Vinay Pawar TYITB098

# Assignment No. 6

**Problem Statement:**
Perform data cleaning and error-correction operations on heart disease datasets using Python.

**Objective:**
The aim of this project is to preprocess heart disease data to ensure it is clean, consistent, and ready for analysis. Key goals include:

1. **Data Cleaning**: Remove inconsistencies, handle missing values, and address outliers.
2. **Error Correction**: Identify and correct erroneous data entries to ensure accurate predictive modeling and statistical analysis.

Effective data cleaning and correction enable better understanding of factors associated with heart disease, enhancing predictability and treatment insights.

**Prerequisites:**

1. **Python and Libraries for Data Manipulation**: Knowledge of Python libraries such as Pandas (for data manipulation) and NumPy (for numerical computations), along with Matplotlib or Seaborn for data visualization.
2. **Data Preprocessing Techniques**: Familiarity with data cleaning techniques, including handling missing values, detecting outliers, and correcting inconsistencies.
3. **Understanding of Heart Disease Indicators**: Knowledge of key indicators for heart disease, allowing for accurate interpretation and correction of data errors.

---

## Theory:

Data cleaning and error correction are essential in heart disease data analysis to ensure the data's accuracy and reliability, as errors and inconsistencies could result in inaccurate predictions.

**Data Cleaning**:

- **Missing Data**: Missing values may result from incomplete records. These can be addressed by deletion, imputation (e.g., using mean or median), or interpolation.
- **Outliers**: These could stem from data entry errors or represent rare cases. Outlier detection techniques, like Z-score and IQR, can help identify and decide how to handle them.

- **Duplicates and Consistency**: Duplicate records and inconsistent formats (e.g., date format, units) can distort analysis, but these can be corrected by removing duplicates and ensuring standardized formats.

**Error Correction**:

- **Logical Errors**: Values outside medically realistic ranges, such as negative cholesterol or extremely high blood pressure, can be flagged and corrected.
- **Standardization of Categorical Values**: Ensures consistent labeling of variables, like gender or diagnosis, by mapping categories to uniform labels.
- **Date-Time Consistency**: Proper date-time formatting is essential for tracking patient history and ensuring chronological accuracy.

---

## Algorithm:

**Data Cleaning**:

1. **Load Dataset**: Use Pandas to load the dataset.
2. **Identify Missing Values**: Use `df.isnull().sum()` to locate columns with missing values, then apply imputation or row removal as necessary.
3. **Remove Duplicates**: Use `df.duplicated()` to identify duplicate rows and remove them with `df.drop_duplicates()` if they are redundant.
4. **Outlier Detection**: Apply statistical techniques, like Z-score or IQR, to detect and decide on handling outliers (e.g., truncation or transformation).

**Error Correction**:

1. **Identify Inconsistencies**: Locate impossible values (e.g., negative cholesterol or unrealistic ages) and correct or remove them.
2. **Standardize Formats**: Ensure categorical variables and dates follow consistent formats.
3. **Logic Testing**: Check logical relationships in the data (e.g., high cholesterol values should correlate with heart disease indicators).

---

## References:

1. McKinney, W. (2017). *Python for Data Analysis.* O'Reilly Media.
2. VanderPlas, J. (2016). *Python Data Science Handbook.* O'Reilly Media.
3. [Pandas Documentation](): A comprehensive guide to data manipulation with Pandas.
4. Relevant guidelines on heart disease risk factors (e.g., American Heart Association).

---

**Conclusion:**
The cleaning and correction of heart disease datasets provide a strong foundation for accurate healthcare analytics. By ensuring the data's reliability, we enhance the potential for accurate predictions, enabling more effective healthcare research and decision-making. Python's robust data manipulation libraries facilitate a streamlined approach to data preparation, fostering data integrity essential in healthcare data analysis.