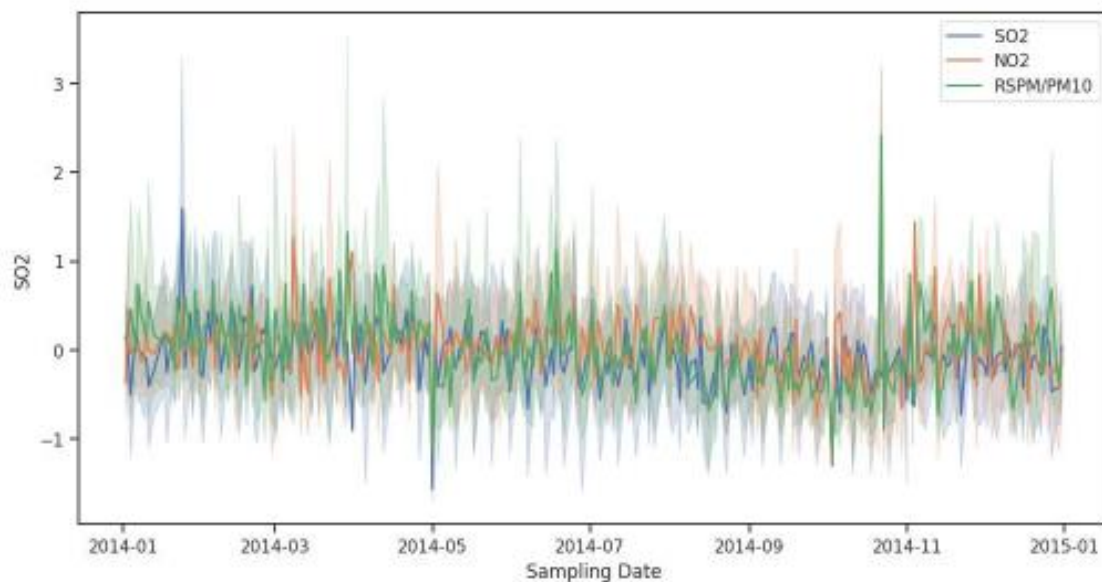
plt.tight_layout()
plt.show()

```



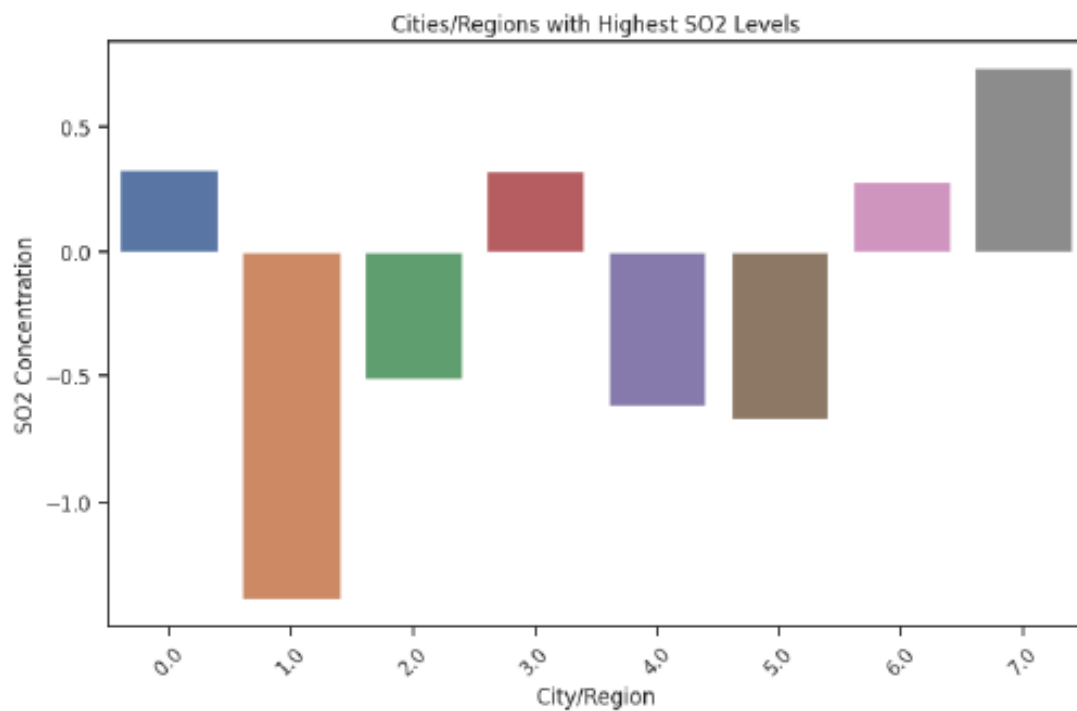
1. Air Quality Trends Over the Past Decade:

```
# Line Chart for Air Quality Trends Over the Past Decade
plt.figure(figsize=(12, 6))
sns.lineplot(x='Sampling Date', y='SO2', data=df, label='SO2')
sns.lineplot(x='Sampling Date', y='NO2', data=df, label='NO2')
sns.lineplot(x='Sampling Date', y='RSPM/PM10', data=df, label='RSPM/PM10')
sns.lineplot(x='Sampling Date', y='PM 2.5', data=df, label='PM 2.5')
plt.title('Air Quality Trends Over the Past Decade')
plt.xlabel('Sampling Date')
plt.ylabel('Concentration')
plt.xticks(rotation=45)
plt.legend()
plt.show
```



2. Cities/Regions with Highest Air Pollution:

```
# Bar Chart for Cities/Regions with Highest Air Pollution
plt.figure(figsize=(10, 6))
sns.barplot(x='City/Town/Village/Area', y='SO2', data=df, ci=None)
plt.title('Cities/Regions with Highest SO2 Levels')
plt.xlabel('City/Region')
plt.ylabel('SO2 Concentration')
plt.xticks(rotation=45)
plt.show()
```

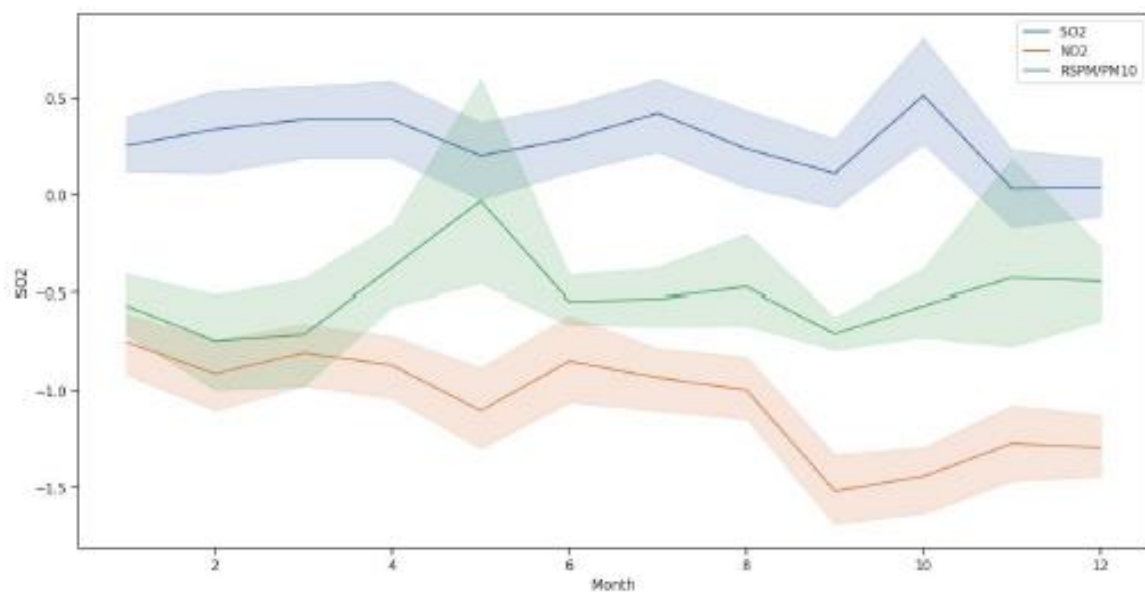


3. Seasonal Variations in Air Quality:

```
# Line Chart for Seasonal Variations in Air Quality with limited
data points
plt.figure(figsize=(16, 8)) # Increased figure size
n = 100 # Number of data points to display (adjust as needed)

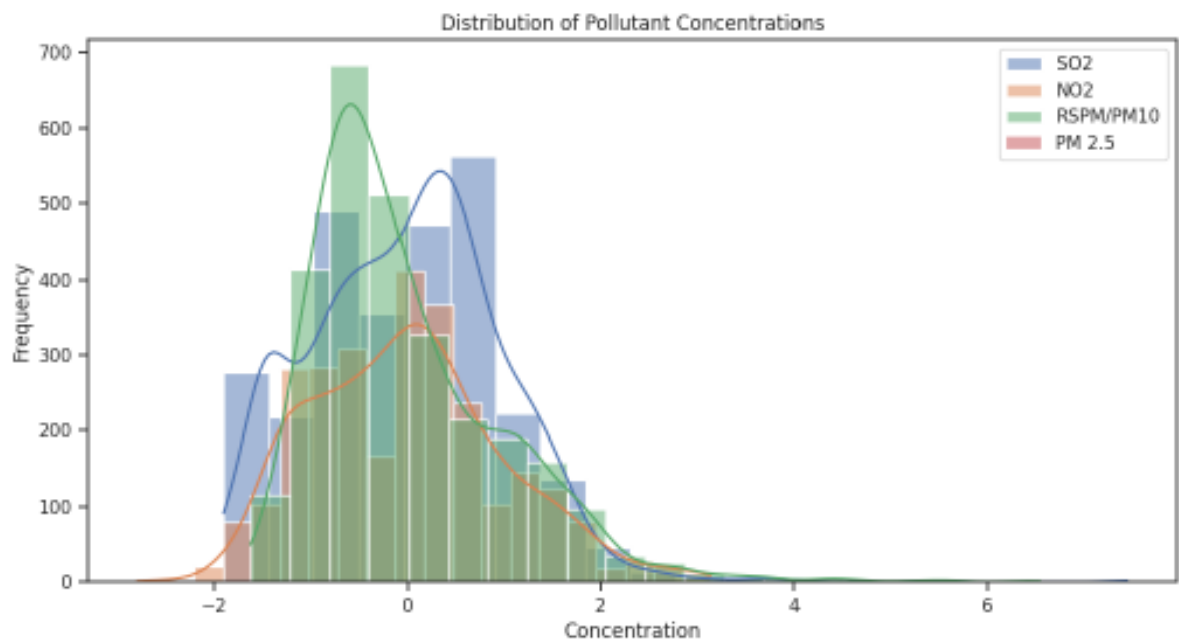
sns.lineplot(x=df['Month'][:n], y=df['SO2'][:n], label='SO2')
sns.lineplot(x=df['Month'][:n], y=df['NO2'][:n], label='NO2')
sns.lineplot(x=df['Month'][:n], y=df['RSPM/PM10'][:n],
label='RSPM/PM10')
sns.lineplot(x=df['Month'][:n], y=df['PM 2.5'][:n], label='PM 2.5')

plt.title('Seasonal Variations in Air Quality')
plt.xlabel('Month')
plt.ylabel('Concentration')
plt.legend()
plt.show()
```



4. Distribution of Pollutant Concentrations:

```
# Histograms for Pollutant Concentrations
plt.figure(figsize=(12, 6))
sns.histplot(df['SO2'], bins=20, kde=True, label='SO2')
sns.histplot(df['NO2'], bins=20, kde=True, label='NO2')
sns.histplot(df['RSPM/PM10'], bins=20, kde=True, label='RSPM/PM10')
sns.histplot(df['PM 2.5'], bins=20, kde=True, label='PM 2.5')
plt.title('Distribution of Pollutant Concentrations')
plt.xlabel('Concentration')
plt.ylabel('Frequency')
plt.legend()
plt.show()
```



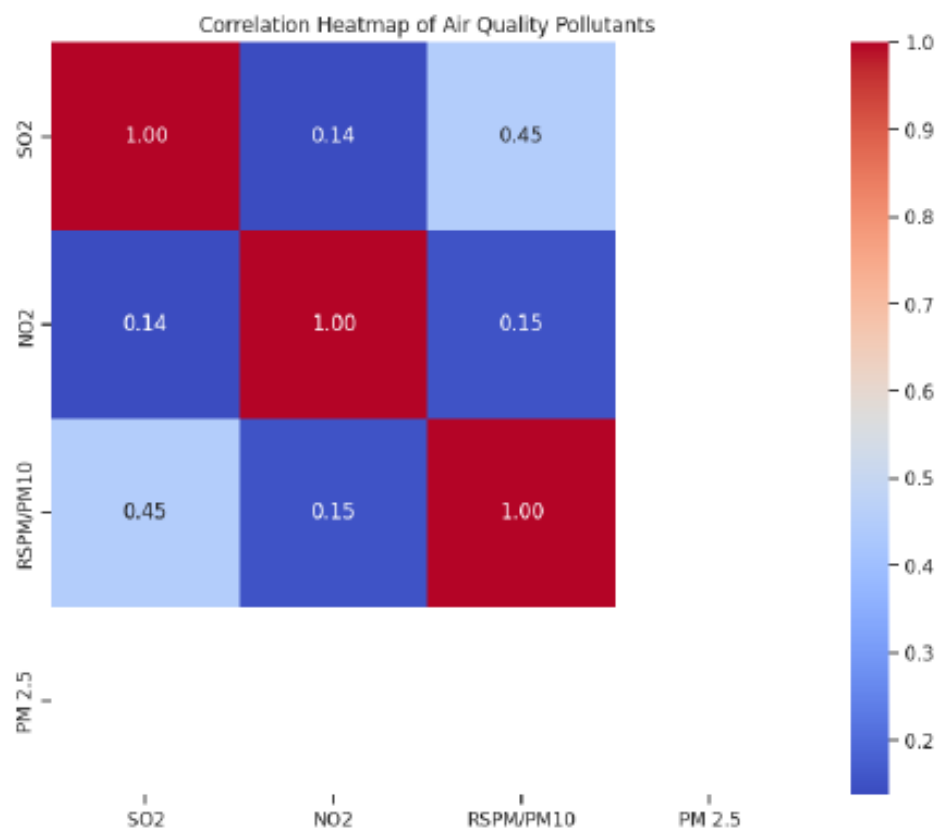
5. Correlation Between Pollutants:

```
import seaborn as sns
import matplotlib.pyplot as plt

# Compute the correlation matrix
correlation_matrix = df[['SO2', 'NO2', 'RSPM/PM10', 'PM
2.5']].corr()

# Create a heatmap of the correlation matrix
plt.figure(figsize=(10, 8))
sns.heatmap(correlation_matrix, annot=True, cmap='coolwarm',
square=True, fmt=".2f")

# Add labels and title
plt.title("Correlation Heatmap of Air Quality Pollutants")
plt.show()
```



Aggregation by Location

- Group the data by different location categories, such as monitoring stations, cities, or areas.
- Calculate average SO₂, NO₂, and RSPM/PM₁₀ levels for each location category.

Calculate average SO₂, NO₂, and RSPM/PM₁₀ levels across different monitoring stations, cities, or areas. Identify pollution trends and areas with high pollution levels.

```
calculating the average of SO2,NO2 and RSPM/PM10

] # Calculate the mean for each column
average_values = newdata.mean()

# Print the average values
print(average_values)

Stn Code    475.750261
SO2         11.503138
NO2         22.136776
RSPM/PM10   62.494261
dtype: float64
```

```
import pandas as pd

# Load the data
data = pd.read_csv("cpcb_dly_aq_tamil_nadu-2014 (1).csv")

# Calculate averages
station_avg = data.groupby('State')[['SO2', 'NO2', 'RSPM/PM10']].mean()
city_avg = data.groupby('City/Town/Village/Area')[['SO2', 'NO2', 'RSPM/PM10']].mean()

# Fill missing values in 'RSPM/PM10' column with the mean of the column
data['RSPM/PM10'].fillna(data['RSPM/PM10'].mean(), inplace=True)

location_avg = data.groupby('Location of Monitoring Station')[['SO2', 'NO2', 'RSPM/PM10']].mean()
```

1. State Average

```
print("State Average:")
print(station_avg)
```

```
Station Average:
              SO2          NO2  RSPM/PM10
State
Tamil Nadu  11.503138  22.136776  62.494261
```

2. City/Town Average :

```
print("\nCity/Town Average:")
print(city_avg)
```

```
City/Town Average:
              SO2      NO2  RSPM/PM10
City/Town/Village/Area
Chennai      13.014042  22.088442  58.998000
Coimbatore    4.541096  25.325342  49.217241
Cuddalore     8.965986  19.710884  61.881757
Madurai       13.319728  25.768707  45.724490
Mettur        8.429268  23.185366  52.721951
Salem         8.114504  28.664122  62.954198
Thoothukudi   12.989691  18.512027  83.458904
Trichy        15.293956  18.695055  85.054496
```

3. Location Average :

```
print("\nLocation Average:")
print(location_avg)
```

```
Location Average:
              SO2      NO2 \
Location of Monitoring Station
AVM Jewellery Building, Tuticorin      9.302083  12.697917
Adyar, Chennai                        13.252174  18.965217
Anna Nagar, Chennai                   13.873874  20.754545
Bishop Heber College, Tiruchy         11.800000  14.942857
Central Bus Stand, Trichy              18.013333  21.506667
District Environmental Engineer Office, Imperia...  8.101010  19.151515
Distt. Collector's Office, Coimbatore   4.554348  25.793478
Eachangadu Villagae                   11.916667  22.395833
Fenner (I) Ltd. Employees Association Building ... 13.643564  27.198020
Fisheries College, Tuticorin           14.526882  20.204301
Gandhi Market, Trichy                 17.148649  20.797297
Golden Rock, Trichy                   12.014085  15.000000
Govt. High School, Manali, Chennai.     13.043011  15.408602
Highway (Project -I) Building, Madurai  11.947917  24.458333
Kathivakkam, Municipal Kalyana Mandapam, Chennai 12.925532  15.170213
Kilpauk, Chennai                      19.232759  27.172414
Kunnathur Chatram East Avani Mollai Street, Mad... 14.340206  25.577320
Madras Medical College, Chennai         7.418605  27.465116
Main Guard Gate, Tiruchy               17.135135  20.837838
NEERI, CSIR Campus Chennai             5.931034  23.758621
Poniarajapuram, On the top of DEL, Coimbatore    4.126214  23.019417
Raja Agencies, Tuticorin               15.058824  22.441176
Raman Nagar, Mettur                    7.572816  20.407767
SIDCO Industrial Complex, Mettur        9.294118  25.990196
SIDCO Office, Coimbatore                4.969072  27.329897
SIPCOT Industrial Complex, Cuddalore     6.969697  17.666667
Sowdeswari College Building, Salem     8.114504  28.664122
Thiruvottiyur Municipal Office, Chennai  8.360465  28.069767
Thiruvottiyur, Chennai                 13.010417  15.583333
Thiyagaraya Nagar, Chennai             18.849558  28.250000
```

```
              RSPM/PM10
Location of Monitoring Station
AVM Jewellery Building, Tuticorin      70.175258
Adyar, Chennai                        57.068966
Anna Nagar, Chennai                   72.187500
Bishop Heber College, Tiruchy         45.633803
Central Bus Stand, Trichy             120.546667
District Environmental Engineer Office, Imperia...  64.020202
Distt. Collector's Office, Coimbatore   42.972933
Eachangadu Villagae                   75.591837
Fenner (I) Ltd. Employees Association Building ... 40.732673
Fisheries College, Tuticorin           85.255319
Gandhi Market, Trichy                 101.743243
Golden Rock, Trichy                   46.222222
Govt. High School, Manali, Chennai.     44.612903
Highway (Project -I) Building, Madurai  46.427083
Kathivakkam, Municipal Kalyana Mandapam, Chennai 46.851064
Kilpauk, Chennai                      88.103448
Kunnathur Chatram East Avani Mollai Street, Mad... 50.226804
Madras Medical College, Chennai        35.837209
Main Guard Gate, Tiruchy               107.693333
NEERI, CSIR Campus Chennai             43.678161
Poniarajapuram, On the top of DEL, Coimbatore    48.883495
Raja Agencies, Tuticorin               94.230336
Raman Nagar, Mettur                    51.106796
SIDCO Industrial Complex, Mettur        54.352941
SIDCO Office, Coimbatore                55.969072
SIPCOT Industrial Complex, Cuddalore     46.171717
Sowdeswari College Building, Salem     62.954198
Thiruvottiyur Municipal Office, Chennai  34.310345
Thiruvottiyur, Chennai                 42.604167
Thiyagaraya Nagar, Chennai             102.327434
```

Calculate Average SO2, NO2, and RSPM/PM10 Levels

4. Visualization for Stations- average between SO2, NO2, RSPM/PM10:

```
import matplotlib.pyplot as plt

# Data
stations = station_averages.index
so2_station = station_averages['SO2']
no2_station = station_averages['NO2']
rspm_station = station_averages['RSPM/PM10']

# Create subplots
plt.figure(figsize=(15, 6))

# Plot 1 - Box Plot for Average SO2 Levels at Monitoring Stations
plt.subplot(131)
plt.boxplot(so2_station)
plt.title('Average SO2 Levels at Monitoring Stations (Box Plot)')
plt.ylabel('Average Levels')

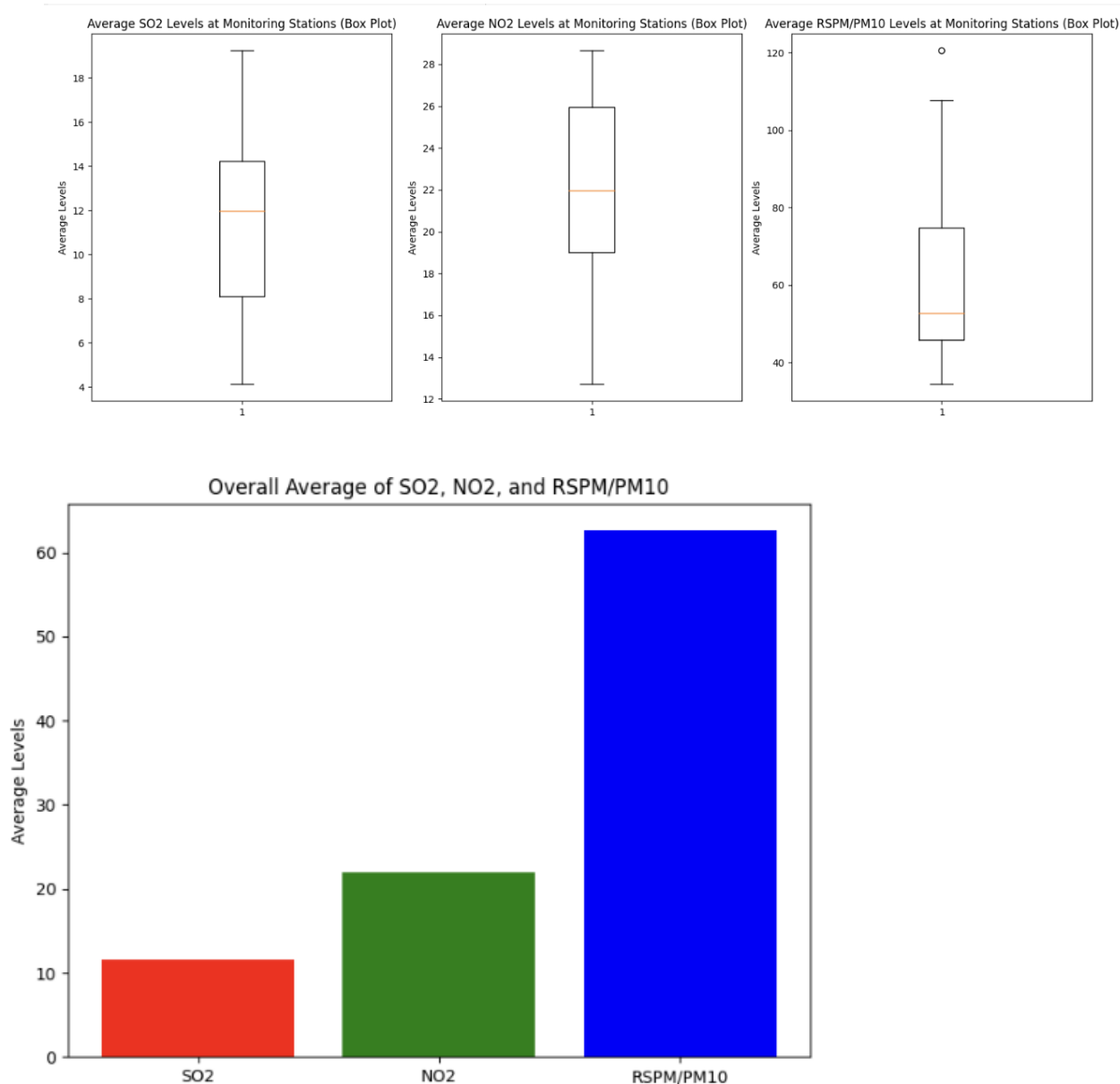
# Plot 2 - Box Plot for Average NO2 Levels at Monitoring Stations
plt.subplot(132)
plt.boxplot(no2_station)
plt.title('Average NO2 Levels at Monitoring Stations (Box Plot)')
plt.ylabel('Average Levels')

# Plot 3 - Box Plot for Average RSPM/PM10 Levels at Monitoring Stations
plt.subplot(133)
plt.boxplot(rspm_station)
plt.title('Average RSPM/PM10 Levels at Monitoring Stations (Box Plot)')
plt.ylabel('Average Levels')

plt.tight_layout()
plt.show()

# Overall Bar Chart - Average of All Three Pollutants
plt.figure(figsize=(8, 6))
overall_average = (so2_station.mean(), no2_station.mean(), rspm_station.mean())
pollutants = ['SO2', 'NO2', 'RSPM/PM10']

plt.bar(pollutants, overall_average, color=['r', 'g', 'b'])
plt.title('Overall Average of SO2, NO2, and RSPM/PM10')
plt.ylabel('Average Levels')
plt.show()
```



5. Visualization for Cities - average between SO2, NO2, RSPM/PM10:

```
import matplotlib.pyplot as plt

# Data
city_areas = city_averages.index
so2_city = city_averages['SO2']
no2_city = city_averages['NO2']
rspm_city = city_averages['RSPM/PM10']

# Create subplots
plt.figure(figsize=(15, 6))

# Plot 1 - Pie Chart for Average SO2 Levels in Cities
plt.subplot(131)
plt.pie(so2_city, labels=city_areas, autopct='%1.1f%%', colors=['r', 'g', 'b', 'y'])
plt.title('Average SO2 Levels in Cities (Pie Chart)')

# Plot 2 - Pie Chart for Average NO2 Levels in Cities
plt.subplot(132)
plt.pie(no2_city, labels=city_areas, autopct='%1.1f%%', colors=['r', 'g', 'b', 'y'])
plt.title('Average NO2 Levels in Cities (Pie Chart)')
```

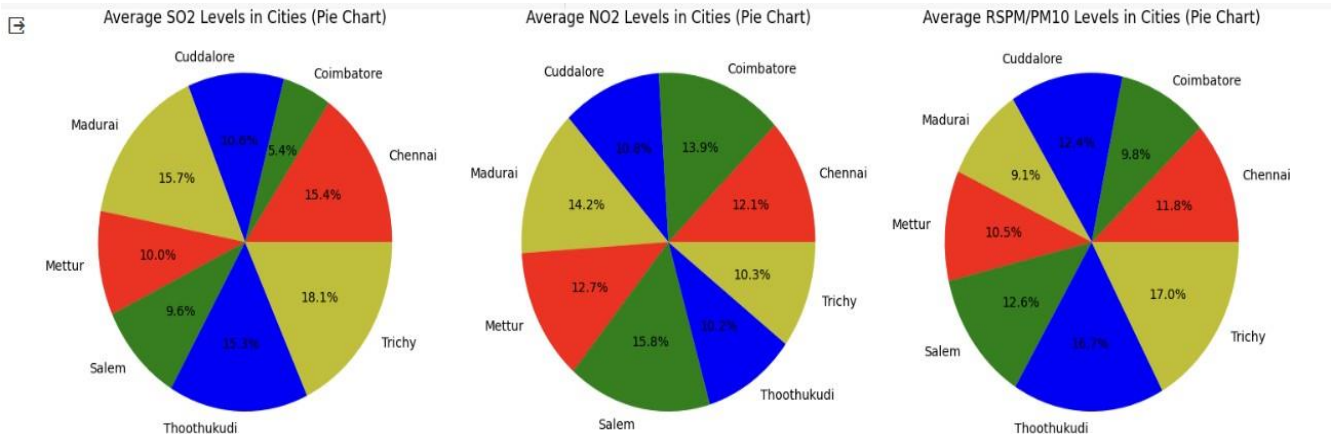
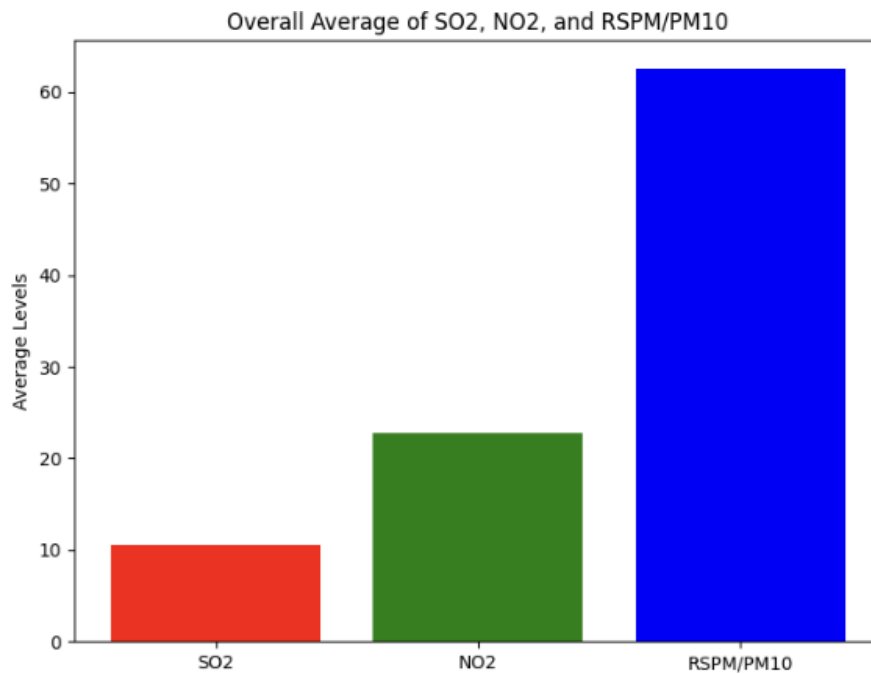


```
# Plot 3 - Pie Chart for Average RSPM/PM10 Levels in Cities
plt.subplot(133)
plt.pie(rspm_city, labels=city_areas, autopct='%1.1f%%', colors=['r', 'g', 'b', 'y'])
plt.title('Average RSPM/PM10 Levels in Cities (Pie Chart)')

plt.tight_layout()
plt.show()

# Overall Bar Chart - Average of All Three Pollutants
plt.figure(figsize=(8, 6))
overall_average = (so2_city.mean(), no2_city.mean(), rspm_city.mean())
pollutants = ['SO2', 'NO2', 'RSPM/PM10']

plt.bar(pollutants, overall_average, color=['r', 'g', 'b'])
plt.title('Overall Average of SO2, NO2, and RSPM/PM10')
plt.ylabel('Average Levels')
plt.show()
```



5. Visualization for Area wise average between SO2, NO2, PSPM/PM10:

```
import matplotlib.pyplot as plt

# Data
areas = area_averages.index
so2_area = area_averages['SO2']
no2_area = area_averages['NO2']
rspm_area = area_averages['RSPM/PM10']

# Create subplots
plt.figure(figsize=(15, 6))

# Plot 1 - Scatter Plot for Average SO2 Levels in Areas
plt.subplot(131)
plt.scatter(areas, so2_area, label='SO2', color='r', marker='o')
plt.title('Average SO2 Levels in Areas')
plt.xlabel('Areas')
plt.ylabel('Average Levels')
plt.xticks(rotation=90)
plt.legend()

# Plot 2 - Scatter Plot for Average NO2 Levels
plt.subplot(132)
plt.scatter(areas, no2_area, label='NO2', color='g', marker='x')
plt.title('Average NO2 Levels in Areas')
plt.xlabel('Areas')
plt.ylabel('Average Levels')
plt.xticks(rotation=90)
plt.legend()

# Plot 3 - Scatter Plot for Average RSPM/PM10 Levels
plt.subplot(133)
plt.scatter(areas, rspm_area, label='RSPM/PM10', color='b', marker='s')
plt.title('Average RSPM/PM10 Levels in Areas')
plt.xlabel('Areas')
plt.ylabel('Average Levels')
plt.xticks(rotation=90)
plt.legend()

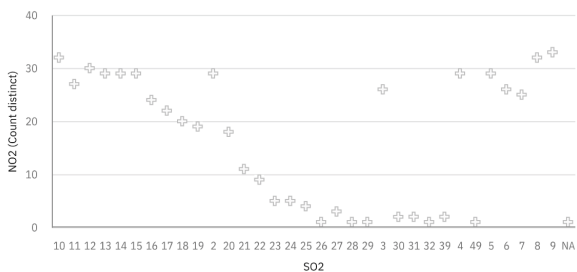
plt.tight_layout()
plt.show()

# Overall Bar Chart - Average of All Three Pollutants
plt.figure(figsize=(8, 6))
overall_average = (so2_area.mean(), no2_area.mean(), rspm_area.mean())
pollutants = ['SO2', 'NO2', 'RSPM/PM10']

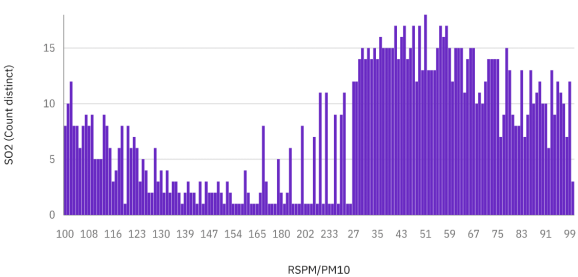
plt.bar(pollutants, overall_average, color=['r', 'g', 'b'])
plt.title('Overall Average of SO2, NO2, and RSPM/PM10')
plt.ylabel('Average Levels')
plt.show()
```


improve the region's environmental well-being and public health.

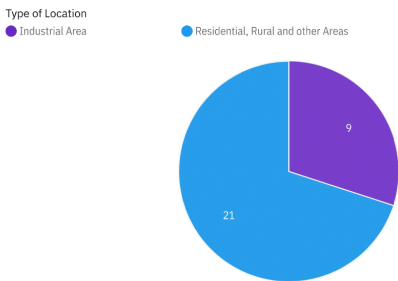
SO2, NO2



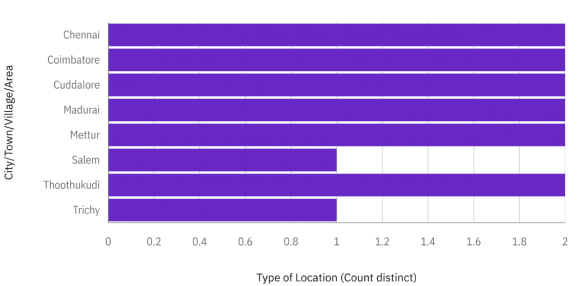
SO2 by RSPM/PM10



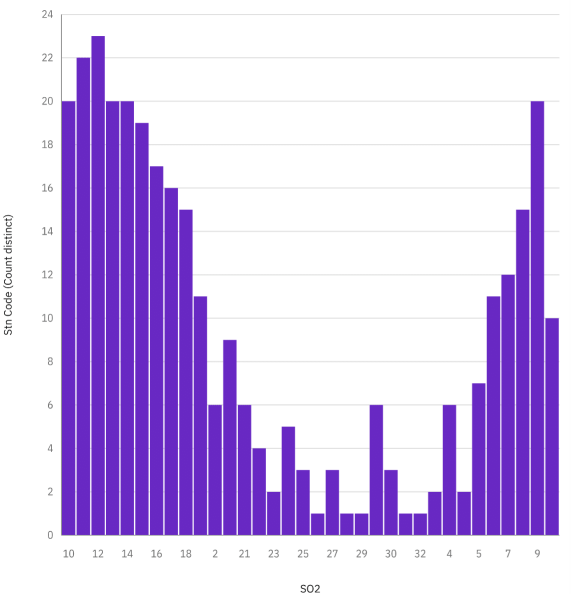
Stn Code by Type of Location



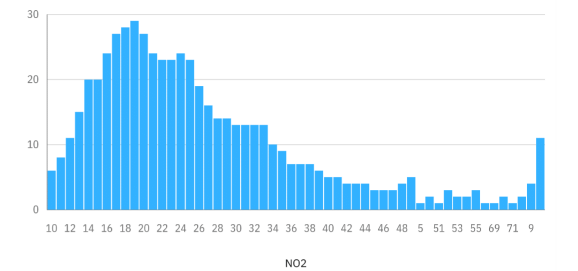
Type of Location by City/Town/Village/Area



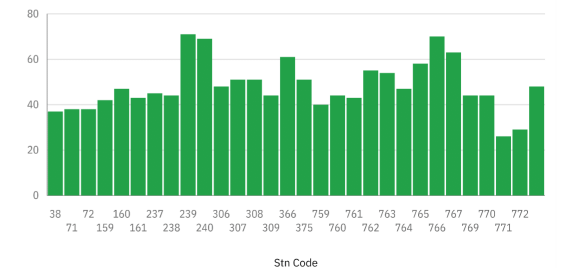
Stn Code by SO2



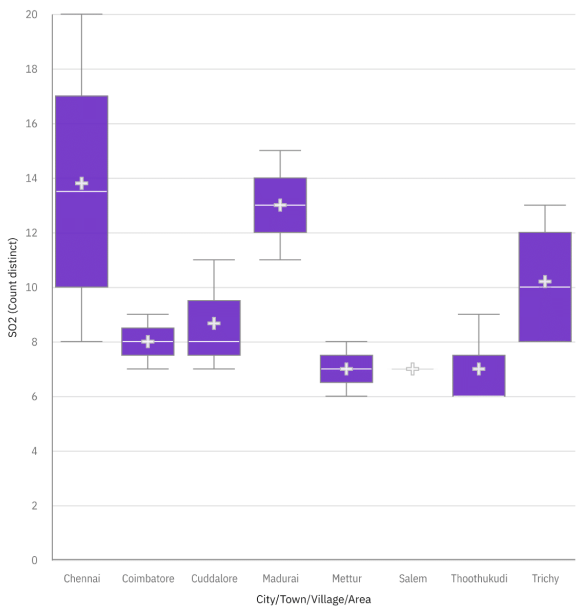
Stn Code by NO2



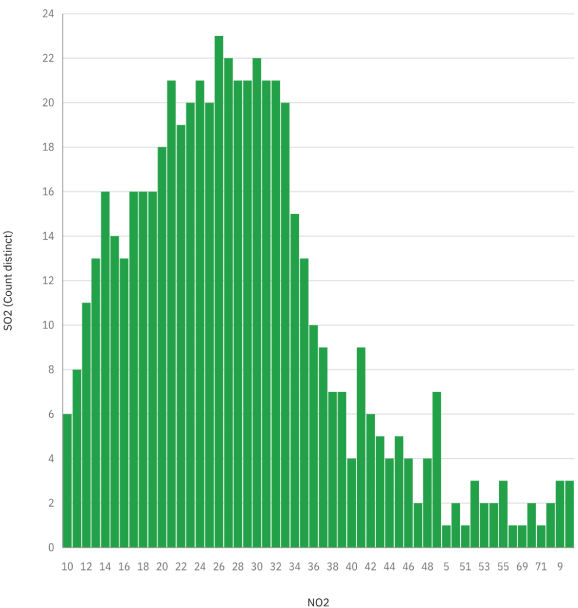
RSPM/PM10 by Stn Code



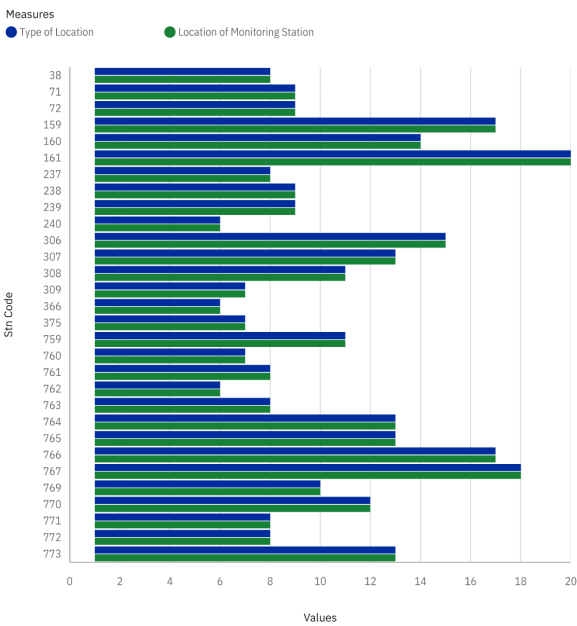
City/Town/Village/Area, SO2, Location of Monitoring Station



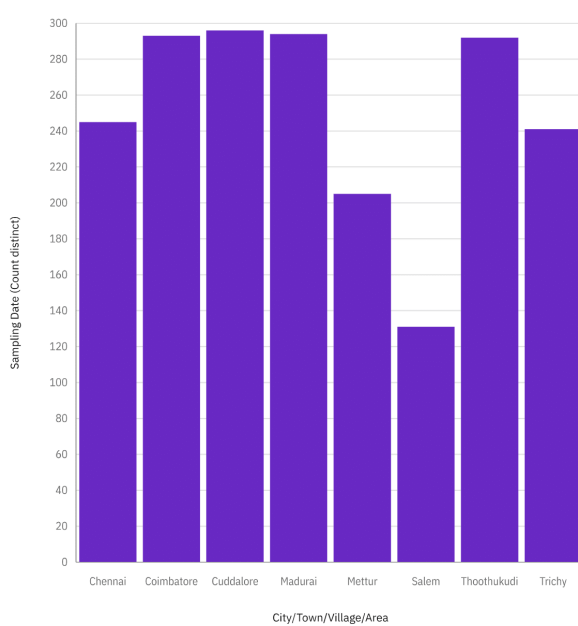
SO2 by NO2



Type of Location and Location of Monitoring Station by Stn Code



Sampling Date by City/Town/Village/Area



Identify Pollution Trends and High Pollution Areas

- Analyse the calculated average pollutant levels to identify trends.
- Use visualizations to identify areas

SO₂, NO₂ & RSPM/PM₁₀ Pollution Levels by Location. (here green shows highest and red shows lowest)

```
import matplotlib.pyplot as plt

# Data
locations = [
    "AVM Jewellery Building, Tuticorin", "Adyar, Chennai", "Anna Nagar, Chennai",
    "Bishop Heber College, Tiruchy", "Central Bus Stand, Trichy",
    "District Environmental Engineer Office, Imperia",
    "Distt. Collector's Office, Coimbatore", "Eachangadu Villagae",
    "Fenner (I) Ltd. Employees Association Building", "Fisheries College, Tuticorin",
    "Gandhi Market, Trichy", "Golden Rock, Trichy",
    "Govt. High School, Manali, Chennai.", "Highway (Project -I) Building, Madurai",
    "Kathivakkam, Municipal Kalyana Mandapam, Chennai", "Kilpauk, Chennai",
    "Kunnathur Chatram East Avani Mollai Street, Mad",
    "Madras Medical College, Chennai", "Main Guard Gate, Tiruchy",
    "NEERI, CSIR Campus Chennai", "Poniarajapuram, On the top of DEL, Coimbatore",
    "Raja Agencies, Tuticorin", "Raman Nagar, Mettur",
    "SIDCO Industrial Complex, Mettur", "SIDCO Office, Coimbatore",
    "SIPCOT Industrial Complex, Cuddalore",
    "Sowdeswari College Building, Salem", "Thiruvottiyur Municipal Office, Chennai",
    "Thiruvottiyur, Chennai", "Thiyagaraya Nagar, Chennai"
```

```
so2_levels = [
    9.302083, 13.252174, 13.873874, 11.800000, 18.013333, 8.101010, 4.554348,
    11.916667, 13.643564, 14.526882, 17.148649, 12.014085, 13.043011, 11.947917,
    12.925532, 19.232759, 14.340206, 7.418605, 17.135135, 5.931034, 4.126214,
    15.058824, 7.572816, 9.294118, 4.969072, 6.969697, 8.114504, 8.360465,
    13.010417, 18.849558
]

# Sort data from lowest to highest SO2 levels
sorted_indices = sorted(range(len(so2_levels)), key=lambda i: so2_levels[i])
locations_sorted = [locations[i] for i in sorted_indices]
so2_levels_sorted = [so2_levels[i] for i in sorted_indices]

# Define colors for highest and lowest values
colors = ['red' if i == min(so2_levels_sorted) else 'green' if i == max(so2_levels_sorted) else 'skyblue' f

# Create a figure
plt.figure(figsize=(8, 10))

# Horizontal bar chart
plt.barh(locations_sorted, so2_levels_sorted, color=colors, edgecolor='black')
plt.xlabel('SO2 Levels')
plt.ylabel('Location of Monitoring Station')
plt.title('SO2 Pollution Levels by Location')

# Display the plot
plt.tight_layout()
plt.show()
```



```

import matplotlib.pyplot as plt

# Data
locations = [
    "AVM Jewellery Building, Tuticorin", "Adyar, Chennai", "Anna Nagar, Chennai",
    "Bishop Heber College, Tiruchy", "Central Bus Stand, Trichy",
    "District Environmental Engineer Office, Imperia",
    "Distt. Collector's Office, Coimbatore", "Eachangadu Villagae",
    "Fenner (I) Ltd. Employees Association Building",
    "Fisheries College, Tuticorin", "Gandhi Market, Trichy",
    "Golden Rock, Trichy", "Govt. High School, Manali, Chennai.",
    "Highway (Project -I) Building, Madurai",
    "Kathivakkam, Municipal Kalyana Mandapam, Chennai",
    "Kilpauk, Chennai",
    "Kunnathur Chatram East Avani Mollai Street, Mad",
    "Madras Medical College, Chennai", "Main Guard Gate, Tiruchy",
    "NEERI, CSIR Campus Chennai",
    "Poniarajapuram, On the top of DEL, Coimbatore",
    "Raja Agencies, Tuticorin", "Raman Nagar, Mettur",
    "SIDCO Industrial Complex, Mettur", "SIDCO Office, Coimbatore",
    "SIPCOT Industrial Complex, Cuddalore",
    "Sowdeswari College Building, Salem",
    "Thiruvottiyur Municipal Office, Chennai",
    "Thiruvottiyur, Chennai", "Thiyagaraya Nagar, Chennai"
]

no2_levels = [
    12.697917, 18.965217, 20.754545, 14.942857, 21.506667,
    19.151515, 25.793478, 22.395833, 27.198020, 20.204301,
    20.797297, 15.000000, 15.408602, 24.458333, 15.170213,
    27.172414, 25.577320, 27.465116, 20.837838, 23.758621,
    23.019417, 22.441176, 20.407767, 25.990196, 27.329897,
    17.666667, 28.664122, 28.069767, 15.583333, 28.250000
]

```

```

# Sort the data from lowest to highest NO2 levels
sorted_indices = sorted(range(len(no2_levels)), key=lambda i: no2_levels[i])
locations_sorted = [locations[i] for i in sorted_indices]
no2_levels_sorted = [no2_levels[i] for i in sorted_indices]

# Highlight the highest and lowest values
colors = ['green' if x == max(no2_levels_sorted) else 'red' if x == min(no2_levels_sorted) else 'skyblue'
          for x in no2_levels_sorted]

# Create a figure
plt.figure(figsize=(8, 12))

# Horizontal bar chart
plt.barh(locations_sorted, no2_levels_sorted, color=colors)
plt.xlabel('NO2 Levels')
plt.ylabel('Location of Monitoring Station')
plt.title('NO2 Pollution Levels by Location')
plt.show()

```

```

import matplotlib.pyplot as plt

# Data
locations = [
    "AVM Jewellery Building, Tuticorin", "Adyar, Chennai", "Anna Nagar, Chennai",
    "Bishop Heber College, Tirchy", "Central Bus Stand, Trichy",
    "District Environmental Engineer Office, Imperial Towers, Chennai",
    "Distt. Collector's Office, Coimbatore", "Eachangadu Village",
    "Fenner (I) Ltd. Employees Association Building (at Entrance), Madurai",
    "Fisheries College, Tuticorin", "Gandhi Market, Trichy",
    "Golden Rock, Trichy", "Govt. High School, Manali, Chennai.",
    "Highway (Project -I) Building, Madurai",
    "Kathivakkam, Municipal Kalyana Mandapam, Chennai", "Kilpauk, Chennai",
    "Kunnathur Chatram East Avani Mollai Street, Madurai",
    "Madras Medical College, Chennai", "Main Guard Gate, Tirchy",
    "NEERI, CSIR Campus Chennai", "Poniarajapuram, On the top of DEL, Coimbatore",
    "Raja Agencies, Tuticorin", "Raman Nagar, Mettur",
    "SIDCO Industrial Complex, Mettur", "SIDCO Office, Coimbatore",
    "SIPCOT Industrial Complex, Cuddalore",
    "Sowdeswari College Building, Salem",
    "Thiruvottiyur Municipal Office, Chennai", "Thiruvottiyur, Chennai",
    "Thiyagaraya Nagar, Chennai"
]

rspm_pm10_levels = [
    70.175258, 57.068966, 72.187500, 45.633803, 120.546667,
    64.020202, 42.972933, 75.591837, 40.732673, 85.255319,
    101.743243, 46.222222, 44.612903, 46.427083, 46.851064,
    88.103448, 50.226804, 35.837209, 107.693333, 43.678161,
    48.883495, 94.230336, 51.106796, 54.352941, 55.969072,
    46.171717, 62.954198, 34.310345, 42.604167, 102.327434
]

# Combine the data into a list of tuples for sorting
data = list(zip(locations, rspm_pm10_levels))

# Sort by RSPM/PM10 levels
sorted_data = sorted(data, key=lambda x: x[1])

# Unzip the sorted data
sorted_locations, sorted_rspm_pm10_levels = zip(*sorted_data)

```

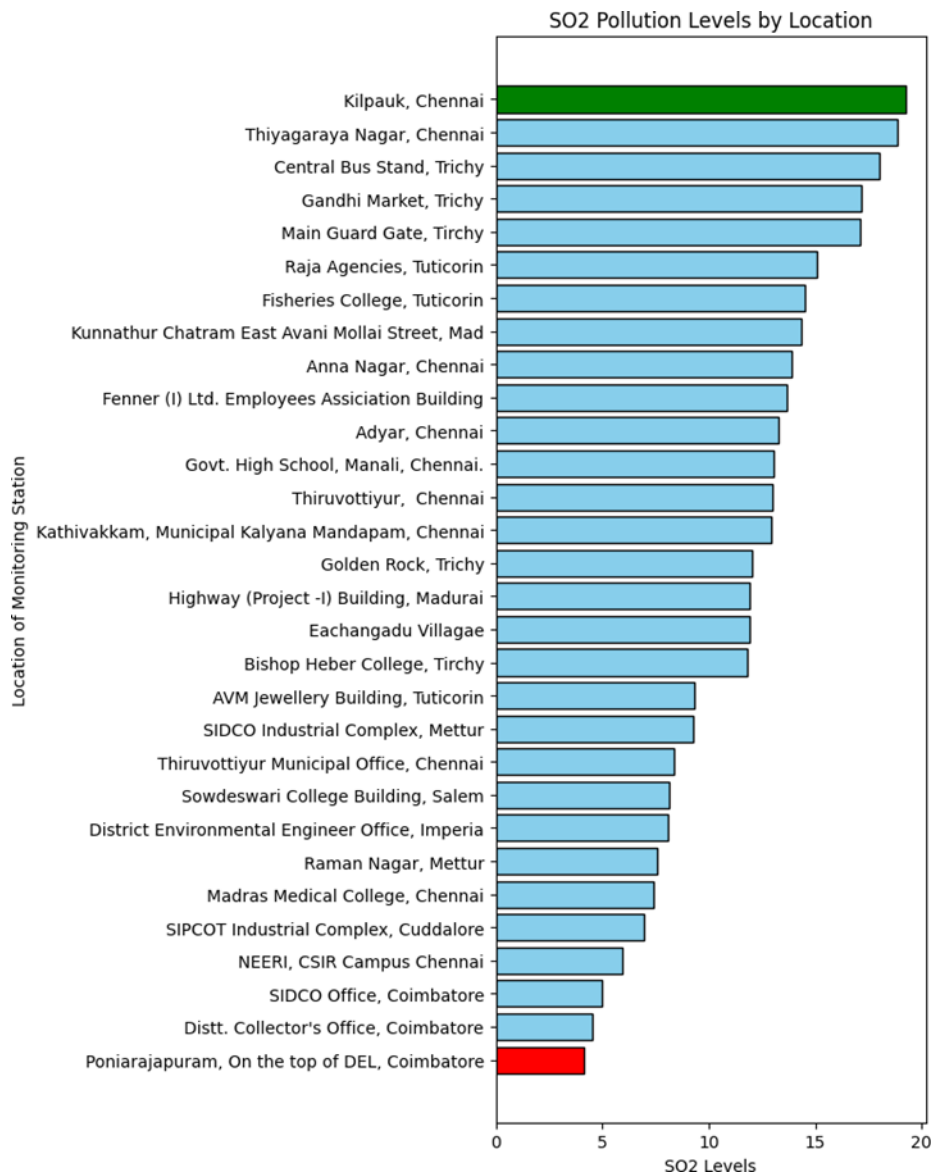


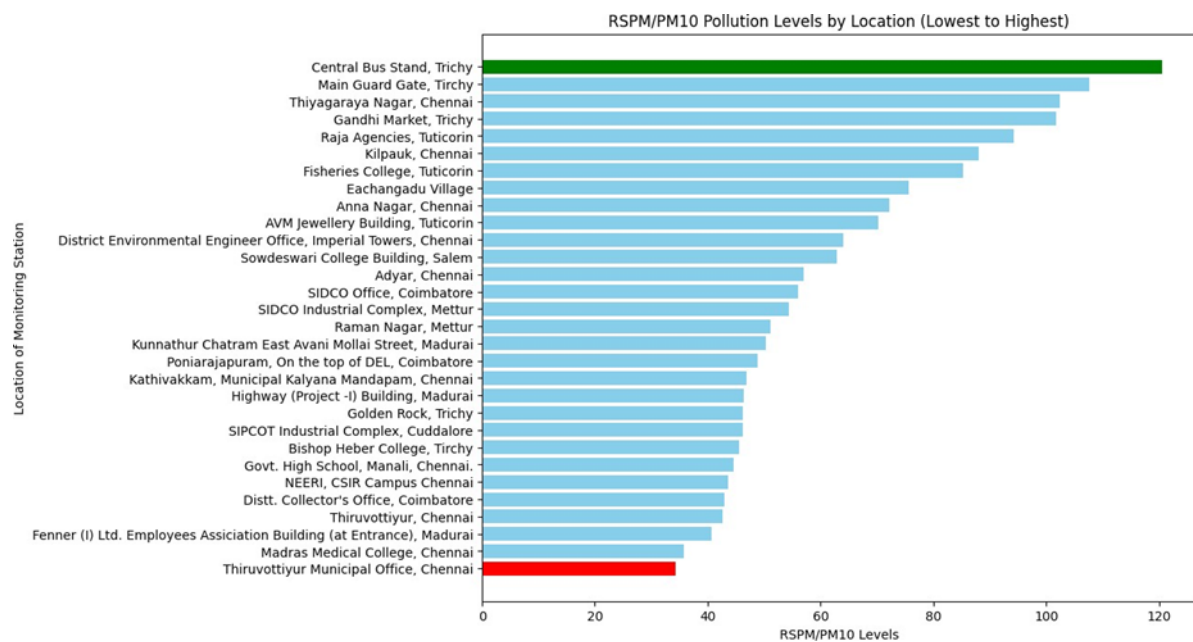
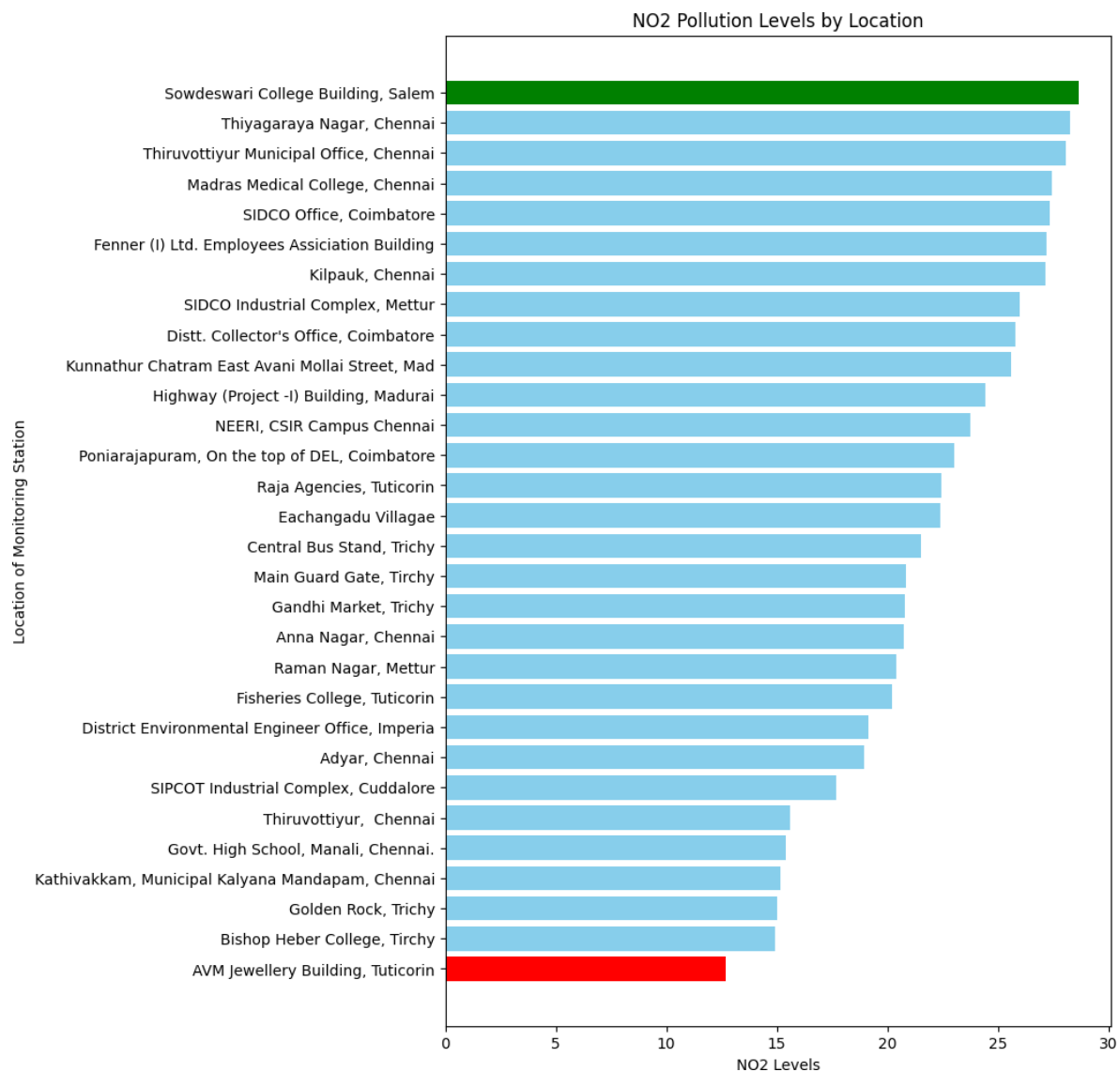
```
# Find the indices of the specified locations
red_location = "Thiruvottiyur Municipal Office, Chennai"
green_location = "Central Bus Stand, Trichy"

red_index = sorted_locations.index(red_location)
green_index = sorted_locations.index(green_location)

# Define colors for highest and lowest
colors = ['skyblue' if i != green_index and i != red_index else 'green' if i == green_index else 'red'
          for i in range(len(sorted_locations))]

# Create a horizontal bar chart
plt.figure(figsize=(10, 8))
plt.barh(sorted_locations, sorted_rspm_pm10_levels, color=colors)
plt.xlabel('RSPM/PM10 Levels')
plt.ylabel('Location of Monitoring Station')
plt.title('RSPM/PM10 Pollution Levels by Location (Lowest to Highest)')
plt.show()
```





Insights:

The Air Quality Analysis project for Tamil Nadu has yielded critical insights into air pollution trends and pollution levels in the region, providing a data-driven basis for informed decision-making and policy formulation:

1. **Temporal Trends:** Through rigorous time series analysis, our project reveals significant temporal trends in air quality over the years. This analysis enables us to identify periods characterized by both high and low pollution levels. By pinpointing these trends, we contribute to a better understanding of how air quality has evolved over time, facilitating the recognition of long-term patterns and variations.
2. **Seasonal Variations:** Seasonal analysis offers a deeper insight into the variations in pollutant levels that occur throughout the year. By considering the influence of weather conditions, local activities, and other seasonal factors, we gain a comprehensive understanding of how different seasons impact air quality. This information is invaluable for designing season-specific environmental strategies and interventions.
3. **Geographical Disparities:** One of the most significant findings from our analysis is the revelation of geographical disparities in air quality. Aggregating and visualizing data by location, such as monitoring stations, cities, or areas, exposes areas that consistently exhibit higher pollution levels. These disparities serve as a crucial indicator of potential pollution hotspots, directing attention to regions in need of targeted interventions and rigorous monitoring.
4. **Early Warning System:** The project's identification of high pollution trends functions as an early warning system for emerging pollution issues. By recognizing patterns of increasing pollution, we enable swift and proactive responses from environmental protection agencies. This early warning system is instrumental in mitigating air quality deterioration and protecting public health.
5. **Data-Driven Decision-Making:** The visualization techniques and insights provided by our project empower policymakers, environmental agencies, and the general public to make data-driven decisions. These decisions can range from implementing targeted pollution control measures to raising public awareness about air quality challenges. The actionable intelligence derived from our analysis supports a collective effort to improve

air quality and safeguard public health in Tamil Nadu.

Conclusion :

The Air Quality Analysis project for Tamil Nadu is a comprehensive effort aimed at gaining a deep understanding of air pollution trends and pollution levels in the region. By harnessing the power of data science, machine learning, and data visualization, this project is designed to have a substantial impact on public health and environmental well-being.

The primary objectives of the project encompass the entire process of data collection, integration, analysis, and visualization. It seeks to provide accurate and reliable insights into air quality, using a data-driven approach to ensure that all findings are grounded in empirical evidence. This commitment to transparency, fairness, and ethical considerations ensures that the project's results are not only informative but also trustworthy and inclusive.

The project's extensive use of data visualization techniques simplifies complex data, making it accessible to a broad audience, including policymakers, environmental agencies, and the general public. By leveraging these techniques, the project empowers stakeholders to make informed decisions that can improve air quality, safeguard public health, and create a more sustainable and livable environment in Tamil Nadu.

The impact of the project extends to several critical areas. First and foremost, it aims to improve public health by identifying pollution trends and high pollution levels. By offering information about areas that require immediate attention, it enables timely interventions to protect residents' health.

Additionally, the project plays a pivotal role in environmental protection. Through the identification of pollution hotspots, seasonal variations, and geographical disparities, resources and measures can be directed more effectively, thereby enhancing pollution control and mitigation. Ultimately, the project contributes to the sustainability of Tamil Nadu by addressing air quality challenges. It promotes a higher quality of life, healthier living conditions, and a greener future for the region.