

## Practical No.2

**Title:** Create a data frame with columns at least 5 observations

- a) Retrieve a particular column from the Data Frame
- b) Summarize the data frame and observe the statistics of the Data Frame created
- c) Observe the mean and standard deviation of the data frame and print the values.

**Objective:** The aim of this practical is to understand how to create a DataFrame, retrieve specific columns, summarize the data, and compute statistical values such as mean and standard deviation using Python. These fundamental operations form the basis for data analysis, allowing for better interpretation and decision-making.

**Introduction:** A DataFrame is a widely used data structure in Python, primarily implemented through the pandas library. It provides a convenient way to store, manipulate, and analyze structured datasets with labeled rows and columns. DataFrames allow for easy access to specific columns, statistical computation, and summary analysis. In this practical, we will create a DataFrame containing multiple observations and perform various operations to extract meaningful insights.

### Procedure:

#### 1. Creating a DataFrame

- A DataFrame is composed of multiple rows (observations) and columns (attributes). Each column represents a specific feature of the dataset, such as numerical scores, categorical values, or textual data. To perform meaningful analysis, we will construct a DataFrame with at least five observations, ensuring sufficient data for statistical evaluation.
- The dataset can include various types of information, such as student names, scores in different subjects, or real-world data like temperature readings, sales figures, or employee performance metrics.

## 2. Retrieving a Specific Column

- DataFrames allow for the retrieval of specific columns, making it easier to analyze particular aspects of a dataset. By selecting a single column, we can perform operations such as sorting, filtering, and aggregating data.
- This step is essential when working with large datasets, as it helps to isolate relevant data for further processing.

## 3. Summarizing the Data

- Summarization of data involves generating key statistical measures that describe the dataset's characteristics. The summary statistics typically include:
  - **Count:** Number of non-null values in each column.
  - **Mean:** The average value of numerical columns.
  - **Standard Deviation:** A measure of data spread around the mean.
  - **Minimum and Maximum:** The lowest and highest values in each column.
  - **Quartiles (25%, 50%, 75%):** Values that help understand the distribution of data.
- These statistics help in identifying trends, detecting anomalies, and gaining insights into the data distribution.

## 4. Computing Mean and Standard Deviation

- The **mean** (average) is calculated by summing all values in a column and dividing by the number of observations. It represents the central tendency of the data.
- The **standard deviation** measures the dispersion of data points from the mean. A high standard deviation indicates significant variation in the dataset, while a low standard deviation suggests that the data points are close to the mean.
- These metrics are crucial for assessing data variability and making predictions based on past trends.

## Observations:

- The summary statistics provide an overview of the dataset's distribution and variability.
- The mean is useful for understanding the typical value in a dataset.
- The standard deviation helps in assessing the spread of data and identifying potential outliers.
- Extracting specific columns simplifies targeted analysis and improves data organization.
- Summarizing data aids in decision-making by revealing important patterns and insights.

**Conclusion:** These practical highlights the importance of working with DataFrames in Python for data analysis. By creating, retrieving, summarizing, and computing statistical measures, we can gain valuable insights into dataset characteristics. Understanding mean and standard deviation is particularly useful in various fields, including business analytics, scientific research, and machine learning. Mastering these operations enables efficient data-driven decision-making and enhances analytical capabilities.