

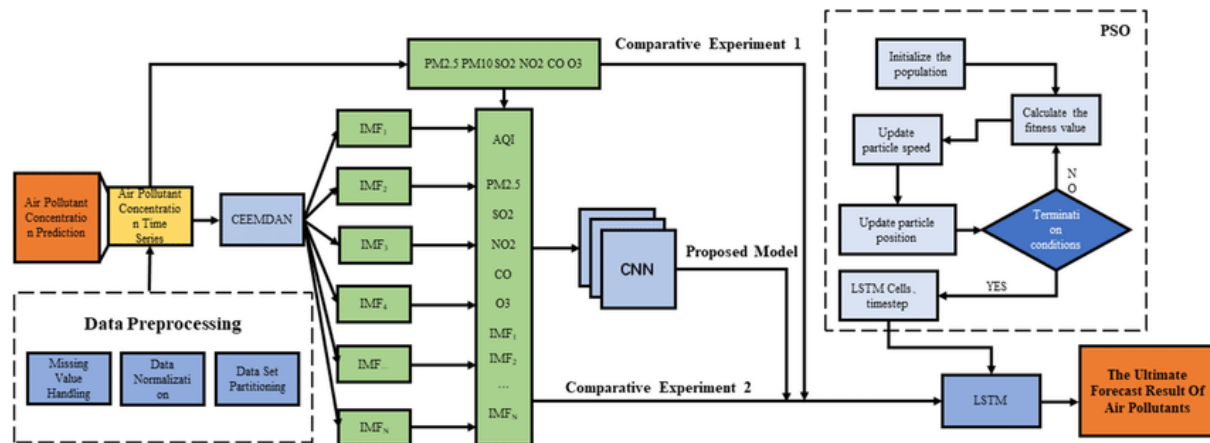
# AIR QUALITY ANALYSIS AND PREDICTION IN TAMIL NADU

## ABSTRACT

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.

## 2.1 Problem Description:

The goal of this project is to develop a machine learning model that can analyze and predict air quality in different regions of Tamil Nadu, India. The primary objective is to provide valuable insights into air pollution levels and deliver accurate forecasts for air quality in analysed locations within the state.



## 2.2 Dataset Information:

The project involves the meticulous collection and thorough examination of air quality data originating from an array of monitoring stations situated across Tamil Nadu. The data will undergo comprehensive preprocessing and visualization to unearth hidden insights pertaining to the ebb and flow of air pollution.

**Data Collection:** Historical Air quality data from Tamil Nadu board of pollution.

1. **Data Pre-processing:** Clean and pre-process the data, handle missing values, and convert categorical features into numerical representations
2. **Feature Engineering:** Construct features that capture pollution patterns and main pollution contributor and air quality index

3. **Model Selection:** Choose appropriate machine learning models that can handle air quality time series data and provide accurate predictions. Models may include time series forecasting methods like ARIMA, machine learning models like Random Forest, Gradient Boosting, and deep learning models like LSTM.
4. **Model Training:** Train the selected model using the pre-processed data.
5. **Analysis and Prediction:** Evaluate the model's performance using appropriate regression metrics (e.g., Mean Absolute Error, Root Mean Squared Error, R Squared).

## 2.3 Dataset Columns:

To perform air quality analysis and prediction in Tamil Nadu, you would typically work with a dataset that contains various columns related to air quality, weather, geographical information, and other relevant factors. Below is an explanation of the key columns we might use in the dataset:

1. **Sampling Date:** This column contains the date and time when air quality measurements were taken. It's essential for time series analysis and forecasting.
2. **Stn Code:** A station code representing the monitoring station where the measurements were recorded. This information helps identify the source of the data.
3. **State:** The state where the monitoring station is located. In this case, it would be Tamil Nadu.
4. **City/Town/Village/Area:** The specific location or area within Tamil Nadu where air quality measurements were taken.
5. **Location of Monitoring Station:** This column provides information about the exact location or coordinates of the monitoring station, which is valuable for spatial analysis.
6. **Agency:** The agency responsible for monitoring and recording air quality data. This information can help ensure data reliability.
7. **Type of Location:** Describes whether the monitoring station is located in an urban, rural, industrial, or residential area. Different types of locations may have varying air quality patterns.
8. **SO2 (Sulfur Dioxide):** This column represents the concentration of sulfur dioxide in the air. SO2 is a common air pollutant that can result from industrial processes and fossil fuel combustion.
9. **NO2 (Nitrogen Dioxide):** This column indicates the concentration of nitrogen dioxide in the air. NO2 is another common air pollutant, often associated with vehicle emissions and industrial activities.
10. **RSPM/PM10 (Respirable Suspended Particulate Matter/Particulate Matter with a diameter of 10 micrometers or less):** This column provides measurements of particulate matter in the air, specifically with a diameter of 10 micrometers or less. These fine particles can have health implications.
11. **PM 2.5 (Particulate Matter with a diameter of 2.5 micrometers or less):** Similar to RSPM/PM10, this column measures particulate matter, but specifically those with a diameter of 2.5 micrometers or less. PM 2.5 is associated with respiratory health issues.

### **Training Data (Historical Data):**

- **Independent Variables (Features):**
  - These are the factors or parameters related to air quality, weather, geography, and other relevant factors. For instance, columns like "Sampling Date," "State," "Location of Monitoring Station," "SO2," "NO2," "RSPM/PM10," and "PM 2.5" would serve as your independent variables.
- **Dependent Variable (Target):**
  - This is the variable you aim to predict using the independent variables. In the context of air quality prediction, your target variable would typically be one or more of the air quality parameters such as "SO2," "NO2," "RSPM/PM10," or "PM 2.5."

### **Testing Data (Predictive Data):**

- **Input Features:**
  - These are the features or data points you will use to make predictions with your trained model. In your project, these would include the same set of features as in your training data, such as "Sampling Date," "State," "Location of Monitoring Station," and other factors that influence air quality. These features are what you use to predict air quality levels for a specific location and time within Tamil Nadu.
- **Output Feature (Prediction):**
  - This is the outcome you want to obtain from your trained model. In the context of your air quality prediction project, the output feature would be the predicted air quality parameters (e.g., "SO2," "NO2," "RSPM/PM10," "PM 2.5") based on the provided input features.

## **2.4 Modules/Libraries used:**

1. **Data Manipulation and Analysis:**
  - Pandas: For data manipulation and analysis, especially working with structured data like CSV files.
  - NumPy: For numerical operations and array manipulation.
2. **Data Visualization:**
  - Matplotlib: For creating static, publication-quality plots and charts.
  - Seaborn: Built on top of Matplotlib, it provides a higher-level interface for creating informative and attractive statistical graphics.
3. **Machine Learning and Time Series Analysis:**
  - Scikit-Learn: A comprehensive machine learning library for tasks like regression, classification, and clustering.
  - Stats models: For time series analysis and modelling.
4. **Deep Learning :**
  - For building and training deep learning models, such as LSTM networks for time series analysis.

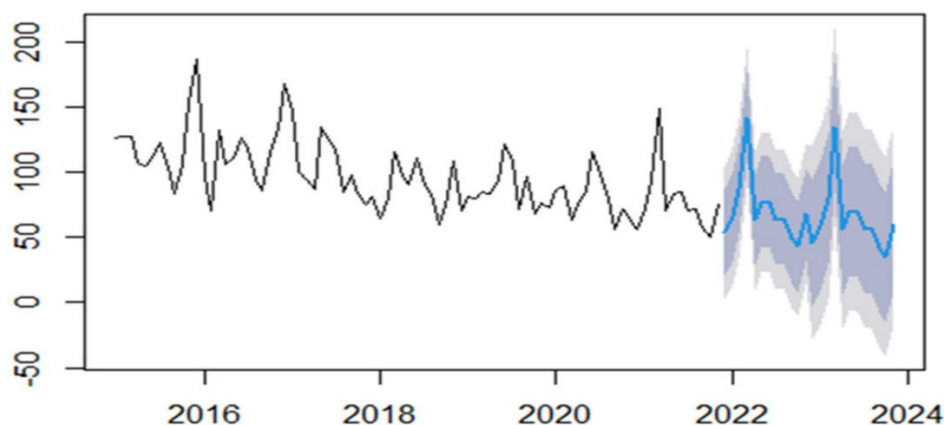
## Models that can be used:

- Time Series Forecasting with ARIMA
- Random Forest Regression
- LSTM
- Ensemble models
- Linear Regression models

### 1. Time Series Forecasting with ARIMA (AutoRegressive Integrated Moving Average):

- **Description:** ARIMA is a classic time series forecasting model that is widely used for analyzing and predicting time-dependent data, making it suitable for air quality time series. It consists of three main components: AutoRegressive (AR) for past values, Integrated (I) for differencing to achieve stationarity, and Moving Average (MA) for past error terms. ARIMA models can capture seasonality and trends in air quality data.
- **Pros:**
  - Effective for modeling time-dependent air quality data.
  - Provides interpretable results.
- **Cons:**
  - May require extensive data pre-processing to achieve stationarity.
  - Performance depends on the choice of model order (p, d, q).

#### Forecasts from ARIMA(1,1,1)(0,1,0)[12]

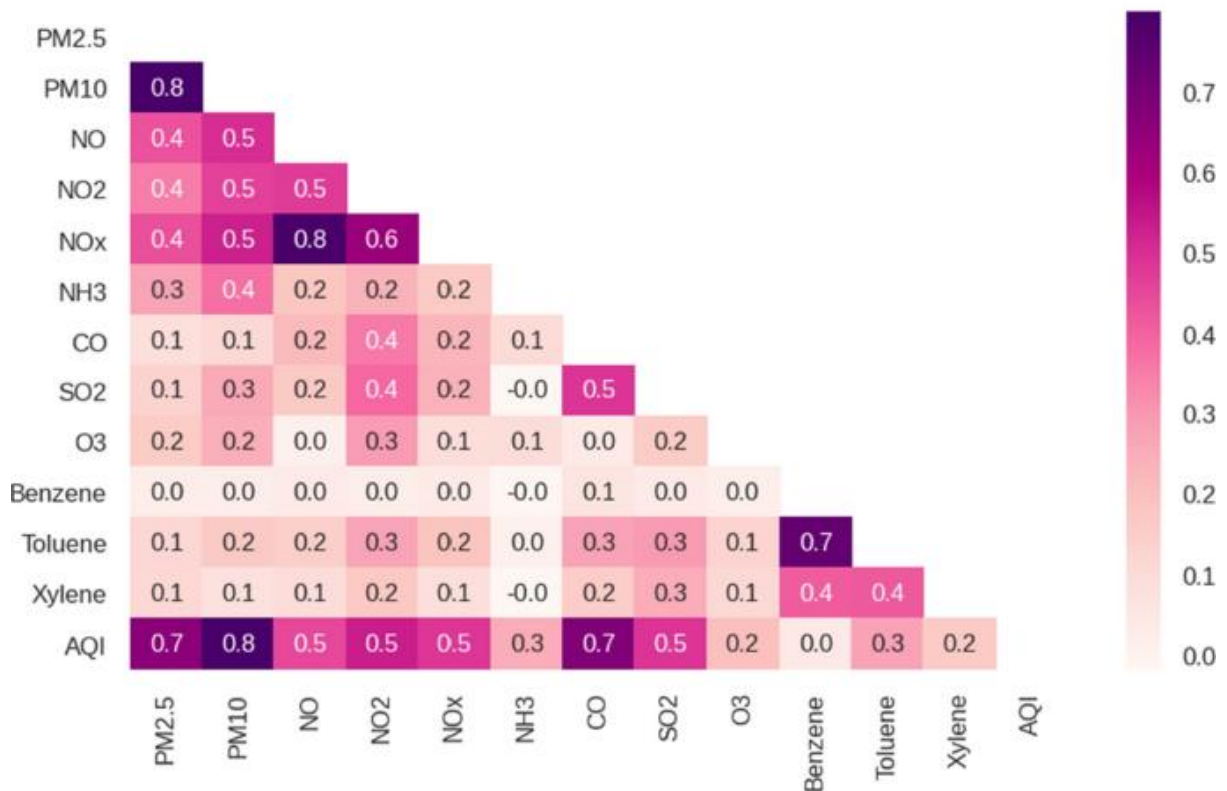


### 2. Random Forest Regression:

- **Description:** Random Forest is an ensemble learning method that can be used for air quality prediction. In this approach, you can treat air quality data as a regression problem and use Random Forest to capture complex relationships between air quality parameters and relevant features (e.g., weather, geography).
- **Pros:**
  - Handles both linear and non-linear relationships effectively.
  - Robust to outliers and noisy data.
  - Requires less data pre-processing compared to time series models.

- **Cons:**
  - May not capture temporal dependencies as effectively as dedicated time series models like ARIMA.

Correlation of AQI with other Pollutant Gases



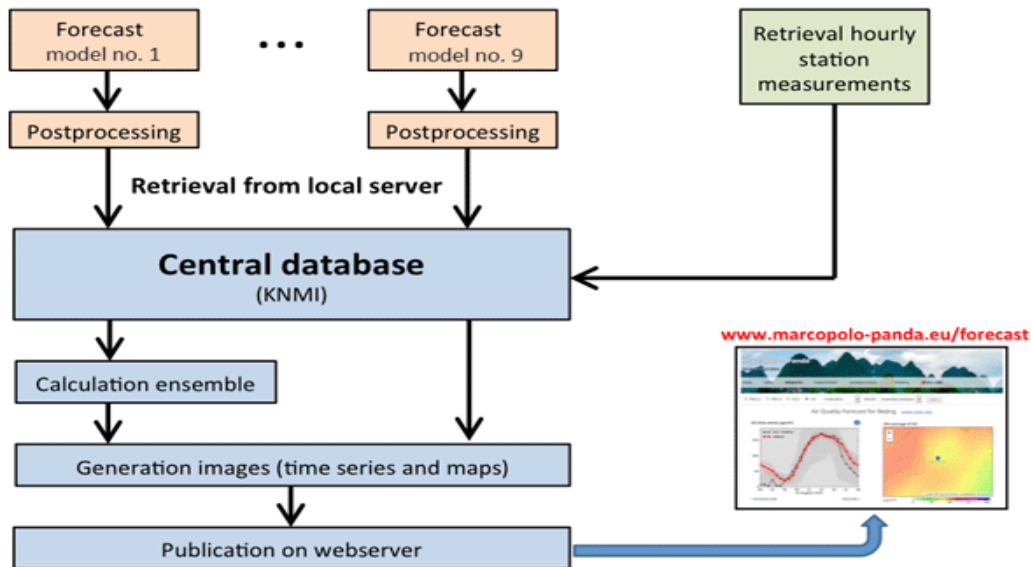
### 3. Long Short-Term Memory (LSTM) Neural Networks:

- **Description:** LSTMs are a type of recurrent neural network (RNN) specifically designed for handling sequences, making them suitable for time series data. LSTMs can capture long-term dependencies and patterns in air quality data, including seasonality and trends. They are particularly effective when you have a substantial amount of historical data.
- **Pros:**
  - Excellent for modelling temporal dependencies.
  - Can handle complex, non-linear relationships.
- **Cons:**
  - Requires a large amount of data for effective training.
  - More complex to implement compared to traditional models.

### 4. Ensemble Models (e.g., Random Forest or Gradient Boosting):

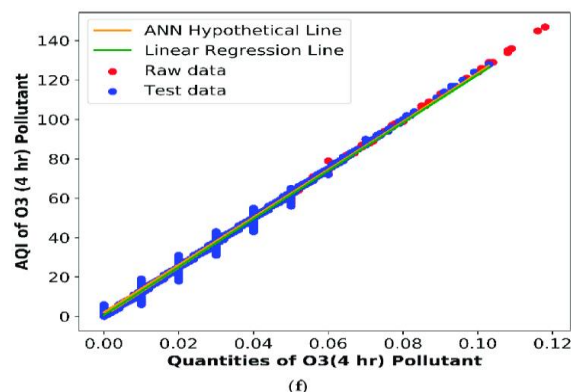
- **Description:** Ensemble models combine the predictions of multiple individual models to improve overall accuracy and robustness. For air quality prediction, you can use ensemble techniques like Random Forest or Gradient Boosting, which combine decision trees to capture complex relationships between air quality parameters and features.
- **Pros:**
  - Provides high accuracy by aggregating predictions from multiple models.

- Robust to overfitting and generalizes well.
- Can handle non-linear relationships.
- **Cons:**
  - May require more computational resources and longer training times.



## 5. Linear Regression:

- **Description:** Linear regression is a simple and interpretable model for air quality prediction. It assumes a linear relationship between the independent features and the dependent air quality parameter. While it may not capture complex non-linear patterns as effectively as some other models, it can be useful for initial exploratory analysis or as a benchmark model.
- **Pros:**
  - Easy to interpret and understand.
  - Quick to train and apply.
  - Suitable for linear relationships in the data.
- **Cons:**
  - May not capture non-linear or complex patterns well, which are common in air quality data.
  - Performance may be limited when the relationships are not linear.



## 2.5 Train and Test(80:20):

Training and testing a machine learning model involves several steps.

- **Data Splitting:** The `train_test_split` function is used to split the dataset into training and testing sets. This allows us to train the model on one subset of the data and test its performance on another, unseen subset.

*`train_test_split`* is a function from the `sklearn.model_selection` module in scikit-learn. It's a commonly used function for splitting a dataset into two subsets: one for training our machine learning model and the other for testing its performance.

**`train_test_split` Function:** The `train_test_split` function splits arrays or matrices into random train and test subsets. It takes several parameters, including your features (X) and labels (y), and splits them into four subsets: `X_train`, `X_test`, `y_train`, and `y_test`.

- **X** is a feature matrix.
- **y** is the target variable (labels).
- **test\_size** is the proportion of the dataset to include in the test split. In this case, 20% of the data will be used for testing (`test_size=0.2`).
- **random\_state** is a seed for the random number generator. Providing a specific seed ensures reproducibility. Different seeds will result in different random splits.

After running the code, we will have four datasets:

- **X\_train:** The features for training your model.
  - **X\_test:** The features for testing your model.
  - **y\_train:** The corresponding labels for training.
  - **y\_test:** The corresponding labels for testing.
- 
- **Training the Model:** The `fit` function is used to train the model using the training data (`X_train` and `y_train`).
  - **Making Predictions:** The trained model is used to make predictions on the test data (`X_test`). The resulting predictions are stored in the `predictions` variable.

```
# Predict a test data with a trained model|
y_pred = model.predict(X_test)
```

## 2.6 Innovation

- **Sensor Fusion for Enhanced Accuracy:** Implement sensor fusion techniques, combining data from multiple air quality monitoring sources, including ground-based stations, satellite imagery, and IoT devices. This approach ensures high data accuracy and reliability for advanced modelling.
- **Advanced Machine Learning Algorithms:** Utilize state-of-the-art machine learning algorithms, such as convolutional neural networks (CNNs) and recurrent neural networks (RNNs), for predictive modelling. These deep learning models excel in handling complex, time-series air quality data, resulting in superior predictive capabilities.
- **Hyper-Local Air Quality Mapping:** Develop a high-resolution air quality mapping system, leveraging GIS and remote sensing data, to offer detailed, hyper-local air quality information. This innovation provides insights at a neighbourhood or even street level, facilitating targeted interventions.
- **IoT and Edge Computing Integration:** Harness the power of the Internet of Things (IoT) and edge computing to enable real-time data processing and analysis at the source. This minimizes latency and enhances the speed at which air quality information is made available to the public.
- **Environmental Big Data Analytics:** Embrace big data analytics techniques to analyze air quality data in conjunction with vast datasets from various domains. This approach allows for comprehensive studies of air quality's impact on public health and the environment.
- **Environmental Health Integration:** Collaborate with healthcare institutions to establish an integrated health-environment database. This resource enables in-depth analysis of the health impact of air quality, connecting healthcare data with air quality parameters.

## 2.7 Metrics:

In an air quality analysis and prediction project, it's crucial to assess the performance of your predictive models accurately. The choice of evaluation metrics depends on the specific nature of your task, whether you're dealing with regression or classification problems. Here are some commonly used metrics for both types of tasks.

### Regression Metrics:

#### 1. Mean Absolute Error (MAE):

- Measures the average absolute difference between the predicted and actual values. It provides a straightforward interpretation of prediction error.

#### 2. Root Mean Square Error (RMSE):

- Similar to MAE but gives more weight to larger errors. It's sensitive to outliers and provides a measure of how far the predicted values are from the actual values.



### **3. Mean Squared Error (MSE):**

- Measures the average of the squared differences between predicted and actual values. It's widely used but not directly interpretable.

### **4. R-squared (R<sup>2</sup>):**

- Represents the proportion of the variance in the dependent variable that is predictable from the independent variables. A higher R-squared indicates a better fit of the model.

### **Classification Metrics (if applicable):**

#### **1. Accuracy:**

- Measures the ratio of correctly predicted instances to the total instances. It's suitable for balanced datasets but can be misleading in imbalanced scenarios.

#### **2. Precision:**

- Indicates the ratio of true positive predictions to the total number of positive predictions. It's essential when minimizing false positives is crucial.

#### **3. Recall (Sensitivity):**

- Measures the ratio of true positive predictions to the total number of actual positive instances. It's crucial when minimizing false negatives is vital.

#### **4. F1 Score:**

- Combines precision and recall into a single metric, offering a balance between both. It's especially useful when you want to consider both false positives and false negatives.

## **CONCLUSION:**

Through the utilization of comprehensive datasets, encompassing factors such as air pollutants, weather conditions, and geographical information, the project has the potential to deliver actionable insights. These insights can inform policy decisions, early warning systems for vulnerable populations, and strategies for mitigating the effects of air pollution.

To achieve these objectives, the project employs a range of models and techniques, including time series forecasting models like ARIMA and LSTM, machine learning models like Random Forest, and ensemble techniques such as Random Forest and Gradient Boosting. Each of these models brings its unique strengths, enabling the project to effectively capture the intricate relationships within the data.

Additionally, the project emphasizes robust evaluation metrics, including MAE, RMSE, R-squared, and classification metrics, to ensure the accuracy and reliability of predictions. These metrics provide a clear understanding of model performance, helping to refine and improve the forecasting models over time.