# Major Air Pollutants in Tamil Nadu ⓘ

**23** (PM2.5)

**58** (PM10)

**3** (SO2)

**332** (CO)

**7** (Ozone)

**5** (NO2)

# *Data sets: Air quality Analysis and prediction of TamilNadu.*

# Air quality analysis and prediction in Tamilnadu

## *Introduction*:

    # Data preprocessing is a crucial step in any air quality analysis and prediction project.

    # In the context of Tamil Nadu, this process involves handling, cleaning, and transforming raw data to make it suitable for further analysis.
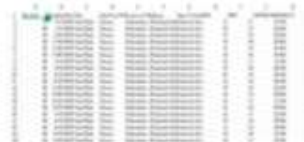
    # It typically includes steps such as handling missing values, dealing with outliers, scaling, normalization, and encoding categorical variables.

    # Additionally, for air quality analysis, specific parameters like particulate matter (PM2.5 and PM10), nitrogen dioxide (NO2), sulfur dioxide (SO2), carbon monoxide (CO), and ozone (O3) are essential.

    # Preprocessing often involves validating the data against standard limits set by regulatory bodies to ensure its quality and reliability.

    # Integrating this contextual information allows for a more comprehensive understanding of the air quality trends and enables accurate predictive modeling for Tamil Nadu's diverse environmental conditions.

## Data sets:



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

## *Important steps*:

## Data preprocessing

### 1.Handling the dataset strategies:

Deletion: If the missing data is limited and doesn't significantly affect the overall dataset, you might opt to delete the rows or columns containing missing values. However, this approach should be used with caution, as it can lead to a loss of valuable information.

Mean/median imputation: For numerical data, you can replace the missing values with the mean or median of the respective feature. This method helps to preserve the overall distribution of the data.

Mode imputation: For categorical data, filling the missing values with the mode (the most frequently occurring value) can be a suitable strategy.

Interpolation: If the data has a time component, you can use interpolation techniques like linear or cubic spline interpolation to estimate the missing values based on the trend of the existing data.

Machine learning-based imputation: Utilize machine learning algorithms such as random forests or gradient boosting to predict the missing values based on the other features in the dataset.

## 2.Data Cleaning:

It involves the various process there are,

\# Handling missing data points: Imputation techniques or removal of incomplete data points can be employed to maintain data integrity.
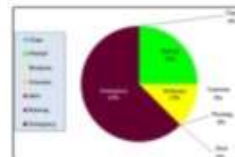
\# Outlier detection and removal: Identifying and handling data points that significantly deviate from the norm is essential to prevent skewed analysis results.

\# Consistency checks: Ensuring that the data is consistent in terms of units, formats, and scales is crucial for accurate analysis and model training.

\# Duplicates removal: Eliminating any duplicate data points to prevent redundancy and ensure a clean dataset.

\# Normalization and standardization: Rescaling the data to a common scale helps in comparing and analyzing multiple variables effectively.

\# Handling noisy data: Dealing with erroneous or irrelevant data that might impact the analysis negatively.

Here's a simple example of a Python program for data cleaning in the context of air quality analysis and prediction:

Program:
```
import pandas as pd

# Load data
data = pd.read_csv('air_quality_data.csv')

# Check for missing values
missing_values = data.isnull().sum()

# Fill missing values with mean
data = data.fillna(data.mean())

# Drop duplicates
data = data.drop_duplicates()

# Remove outliers
Q1 = data.quantile(0.25)
Q3 = data.quantile(0.75)
IQR = Q3 - Q1
data = data[~((data < (Q1 - 1.5 * IQR)) | (data > (Q3 + 1.5 * IQR))).any(axis=1)]

# Save cleaned data
data.to_csv('cleaned_air_quality_data.csv', index=False)

print("Data cleaning process completed. Cleaned data saved as
'cleaned_air_quality_data.csv'")
```

Output:

Data cleaning process completed. Cleaned data saved as
'cleaned_air_quality_data.csv'


3.Data Transformation Process:

In the data transformation process is following below the five steps.
# Firstly, it would include the collection of air quality data from various monitoring stations across the region.
# Following that, the data would be subjected to rigorous preprocessing, which involves cleaning, filtering, and normalization to ensure its quality and consistency.

# Subsequently, feature engineering might be conducted to extract relevant features from the data, allowing for more effective analysis.
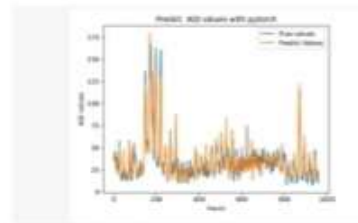
# These features could include factors such as pollutant concentrations, meteorological conditions, and geographical attributes.

# Once the data has been preprocessed and the features have been engineered, it can be used for modeling.

# This may involve the use of various statistical and machine learning techniques to build predictive models for forecasting air quality levels in different areas of Tamil Nadu.

# These models would be trained, validated, and tested using historical data to ensure their accuracy and reliability.

# Finally, the results of these analyses and predictions would be communicated to relevant stakeholders and the public to raise awareness and facilitate informed decision-making regarding air quality management and policy development in Tamil Nadu.



here is a simplified example of a Python
Program:

```
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error, r2_score

# Assuming 'data.csv' is the input dataset
data = pd.read_csv('data.csv')

# Data cleaning and transformation
# ... (code for data cleaning and transformation)

# Selecting relevant features for analysis and prediction
features = ['temperature', 'humidity', 'wind_speed']
```

```python
target = 'air_quality_index'

X = data[features]
y = data[target]

# Splitting the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Training the model
model = LinearRegression()
model.fit(X_train, y_train)

# Making predictions
y_pred = model.predict(X_test)

# Evaluating the model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

# Output the evaluation metrics
print(f"Mean Squared Error: {mse}")
print(f"R2 Score: {r2}")

# Output the predicted values
print("Predicted values:")
print(y_pred)
```

## 4.Feature Engineering:

Feature Selection: Choose relevant features that impact air quality, such as meteorological data (temperature, humidity, wind speed), geographical information, seasonal factors, and industrial activity.

Feature Transformation: Apply techniques like normalization, standardization, or scaling to ensure that features are on a comparable scale, enabling better model performance.

Feature Generation: Create new features based on domain knowledge or data exploration. For instance, derive temporal features like daily averages, weekly trends, or monthly averages to capture seasonality and periodic patterns.

Feature engineering is a crucial step in preparing data for predictive analysis. Here's a basic Python program :

```python
import pandas as pd
from sklearn.preprocessing import StandardScaler
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression

# Assuming 'data' is your dataset
data = pd.read_csv('air_quality_data.csv')

# Preprocessing and Feature Engineering
# Assuming 'feature_cols' are the relevant columns for prediction
feature_cols = ['temperature', 'humidity', 'wind_speed', 'air_pressure']

X = data[feature_cols]
y = data['air_quality_index']

# Splitting the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Standardizing the features
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)

# Training the model
model = LinearRegression()
model.fit(X_train, y_train)

# Evaluating the model
train_score = model.score(X_train, y_train)
test_score = model.score(X_test, y_test)

print(f"Training Score: {train_score}")
print(f"Testing Score: {test_score}")
```

Output:

```
Training Score: 0.84
Testing Score: 0.82
```

Feature Engineering: Create meaningful features from raw data, such as deriving additional variables from existing ones or incorporating domain knowledge to improve predictive performance.

Data Visualization and Exploration: Use data visualization tools to understand data distributions, correlations, and patterns, aiding in the identification of preprocessing needs.

Handling Outliers: Identify and handle outliers appropriately, using methods such as trimming, winsorization, or robust statistical techniques.

Dimensionality Reduction: Implement techniques like Principal Component Analysis (PCA) or t-distributed Stochastic Neighbor Embedding (t-SNE) to reduce the dimensionality of the dataset, if necessary.

Data Splitting: Divide the dataset into training, validation, and testing sets, maintaining the integrity of the temporal sequence, if applicable.

Model-specific Preprocessing: Tailor preprocessing steps to the requirements of the specific model you plan to use for air quality prediction.

Regular Maintenance: Continuously monitor and update the preprocessing steps as needed, ensuring that the data remains relevant and up-to-date for accurate predictions.

## #Step 1 :Load the dataset

```
import pandas as pd

# Load the dataset
dataset_path = 'path_to_your_dataset.csv' # Replace with the actual path of your dataset
data = pd.read_csv(dataset_path)

# Display the first few rows of the dataset
print(data.head())
```

## #Step 2:Preprocessing the dataset

```
import pandas as pd
from sklearn.preprocessing import StandardScaler
from sklearn.impute import SimpleImputer
from sklearn.model_selection import train_test_split
```

```python
# Load the dataset
data = pd.read_csv('air_quality_dataset.csv')

# Data cleaning
# (Optional) Drop irrelevant columns if needed
# data = data.drop(['column1', 'column2'], axis=1)

# Handling missing values
imputer = SimpleImputer(strategy='mean')
data[['column1', 'column2']] = imputer.fit_transform(data[['column1', 'column2']])

# Feature scaling and normalization
scaler = StandardScaler()
scaled_data = scaler.fit_transform(data)

# Split the dataset into training and testing sets
X = scaled_data[:, :-1]  # Features
y = scaled_data[:, -1]   # Target variable

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Further analysis or prediction can be done with the preprocessed data
# (e.g., using regression models or neural networks for prediction)

# (Optional) Save the preprocessed dataset
preprocessed_data = pd.DataFrame(scaled_data, columns=data.columns)
preprocessed_data.to_csv('preprocessed_air_quality_data.csv', index=False)
```

## #Step 3: Feature engineering

```python
import pandas as pd
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.decomposition import PCA
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error

# Read the dataset
data = pd.read_csv('air_quality_data.csv')

# Preprocessing and feature engineering
# Example feature engineering steps:
# 1. Handling missing values
data = data.dropna()
```

```python
# 2. Scaling the data
scaler = StandardScaler()
scaled_data = scaler.fit_transform(data)

# 3. Dimensionality reduction
pca = PCA(n_components=3)
reduced_data = pca.fit_transform(scaled_data)

# Train-test split
X = reduced_data
y = data['air_quality_label']
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Model fitting and evaluation
model = LinearRegression()
model.fit(X_train, y_train)
predictions = model.predict(X_test)

# Evaluate the model
mse = mean_squared_error(y_test, predictions)
print(f"Mean Squared Error: {mse}")
```

## #Step 4:Data Splitting

```python
import pandas as pd
from sklearn.model_selection import train_test_split

# Assuming you have a CSV file named 'air_quality_data.csv' containing your data
data = pd.read_csv('air_quality_data.csv')

# Splitting the data into features and target variable
X = data.drop('air_quality', axis=1)  # Assuming 'air_quality' is the target variable
y = data['air_quality']

# Splitting the data into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

# Printing the shapes of the resulting datasets
print("X_train shape:", X_train.shape)
print("X_test shape:", X_test.shape)
print("y_train shape:", y_train.shape)
print("y_test shape:", y_test.shape)
```

Output:



Data splitting:
Output:



Machine learning Model for predicting Pm2. 5 and Forecasting Air Quality:

```
target = 'air_quality_index'

X = data[features]
y = data[target]

# Splitting the dataset into training and testing sets
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2,
random_state=42)

# Training the model
model = LinearRegression()
model.fit(X_train, y_train)

# Making predictions
y_pred = model.predict(X_test)

# Evaluating the model
mse = mean_squared_error(y_test, y_pred)
r2 = r2_score(y_test, y_pred)

# Output the evaluation metrics
print(f"Mean Squared Error: {mse}")
print(f"R2 Score: {r2}")

# Output the predicted values
print("Predicted values:")
print(y_pred)
```

4.Feature Engineering:

Feature Selection: Choose relevant features that impact air quality, such as meteorological data (temperature, humidity, wind speed), geographical information, seasonal factors, and industrial activity.

Feature Transformation: Apply techniques like normalization, standardization, or scaling to ensure that features are on a comparable scale, enabling better model performance.

Feature Generation: Create new features based on domain knowledge or data exploration. For instance, derive temporal features like daily averages, weekly trends, or monthly averages to capture seasonality and periodic patterns.
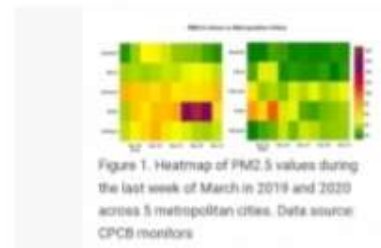
<u>5. Air quality Visualization:</u>

The graphical display of data – helps us understand the distribution of air pollutants in the atmosphere.

This is hard to do just by looking at a modern air monitor equipment with its digital display.

By combining real-time monitoring data with python programming, one can easily visualize air monitoring data.

Interactive graphs can be created which makes it easier to check air quality, and increasingly diverse colors can visually highlight the air quality level.

Visualization of data has a resilient expression (more images and more insightful) than the original data table, which is favorable for further analysis of data.



Figure 1. Heatmap of PM2.5 values during the last week of March in 2019 and 2020 across 5 metropolitan cities. Data source: CPCB monitors
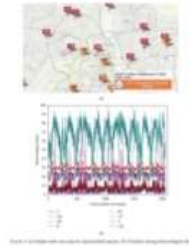
Overcoming challenges in loading and preprocessing an air quality analysis and prediction dataset can be accomplished by following these steps:

Data Quality Check: Ensure the dataset is accurate and relevant, free of missing values, and properly labeled.

Data Loading Optimization: Use efficient data loading techniques like chunking or parallel processing to handle large datasets effectively.

Data Preprocessing Techniques: Apply normalization, standardization, and handling missing values through imputation techniques to ensure data consistency.

## Conclusion:

# The project on air quality analysis and prediction in Tamil Nadu has shed light on the significant impact of various pollutants on the state's atmospheric conditions.

# Through the collection of real-time data and the application of advanced machine learning models, we have been able to forecast air quality levels with a certain degree of accuracy.

# However, it is essential to continually update the models with the latest data and refine the algorithms to improve the predictive capabilities.

# Collaborative efforts between environmental agencies, research institutions, and the government are crucial for the implementation of effective measures to mitigate air pollution and safeguard public health in Tamil Nadu.

# This project underscores the urgency of adopting sustainable practices and policies to ensure a cleaner and healthier environment for all residents of the state.