Vinay Pawar TYITB098

## Assignment No. 5

**Problem Statement:**
Perform data cleaning and data transformation operations on air quality datasets using Python.

**Objective:**
The goal is to preprocess and transform air quality data in Python to make it clean, consistent, and ready for analysis. This involves removing missing or inconsistent values and structuring the dataset to make it suitable for analysis and visualization.

**Prerequisites:**

1. **Python Basics** and familiarity with data manipulation and visualization libraries such as Pandas, NumPy, and Matplotlib.
2. **Data Preprocessing Techniques**: Handling missing values, detecting outliers, and transforming data for consistency.
3. **Air Quality Indicators**: Understanding the significance of indicators like PM2.5, PM10, NO2, and O3.

---

# Theory:

Data cleaning and transformation are foundational steps in any data processing pipeline. These processes involve preparing raw data into a structured format, making it suitable for analysis.

**Data Cleaning:**
Data cleaning focuses on identifying and correcting inconsistencies, missing values, and anomalies to make data accurate and reliable. In air quality datasets, common issues include:

- **Missing Values**: Caused by sensor malfunctions; can be handled using techniques like mean or median imputation, forward or backward filling, or row removal.
- **Outliers**: Extreme values may skew analysis; they are often identified using statistical methods (like Z-scores or IQR) or domain knowledge.
- **Inconsistent Formatting**: Variations in formats for dates, units, and other entries can disrupt analysis.

**Data Transformation:**
Data transformation tailors data for analytical requirements. This typically includes scaling, aggregating, and feature engineering.

- **Scaling and Normalization**: Ensures all numeric values fall within the same range, often improving the accuracy of models.

- **Encoding**: Converts categorical variables into numeric form (e.g., encoding regions for further analysis).
- **Aggregation**: Reduces highly granular data, for instance, hourly data can be averaged daily for broader insights.
- **Feature Engineering**: Deriving new features from existing data, such as calculating AQI or aggregating pollutant concentrations.

---

## Algorithm:

### Data Cleaning:

1. **Import Required Libraries**: Load necessary libraries like Pandas, NumPy, and Matplotlib.
2. **Load the Dataset**: Use Pandas to load the air quality dataset.
3. **Identify and Handle Missing Values**: Use imputation techniques or remove rows with missing values.
4. **Remove Duplicates**: Check for duplicated rows and eliminate them.
5. **Detect and Remove Outliers**: Use statistical methods to identify and address outliers.
6. **Standardize Formats**: Ensure consistency in formats, such as dates and units.

### Data Transformation:

1. **Scale Numeric Features**: Apply scaling (e.g., MinMaxScaler or StandardScaler) if needed.
2. **Encode Categorical Variables**: Convert categorical features to numeric format for analysis.
3. **Aggregate Data**: Summarize data to reduce granularity (e.g., calculate daily averages from hourly data).
4. **Create New Features**: Generate relevant features, like aggregating pollutant values for AQI.

---

## References:

1. McKinney, W. (2017). *Python for Data Analysis: Data Wrangling with Pandas, NumPy, and IPython.* O'Reilly Media.
2. VanderPlas, J. (2016). *Python Data Science Handbook: Essential Tools for Working with Data.* O'Reilly Media.
3. Pandas Documentation: Comprehensive guide for data manipulation using Pandas.

---

**Conclusion:**
Data cleaning and transformation are critical steps in preparing air quality data for accurate analysis. These processes ensure data reliability, enhance the accuracy of analytical models, and provide deeper insights into air quality trends. Python's data manipulation libraries streamline these preprocessing steps, laying a solid foundation for meaningful data analysis and potential machine learning applications.