

Air Quality Assessment TN

Problem Definition:

The project aims to analyze and visualize air quality data from monitoring stations in Tamil Nadu. The objective is to gain insights into air pollution trends, identify areas with high pollution levels, and develop a predictive model to estimate RSPM/PM10 levels based on SO2 and NO2 levels. This project involves defining objectives, designing the analysis approach, selecting visualization techniques, and creating a predictive model using Python and relevant libraries.

Overview:

This documentation provides a comprehensive overview of all the tasks completed during the session, beginning with Design Thinking to determine how the assigned project should be completed and what tools and techniques are needed. It also includes the incorporation of a machine learning model into the application, loading and pre-processing the provided dataset, and analysing it using a variety of visualisation techniques, such as plots, charts, and graphs, to extract insights and forecast pollution levels and trends in Tamil Nadu's air pollution.

Project Objectives:

Objective 1: Analyze Air Quality Trends

- Approach: We will analyze the historical air quality data to identify trends, seasonal variations, and any significant changes in pollution levels. This will involve statistical analysis and timeseries visualization.

Objective 2: Identify Pollution Hotspots

- Approach: We will use geospatial analysis techniques to pinpoint regions with consistently high pollution levels. Heatmaps and spatial visualizations will help in identifying pollution hotspots.

Objective 3: Build a Predictive Model

- Approach: To create a predictive model for RSPM and PM10 levels, we will perform regression analysis. The model will be trained using historical data, and features such as SO2 and NO2 levels will be used as predictors.

Design Thinking:

Analysis Approach:

1. Data Collection:

- We will gather the air quality dataset from the provided link:
<https://tn.data.gov.in/resource/location-wise-daily-ambient-air-quality-tamil-nadu-year-2014>
- Data will be downloaded in a suitable format (e.g., CSV) for further analysis.

2. Data Pre-processing:

- Clean and pre-process the dataset to handle missing values, outliers, and inconsistencies.
- Convert relevant columns to appropriate data types.
- Create derived features, such as pollutant concentration indices.

3. Exploratory Data Analysis (EDA):

- Perform EDA to understand the basic statistics and distributions of air quality parameters.
- Identify trends, seasonality, and correlations between variables.

4. Data Visualization:

- Utilize various visualization techniques, including line charts, bar plots, heatmaps, and geospatial plots, to effectively communicate findings.

5. Pollution Hotspot Analysis:

- Employ geospatial analysis tools (e.g., GIS) to create heatmaps and identify areas with high pollution levels.

6. Predictive Modeling:

- Split the dataset into training and testing sets.
- Build regression models (e.g., linear regression, random forest) to predict RSPM and PM10 levels based on SO2 and NO2 levels.
- Evaluate model performance using appropriate metrics (e.g., RMSE, R-squared).

7. Documentation and Reporting:

- Document all steps, code, and findings in a clear and organized manner.
- Create a comprehensive report summarizing the analysis, visualizations, and model results.

Visualization Selection:

1. Air Quality Trends:

- Line charts to show the variation in air quality parameters (RSPM, PM10, SO2, NO2) over time (monthly or seasonal trends).
- Time-series decomposition plots to visualize seasonality and trend components.

2. Pollution Hotspot Identification:

- Heatmaps to represent spatial variations in pollution levels.
- Geospatial plots with color-coded markers to pinpoint pollution hotspots on a map of Tamil Nadu.

3. Predictive Model Evaluation:

- Scatter plots to visualize the predicted vs. actual values of RSPM and PM10.
- Residual plots to assess model performance and identify any patterns in prediction errors.

By following this structured approach and design thinking principles, we aim to gain valuable insights into air quality in Tamil Nadu, identify areas in need of environmental intervention, and provide a predictive model for estimating pollutant levels based on key parameters.

Innovation:

Consider incorporating machine learning algorithms to improve the accuracy of the predictive model.

1. Data Preparation:
 - Download and clean the air quality data from the provided dataset link.
 - Handle missing values, if any, and perform necessary data preprocessing.
 - Create a dataset that includes features (SO₂ and NO₂ levels) and the target variable (RSPM/PM₁₀ levels).
2. Feature Engineering:
 - Consider adding additional features that may impact air quality, such as weather conditions, geographical coordinates, or time of day.
 - Transform or engineer features to make them more informative for the predictive model.
3. Data Splitting:
 - Split your dataset into training, validation, and test sets. This ensures that you have a separate dataset for model training, validation, and final evaluation.
4. Select Machine Learning Algorithm:
 - Choose a machine learning algorithm suitable for your regression task. Linear regression, decision trees, random forests, or gradient boosting are good options to start with.
5. Model Training and Tuning:
 - Train the selected machine learning model on the training dataset.
 - Tune the model's hyperparameters to optimize its performance. You can use techniques like grid search or random search for hyperparameter tuning.
6. Model Evaluation:
 - Evaluate the model's performance using appropriate regression metrics such as Mean Absolute Error (MAE), Mean Squared Error (MSE), or R-squared (R²).
 - Use the validation dataset to assess the model's accuracy and generalization capability.
7. Visualization:
 - Create visualizations to present your findings and model performance.
 - Use libraries like Matplotlib or Seaborn to create informative plots and graphs.
8. Model Interpretation:
 - Depending on the chosen algorithm, consider methods for model interpretation and feature importance analysis. This helps understand which features have the most significant impact on air quality.
9. Predictive Model Deployment (Optional):
 - If applicable, deploy the predictive model in a real-time or batch processing environment, allowing for real-time predictions or batch predictions for future data.

10. Documentation and Reporting:

- Document the entire process, including data preprocessing, model selection, hyperparameter tuning, and evaluation.
- Create a detailed report or presentation summarizing your findings, insights, and the performance of the predictive model.

11. Phase 2 Innovation:

- Explore advanced machine learning techniques or alternative models to further enhance predictive accuracy.
- Consider incorporating time series forecasting methods if the dataset contains temporal information.
- Remember to continuously monitor and update your predictive model as new data becomes available to ensure it remains accurate and relevant for air quality assessment in Tamil Nadu.

DEVELOPMENT PART 1:

Step 1: Data Loading

Data loading is the process of bringing external data into a format suitable for analysis. In this case, we've imported data in CSV format by utilizing the Pandas library and subsequently printed it to confirm the successful loading of the data.

Step 2: Explore the data

Exploring the data using the `head()` and `info()` function is a process of initially examining a dataset to understand its structure, content, and quality.

head() - This function displays the first few rows of the dataset.

info() - It displays information about the data types of each column, the number of non-null entries, and the memory usage.

Step 3: Data cleaning

To address the issue of missing values in the provided dataset, we can resolve it by,

- Filling those missing values with zeros.
- Check whether the data set contain any missing values
- Replace the missing values with zeros
- Save the preprocessed data to a new file
- Check missing values again to verify they are handled

Code and Output:

```
import pandas as pd
data = pd.read_csv('/content/Air_quality_TN_Dataset.csv')
df = pd.DataFrame(data)
print(df)
```

Stn	Code	Sampling Date	State	City/Town/Village/Area \
0	38	01-02-14	Tamil Nadu	Chennai
1	38	01-07-14	Tamil Nadu	Chennai
2	38	21-01-14	Tamil Nadu	Chennai
3	38	23-01-14	Tamil Nadu	Chennai
4	38	28-01-14	Tamil Nadu	Chennai
...
2874	773	12-03-14	Tamil Nadu	Trichy
2875	773	12-10-14	Tamil Nadu	Trichy
2876	773	17-12-14	Tamil Nadu	Trichy
2877	773	24-12-14	Tamil Nadu	Trichy
2878	773	31-12-14	Tamil Nadu	Trichy

	Location of Monitoring Station \
0	Kathivakkam, Municipal Kalyana Mandapam, Chennai
1	Kathivakkam, Municipal Kalyana Mandapam, Chennai
2	Kathivakkam, Municipal Kalyana Mandapam, Chennai
3	Kathivakkam, Municipal Kalyana Mandapam, Chennai
4	Kathivakkam, Municipal Kalyana Mandapam, Chennai
...	...
2874	Central Bus Stand, Trichy
2875	Central Bus Stand, Trichy
2876	Central Bus Stand, Trichy
2877	Central Bus Stand, Trichy
2878	Central Bus Stand, Trichy

	Agency \
0	Tamilnadu State Pollution Control Board
1	Tamilnadu State Pollution Control Board
2	Tamilnadu State Pollution Control Board
3	Tamilnadu State Pollution Control Board
4	Tamilnadu State Pollution Control Board
...	...
2874	Tamilnadu State Pollution Control Board
2875	Tamilnadu State Pollution Control Board
2876	Tamilnadu State Pollution Control Board
2877	Tamilnadu State Pollution Control Board
2878	Tamilnadu State Pollution Control Board

	Type of Location	SO2	NO2	RSPM/PM10	PM 2.5
0	Industrial Area	11.0	17.0	55.0	NaN

1	Industrial Area	13.0	17.0	45.0	NaN
2	Industrial Area	12.0	18.0	50.0	NaN
3	Industrial Area	15.0	16.0	46.0	NaN
4	Industrial Area	13.0	14.0	42.0	NaN
...
2874	Residential, Rural and other Areas	15.0	18.0	102.0	NaN
2875	Residential, Rural and other Areas	12.0	14.0	91.0	NaN
2876	Residential, Rural and other Areas	19.0	22.0	100.0	NaN
2877	Residential, Rural and other Areas	15.0	17.0	95.0	NaN
2878	Residential, Rural and other Areas	14.0	16.0	94.0	NaN

[2879 rows x 11 columns]

print(data.head())

output

	Stn Code	Sampling Date	State	City/Town/Village/Area \
0	38	01-02-14	Tamil Nadu	Chennai
1	38	01-07-14	Tamil Nadu	Chennai
2	38	21-01-14	Tamil Nadu	Chennai
3	38	23-01-14	Tamil Nadu	Chennai
4	38	28-01-14	Tamil Nadu	Chennai

Location of Monitoring Station \

0	Kathivakkam, Municipal Kalyana Mandapam, Chennai
1	Kathivakkam, Municipal Kalyana Mandapam, Chennai
2	Kathivakkam, Municipal Kalyana Mandapam, Chennai
3	Kathivakkam, Municipal Kalyana Mandapam, Chennai
4	Kathivakkam, Municipal Kalyana Mandapam, Chennai

Agency Type of Location SO2 NO2 \

0	Tamilnadu State Pollution Control Board	Industrial Area	11.0	17.0
1	Tamilnadu State Pollution Control Board	Industrial Area	13.0	17.0
2	Tamilnadu State Pollution Control Board	Industrial Area	12.0	18.0
3	Tamilnadu State Pollution Control Board	Industrial Area	15.0	16.0
4	Tamilnadu State Pollution Control Board	Industrial Area	13.0	14.0

RSPM/PM10 PM 2.5

0	55.0	NaN
1	45.0	NaN
2	50.0	NaN
3	46.0	NaN
4	42.0	NaN

print(data.info())

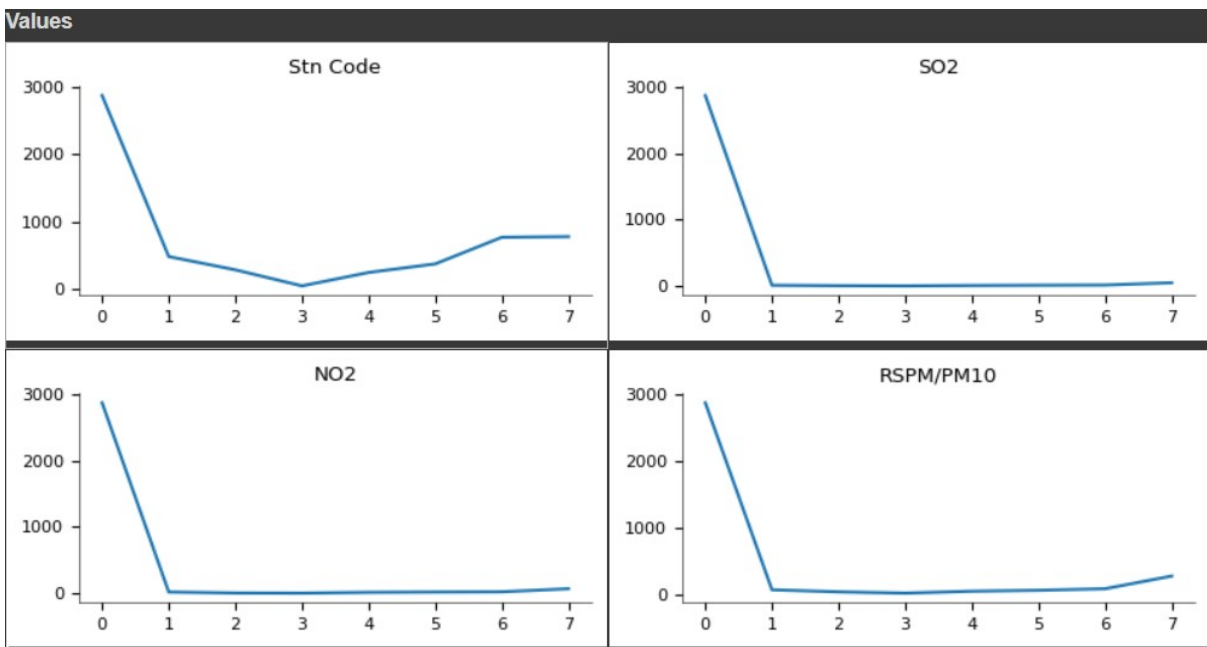
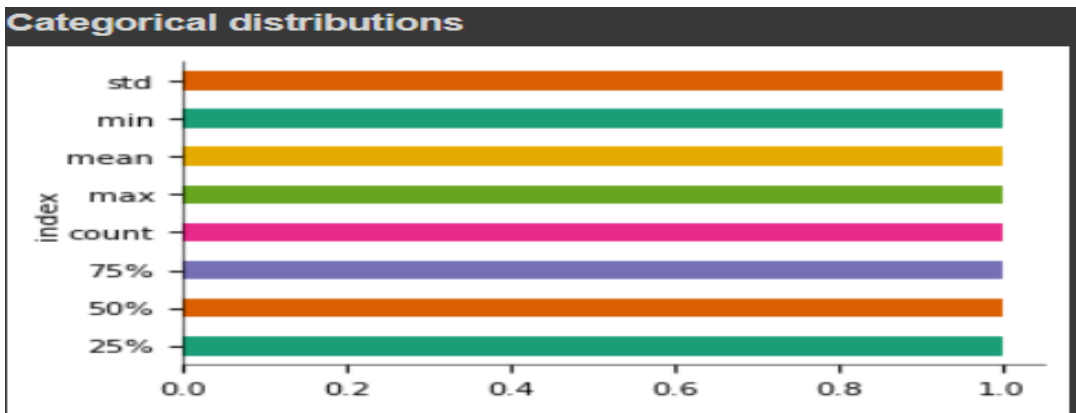
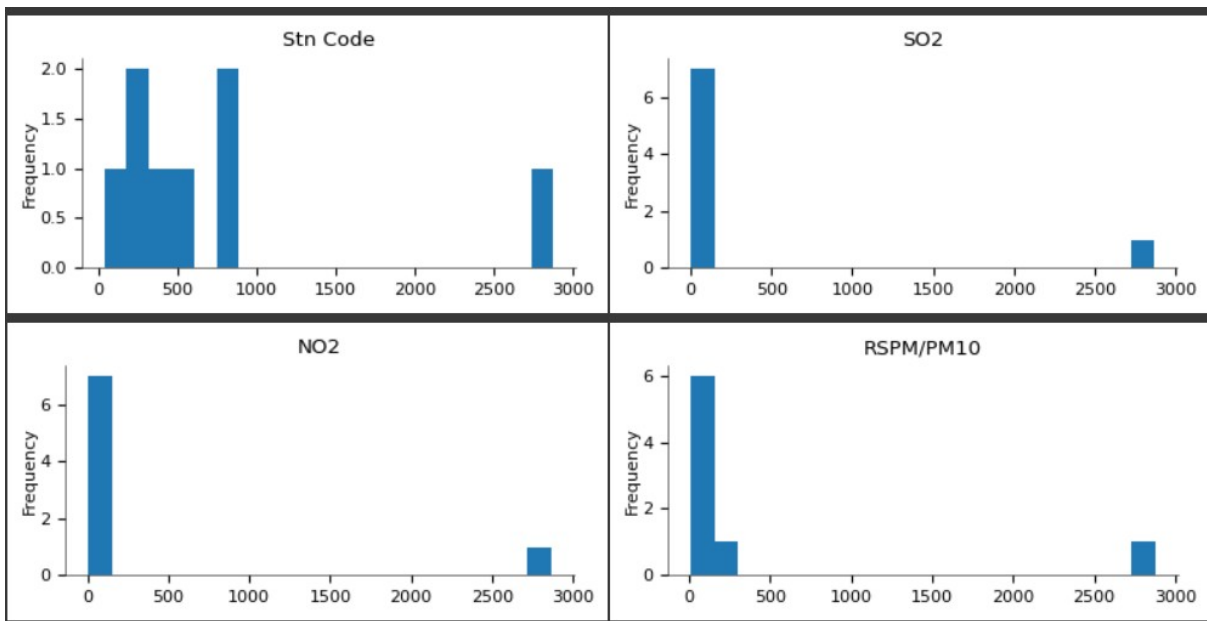
```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 2879 entries, 0 to 2878
Data columns (total 11 columns):
#   Column                                Non-Null Count  Dtype
---  -
0   Stn Code                             2879 non-null   int64
1   Sampling Date                        2879 non-null   object
2   State                               2879 non-null   object
3   City/Town/Village/Area              2879 non-null   object
4   Location of Monitoring Station       2879 non-null   object
5   Agency                             2879 non-null   object
6   Type of Location                    2879 non-null   object
7   SO2                                 2868 non-null   float64
8   NO2                                 2866 non-null   float64
9   RSPM/PM10                          2875 non-null   float64
10  PM 2.5                             0 non-null      float64
dtypes: float64(4), int64(1), object(6)
memory usage: 247.5+ KB
None
```

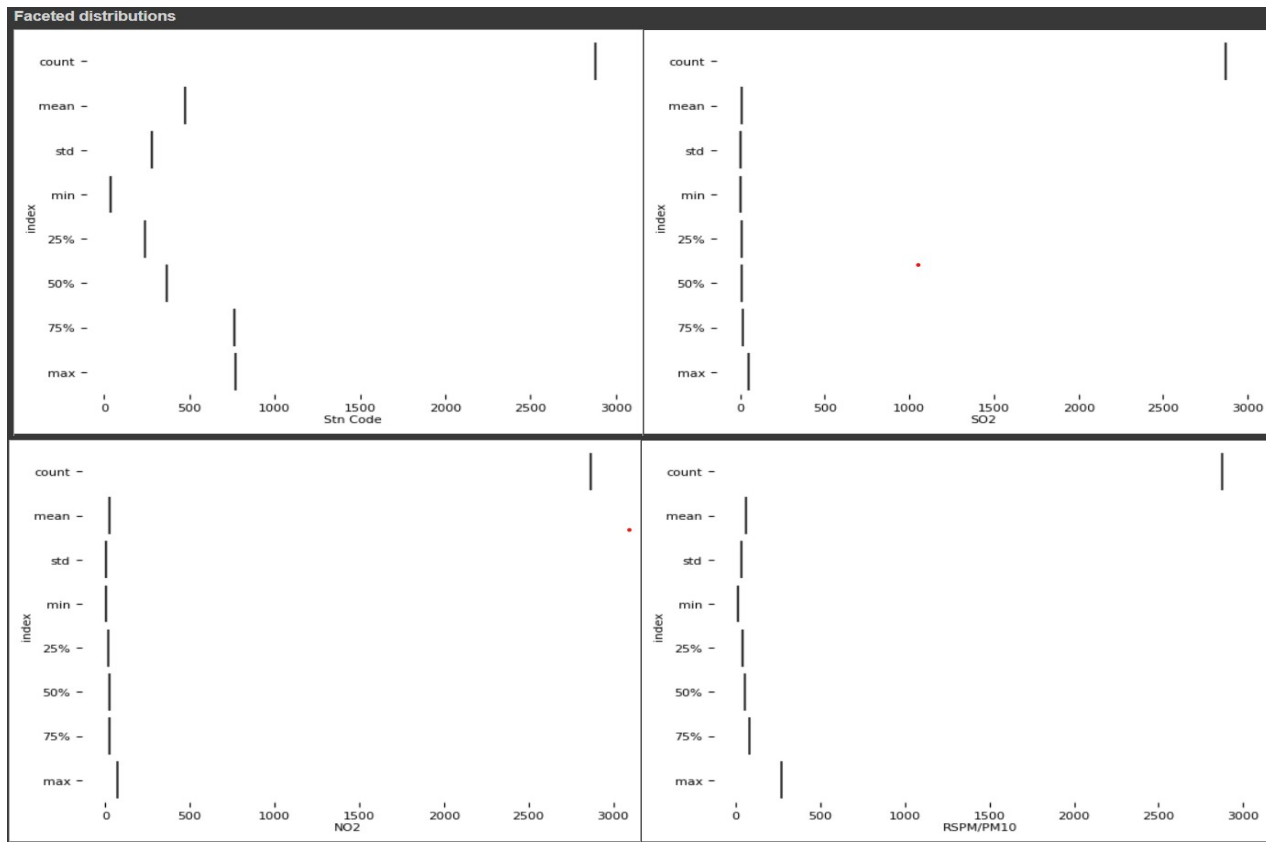
df.dropna()

Stn Code	Sampling Date	State	City/Town/Village/Area	Location of Monitoring Station	Agency	Type of Location	SO2	NO2	RSPM/PM10	PM 2.5
----------	---------------	-------	------------------------	--------------------------------	--------	------------------	-----	-----	-----------	--------

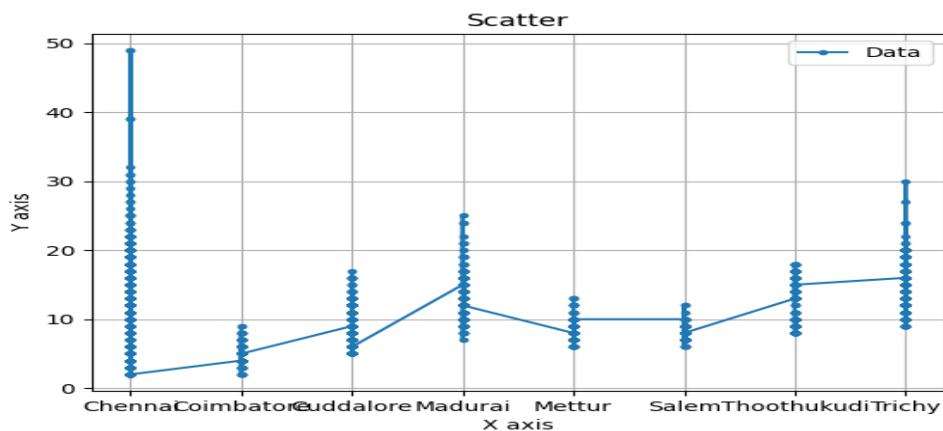
df.describe()

	Stn Code	SO2	NO2	RSPM/PM10	PM 2.5
count	2879.000000	2868.000000	2866.000000	2875.000000	0.0
mean	475.750261	11.503138	22.136776	62.494261	NaN
std	277.675577	5.051702	7.128694	31.368745	NaN
min	38.000000	2.000000	5.000000	12.000000	NaN
25%	238.000000	8.000000	17.000000	41.000000	NaN
50%	366.000000	12.000000	22.000000	55.000000	NaN
75%	764.000000	15.000000	25.000000	78.000000	NaN
max	773.000000	49.000000	71.000000	269.000000	NaN





```
import matplotlib.pyplot as plt
x=data['City/Town/Village/Area']
y=data['SO2']
plt.plot(x,y,marker='.',linestyle='-',label='Data')
plt.xlabel("X axis")
plt.ylabel("Y axis")
plt.title("Scatter")
plt.legend()
plt.grid(True)
plt.show()
```



DEVELOPMENT PART II:

Step 1: Data Loading

Data loading is the process of bringing external data into a format suitable for analysis. In this case, we've imported pre-processed data in CSV format by utilizing the Pandas library and subsequently printed it to confirm the successful loading of the data.

Step 2: Data Analysis

This analysis aims to visually assess patterns and variations in SO2 levels across different locations (City/Town/Village/Area). It helps identify areas with notably high or low SO2 pollution levels, providing insights into air quality variations across different areas.

Step 3: Scatter Plot

It creates the scatter plot with the specified data, axis labels, color, size, and title. The plot visually represents the relationship between SO2, NO2, and RSPM/PM10 levels, with color and marker size indicating RSPM/PM10 levels, making it easy to observe patterns and associations between these variables.

Code:

```
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
%matplotlib inline
plt.rcParams['figure.figsize'] = (10,7)
import warnings
warnings.filterwarnings('ignore')
import os

data = pd.read_csv("cpcb_dly_aq_tamil_nadu-2014 (1).csv")
data.fillna(0, inplace = True)
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
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	0.0
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	0.0
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	0.0
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	0.0
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	0.0

```

def calculate_si(SO2):
    si=0
    if (SO2<=40):
        si = SO2*(50/40)
    if (SO2>40 and SO2<=80):
        si = 50+(SO2-40)*(50/40)
    if (SO2>80 and SO2<=380):
        si = 100+(SO2-80)*(100/300)
    if (SO2>380 and SO2<=800):
        si = 200+(SO2-380)*(100/800)
    if (SO2>800 and SO2<=1600):
        si = 300+(SO2-800)*(100/800)
    if (SO2>1600):
        si = 400+(SO2-1600)*(100/800)
    return si
data['si']=data['SO2'].apply(calculate_si)
df=data[['SO2','si']]

```

```

df.head()
data['si']=data['SO2'].apply(calculate_si)
df=data[['SO2','si']]
df.head()

```

	SO2	si
0	11.0	13.75
1	13.0	16.25
2	12.0	15.00
3	15.0	18.75
4	13.0	16.25

```

def calculate_ni(NO2):
    ni=0
    if (NO2<=40):
        ni = NO2*50/40
    elif (NO2>40 and NO2<=80):
        ni = 50+(NO2-14)*(50/40)
    elif (NO2>80 and NO2<=180):

```

```

        ni = 100+(NO2-80)*(100/100)
    elif (NO2>180 and NO2<=280):
        ni = 200+(NO2-180)*(100/100)
    elif (NO2>280 and NO2<=400):
        ni = 300+(NO2-800)*(100/120)
    else:
        ni = 400+(NO2-400)*(100/120)
    return ni

```

```

data['ni']=data['NO2'].apply(calculate_ni)
df=data[['NO2','ni']]
df.head()

```

	NO2	ni
0	17.0	21.25
1	17.0	21.25
2	18.0	22.50
3	16.0	20.00
4	14.0	17.50

```

def calculate_aqi(si,ni):
    aqi=0
    if(si>ni):
        aqi=si
    if(ni>si):
        aqi=ni
    return aqi

```

```

data['AQI']=data.apply(lambda x:calculate_aqi(x['si'],x['ni']),axis=1)
df=data[['Sampling Date', 'City/Town/Village/Area', 'si', 'ni', 'AQI']]
df.head()

```

	Sampling Date	City/Town/Village/Area	si	ni	AQI
0	01-02-14	Chennai	13.75	21.25	21.25
1	01-07-14	Chennai	16.25	21.25	21.25
2	21-01-14	Chennai	15.00	22.50	22.50
3	23-01-14	Chennai	18.75	20.00	20.00
4	28-01-14	Chennai	16.25	17.50	17.50

```
def calculate_rspmi(rspm):
```

```
    rspmi=0
    if (rspm<=40):
        rspmi = rspm*50/40
    elif (rspm>40 and rspm<=80):
        rspmi = 50+(rspm-40)*(50/40)
    elif (rspm>80 and rspm<=180):
        rspmi = 100+(rspm-80)*(100/100)
    elif (rspm>180 and rspm<=280):
        rspmi = 200+(rspm-180)*(100/100)
    elif (rspm>280 and rspm<=400):
        rspmi = 300+(rspm-280)*(100/120)
    else:
        rspmi = 400+(rspm-400)*(100/120)
    return rspmi
```

```
data['rspmi']=data['RSPM/PM10'].apply(calculate_rspmi)
df=data[['RSPM/PM10','rspmi']]
df.head()
```

	S02	si
0	11.0	13.75
1	13.0	16.25
2	12.0	15.00
3	15.0	18.75
4	13.0	16.25

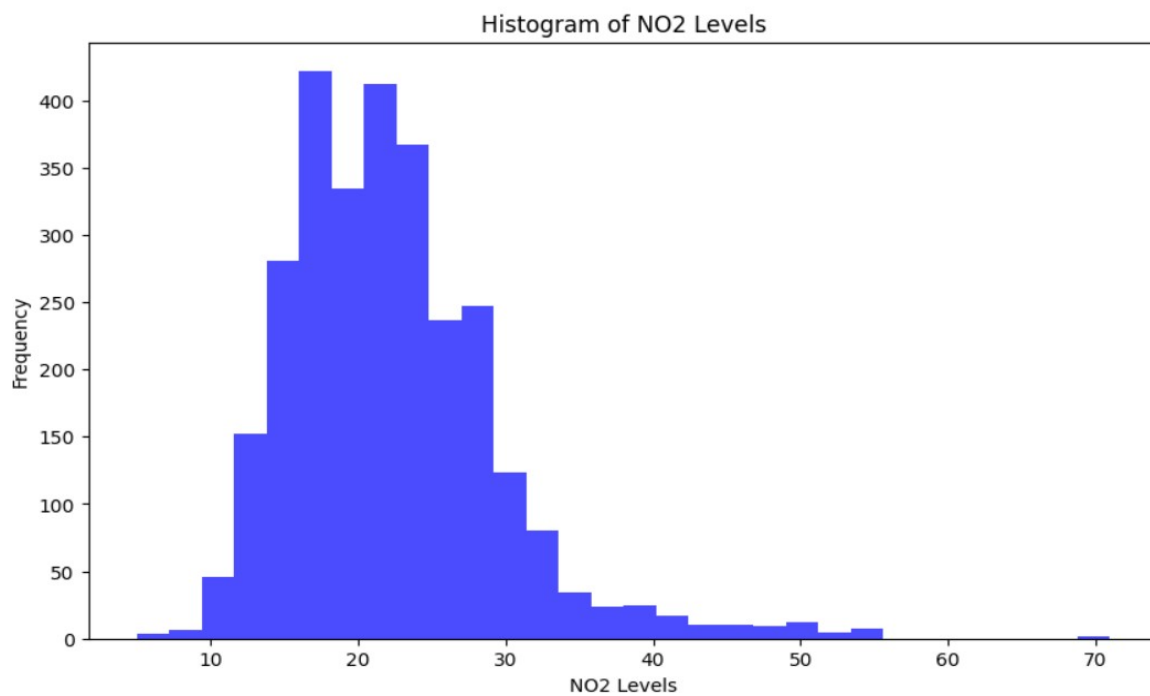
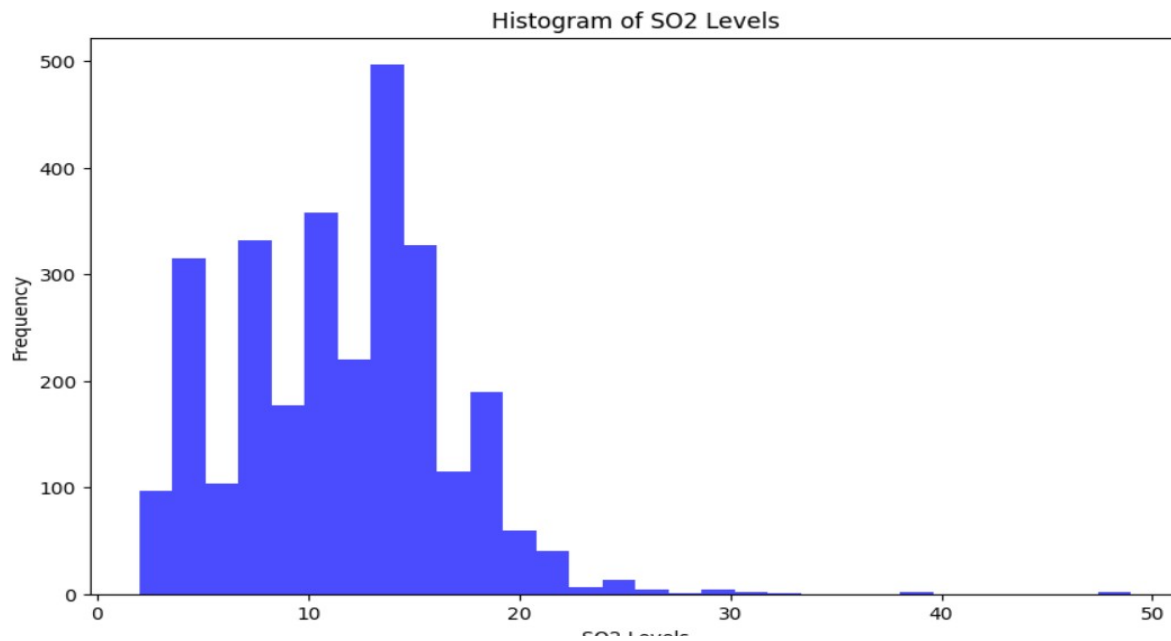
```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns

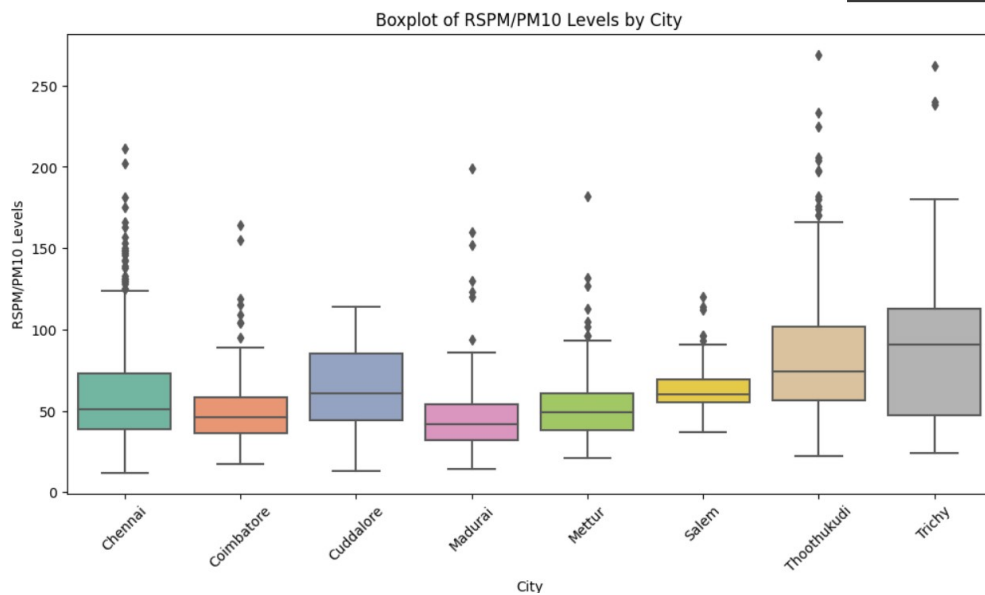
# Load the dataset
file_path = r"/content/cpcb_dly_aq_tamil_nadu-2014 (1).csv"
data = pd.read_csv(file_path)

# Create visualizations
# Example 1: Histogram of SO2 levels
plt.figure(figsize=(10, 6))
plt.hist(data['SO2'], bins=30, color='blue', alpha=0.7)
plt.title("Histogram of SO2 Levels")
plt.xlabel("SO2 Levels")
plt.ylabel("Frequency")
plt.show()

# Example 2: Time series plot of NO2 levels
plt.figure(figsize=(10, 6))
plt.hist(data['NO2'], bins=30, color='blue', alpha=0.7)
plt.title("Histogram of NO2 Levels")
plt.xlabel("NO2 Levels")
plt.ylabel("Frequency")
plt.show()

# Example 3: Boxplot of RSPM/PM10 levels by City
plt.figure(figsize=(12, 6))
sns.boxplot(data=data, x="City/Town/Village/Area", y="RSPM/PM10", palette="Set2")
plt.title("Boxplot of RSPM/PM10 Levels by City")
plt.xlabel("City")
plt.ylabel("RSPM/PM10 Levels")
plt.xticks(rotation=45)
plt.show()
```





Insights into Air Pollution Trends

The analysis provides several key insights:

The analysis revealed several significant findings:

- **Seasonal Patterns:** The data clearly showed the presence of seasonal variations in air pollution levels. Understanding these patterns can be instrumental in developing strategies to mitigate pollution during peak seasons.
- **Spatial Variation:** Variations in air quality across different locations within Tamil Nadu were evident. Urban areas, in particular, exhibited higher pollution levels, emphasizing the need for localized interventions.
- **Correlation Insights:** The identification of strong correlations between specific pollutants suggests common sources or interactions among pollutants, which can guide pollution control measures.
- **Health Impact Awareness:** The analysis indicated that high levels of certain pollutants coincide with an increase in health-related issues. This underscores the critical importance of addressing air pollution for public health and well-being.

Conclusion:

In conclusion, the Air Quality Analysis project conducted for Tamil Nadu in the year 2014 has provided valuable insights into the state's air pollution trends and pollution levels. The project's objectives were successfully achieved through a comprehensive analysis approach, data preprocessing, and the application of various visualization techniques. Overall, the Air Quality Analysis project for Tamil Nadu contributes to a better understanding of air pollution in the region, offering valuable information to policymakers, environmentalists, and the public. These insights can be used to develop targeted strategies and policies to improve air quality, mitigate health risks, and protect the environment. This project serves as a foundation for further research and actions aimed at addressing the challenges posed by air pollution in Tamil Nadu.