

## Assignment No. 5

**Problem Statement:** Perform the following operations using Python on the Air quality data sets  
a. Data cleaning b. Data transformation

**Objective:** The objective of this project is to perform data cleaning and transformation on air quality datasets using Python. This process aims to enhance the quality of the data by removing inaccuracies and preparing it for analysis.

### Prerequisite :

1. Basic Python Knowledge: Familiarity with Python programming.
2. Libraries: Knowledge of data manipulation libraries such as:
  - Pandas: For data manipulation and analysis.
  - NumPy: For numerical operations.
3. Understanding of Data Cleaning and Transformation Concepts: Familiarity with common techniques used in data cleaning and transformation.

### Theory :

#### Data Cleaning

Data cleaning, also known as data cleansing or scrubbing, is a fundamental step in the data preprocessing phase that ensures the quality and integrity of the dataset. Poor data quality can lead to inaccurate analyses and misinformed decisions, making data cleaning essential for any data-driven project. Key aspects of data cleaning include:

#### 1. Missing Values:

- **Identification:** Missing values can occur due to various reasons, such as sensor malfunctions, data entry errors, or non-responses in surveys. Techniques to identify missing values include using functions like `.isnull()` or `.isna()` in Pandas.
- **Handling Strategies:**
  - **Imputation:** This involves replacing missing values with substitutes. Common methods include using the mean, median, or mode of the respective feature. Advanced techniques include K-Nearest Neighbors (KNN) imputation or predictive models to estimate missing values.
  - **Removal:** If the percentage of missing data in a feature is high (e.g., >30%), it may be prudent to drop that feature entirely using `.dropna()`.

For rows with missing data, if they are a small fraction of the total dataset, they can be removed.

## 2. Removing Duplicates:

- Duplicated entries can skew analysis and lead to incorrect conclusions. Pandas offers the `.drop_duplicates()` method, which allows users to identify and eliminate duplicate rows based on specified columns or the entire dataset.

## 3. Data Type Correction:

- Each feature in a dataset should have the correct data type to ensure appropriate analysis. For example, air quality metrics like PM2.5 and PM10 should be float types, while categorical data (e.g., pollution category) should be of object type. The `.astype()` function in Pandas can convert data types when necessary.

## 4. Outlier Detection and Handling:

- Outliers are data points that deviate significantly from the rest of the dataset. They can arise from measurement errors or real anomalies. Methods to detect outliers include statistical techniques such as the Z-score method or the Interquartile Range (IQR) method. Handling outliers may involve removing them or transforming them (e.g., using logarithmic transformations).

## 5. Consistency Checks:

- Ensuring consistency across the dataset involves checking for uniformity in data entry formats (e.g., date formats, measurement units). This step is crucial for datasets collected from various sources or sensors.

## Data Transformation

Data transformation is the process of converting data from its original format into a more suitable structure for analysis. This step enhances the dataset's usability, ensuring that the data is compatible with the analytical techniques to be applied. Key aspects of data transformation include:

### 1. Normalization and Standardization:

- **Normalization:** This technique rescales the data to a fixed range, usually [0, 1]. It is particularly useful when features have different units or scales. The Min-Max scaling technique is a common method for normalization.

$$X' = \frac{X - X_{\min}}{X_{\max} - X_{\min}}$$

- **Standardization:** This technique transforms the data to have a mean of 0 and a standard deviation of 1, making it useful for algorithms that assume a Gaussian distribution (e.g., many machine learning models).

$$Z = \frac{X - \mu}{\sigma}$$

## 2. Encoding Categorical Variables:

- Machine learning algorithms typically require numerical inputs. Therefore, categorical variables must be transformed into numerical formats. Common techniques include:
  - **One-Hot Encoding:** This creates binary columns for each category in a categorical feature. For example, a "Pollution Category" with values "High," "Medium," and "Low" would result in three new columns.
  - **Label Encoding:** Assigns a unique integer to each category. While this is simpler, it can introduce an ordinal relationship that might not exist.

## 3. Feature Engineering:

- This process involves creating new features from existing ones to improve model performance. In the context of air quality datasets, this might include:
  - **Temporal Features:** Extracting the hour, day, month, or season from timestamps to analyze temporal patterns in pollution levels.
  - **Lag Features:** Including previous time steps as features to help capture trends over time.

## 4. Aggregation:

- Aggregation involves summarizing data to provide insights into trends and patterns. For instance, daily average pollution levels can be computed from hourly data using the `.groupby()` method in Pandas, allowing for easier interpretation and analysis.

## 5. Data Formatting:

- Proper formatting is crucial for further analysis and visualization. This may involve converting dates to a standard format or ensuring numeric columns are formatted consistently.

## Algorithm (if any to achieve the objective )

### 1. Load the Dataset:

- Import necessary libraries (Pandas, NumPy).
- Load the air quality dataset into a DataFrame.

### 2. Data Cleaning:

- **Identify Missing Values:** Use methods like `.isnull()` to detect missing data.
- **Handle Missing Values:** Decide on a strategy (imputation or removal) and apply it using methods like `.fillna()` or `.dropna()`.
- **Remove Duplicates:** Use the `.drop_duplicates()` method to remove duplicate entries.
- **Correct Data Types:** Ensure that each column has the appropriate data type using methods like `.astype()`.

### 3. Data Transformation:

- **Normalization/Standardization:** Use techniques like Min-Max scaling or Z-score normalization.
- **Encode Categorical Variables:** Use methods like `pd.get_dummies()` for one-hot encoding.
- **Aggregate Data (if necessary):** Use methods like `.groupby()` for aggregation.

### 4. Save the Cleaned Data:

Export the cleaned and transformed dataset to a new file format (e.g., CSV).

## References :

<https://medium.com/womenintechology/data-preprocessing-steps-for-machine-learning-in-python-part-1-18009c6f1153>

<https://medium.com/@yogeshojha/data-preprocessing-75485c7188c4>

<https://medium.com/almabetter/data-preprocessing-techniques-6ea145684812>

<https://medium.easyread.co/basics-of-data-preprocessing-71c314bc7188>

## Conclusion

Data cleaning and transformation are essential processes in data science, particularly in working with air quality datasets. Properly cleaned and transformed data enables accurate analysis and insights, leading to better decision-making. By utilizing Python and its powerful libraries, data scientists can efficiently handle and prepare air quality data for further analysis.

