**EXP NO:**    DATA CLEANING: IMPLEMENT VARIOUS MISSING HANDLING MECHANISMS, IMPLEMENT VARIOUS
**DATE:**                    NOISY HANDLING MECHANISMS

**AIM:**

## BACKGROUND THEORY:

### Data Cleaning:

Data cleaning is the process of fixing or removing incorrect, corrupted, incorrectly formatted, duplicate, or incomplete data within a dataset. When combining multiple data sources, there are many opportunities for data to be duplicated or mislabelled.

### Etl Process:

Data Cleaning is an important part of ETL processes as it ensures that only high-quality data is loaded into the Data Warehouse. This helps to improve the accuracy of security decisions.

Data Warehousing is a process of organizing and storing data in a centralized location for easy access and analysis. Data warehousing is used to store historical data from multiple sources in a single location. Data warehouses provide a single view of data that can be used for reporting and analysis. Data warehouses are often used in business intelligence applications.

## PROCEDURE:

### 1. Load and Inspect Data:

o Use the "File" widget to load dataset.csv.

o Connect it to the "Data Table" widget to inspect the data.

### 2. Handle Missing Data:

o Connect the "File" widget to the "Impute" widget.

o In the "Impute" widget, choose "Mean/Median" for numeric features and "Most Frequent" for categorical features.

o Connect the "Impute" widget to a "Data Table" widget to inspect the imputed data.

### 3. Remove Columns with Excessive Missing Data:

o Connect the "Impute" widget to the "Select Columns" widget.

o In the "Select Columns" widget, manually remove columns with a high percentage of

missing values.

### 4. Handle Noisy Data:

o Connect the "Select Columns" widget to the "Smoothing" widget.

o In the "Smoothing" widget, apply "Binning" for continuous variables.

o Connect the "Smoothing" widget to a "Data Table" widget to inspect the smoothed data.

**5. Detect and Remove Outliers:**

o Connect the "Smoothing" widget to the "Outliers" widget.

o In the "Outliers" widget, choose the "Z-Score" method and set a threshold (e.g., 3).

o Connect the "Outliers" widget to a "Data Table" widget to inspect the data with outliers removed.
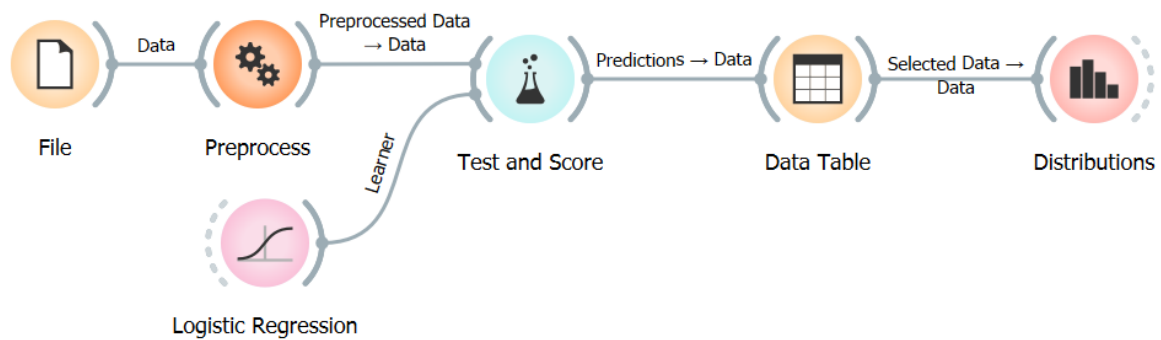
**OUTPUT:**



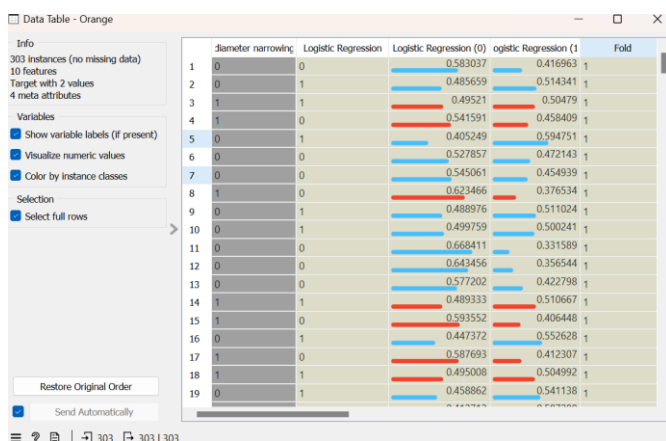FIG 2.1: IMPLEMENTATION OF DATA PREPROCESSING AND LOGISTIC REGRESSION
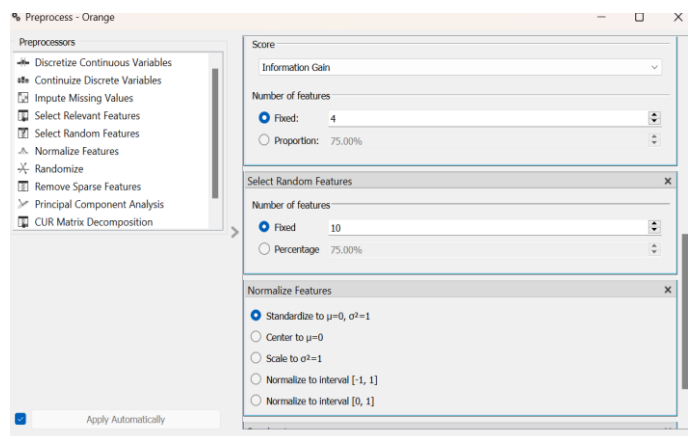


FIG 2.1.1: IMPLEMENTATION OF LOGISTIC REGRESSION    FIG 2.1.2: IMPLEMENTATION OF DATA PREPROCESSING
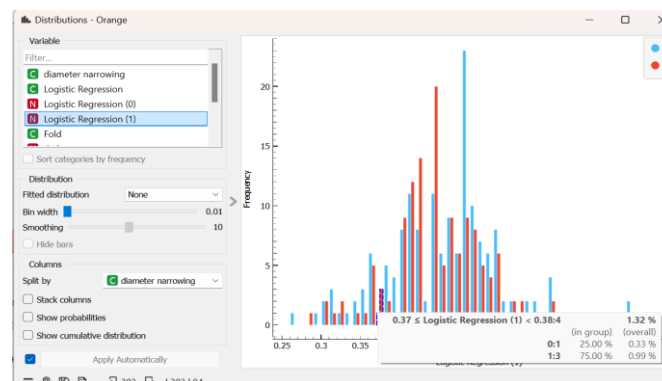


FIG 2.2: DISTRIBUTION OF LOGISTIC REGRESSION

**RESULT:**