

Assignment No 2

Aim

Perform the following operations using R/Python on the data sets:

- a) Compute and display summary statistics for each feature available in the dataset. (e.g. minimum value, maximum value, mean, range, standard deviation, variance and percentiles)
- b) Illustrate the feature distributions using histogram.
- c) Data cleaning, Data integration, Data transformation, Data model building (e.g. Classification)

Objective

The objective of this assignment is to provide practical experience in performing summary statistics, visualizing feature distributions through histograms, applying data cleaning and transformation techniques, integrating datasets, and finally, constructing a basic classification model. This enhances both analytical and modeling skills using R or Python.

Theoretical

In the data science workflow, understanding the dataset is the first essential step. This is achieved through summary statistics and visualizations. Cleaning and transforming data ensures consistency and completeness. Data integration brings multiple sources into a cohesive structure. Finally, model building—especially classification—applies machine learning to predict categorical outcomes based on feature patterns.

Methods and Explanations of Operations

a) Compute and Display Summary Statistics

- - Summary statistics provide an overview of numerical columns in the dataset.
- - Includes: minimum, maximum, mean, median, range, standard deviation, variance, and percentiles (25th, 50th, 75th).
- - Why Important: Helps understand the distribution, spread, and central tendency of the data.

b) Illustrate Feature Distributions using Histogram

- - Histograms are graphical representations of the distribution of a numerical feature.
- - They divide data into bins and count how many data points fall into each bin.
- - Why Important: Helps to detect skewness, modality, and potential outliers in the data.

c) Data Cleaning, Integration, Transformation, and Model Building

- - Data Cleaning: Handling missing values, duplicates, inconsistencies, and errors.
- - Data Integration: Merging datasets from different sources into a unified dataset.
- - Data Transformation: Standardizing formats, normalizing numerical values, and encoding categorical variables.
- - Model Building (Classification): Applying supervised learning to predict categories (e.g. logistic regression, decision tree).
- - Why Important: These processes ensure high-quality input for modeling and enable meaningful predictions.

Outputs:



Conclusion

We have successfully performed descriptive statistics, histogram-based visualizations, data cleaning, integration, and transformation followed by basic classification modeling using R/Python. This assignment has provided a comprehensive understanding of data preparation and how it feeds into effective model building. Through this assignment, I learned how each preprocessing step influences the accuracy and reliability of data analysis and predictions.