

## Lab – 24

**Student Name: Navneet P**

**Student ID: AF0411619**

### **Topic-Pandas Library**

- **What is Pandas?**

Pandas is an open-source data manipulation and analysis library for the Python programming language that is made mainly for working with relational or labeled data both easily and intuitively. It provides various data structures and operations for manipulating numerical data and time series and for working with structured data, making it one of the most popular and widely used libraries for data analysis and manipulation in the Python ecosystem.

- **Installation**

- Pip install pandas

- **Data Cleaning**

- Pandas includes functions for handling missing data, removing duplicates, and transforming data to make it suitable for analysis.

- **Dataframes**

- The DataFrame is one of the central data structures in Pandas. It is a two-dimensional table with rows and columns, similar to a spreadsheet or a SQL table.

- **Series**

- A Series is a one-dimensional array-like object in Pandas. It can be thought of as a single column of data within a DataFrame, with an associated index

- **Grouping**

- Pandas allows you to group data by one or more columns and perform aggregation operations on the groups.

- **Why Pandas?**

- Fast and efficient for manipulating data. Data from different file objects can be easily loaded. Flexible reshaping, time-series functionality.

## Functions used in this assignment:-

1. **Pandas.Series** – A one-dimensional labeled array capable of holding any data type, similar to a column in a DataFrame.
2. **Pandas.DataFrame** – A two-dimensional, size-mutable, and potentially heterogeneous tabular data structure with labeled axes (rows and columns).
3. **to\_string()** – Converts a DataFrame or Series to a string representation for more readable console output.
4. **pandas.read\_excel()** – Reads an Excel file into a DataFrame, allowing you to load data from spreadsheets.
5. **head()** – Returns the first few rows of a DataFrame (default is 5 rows).
6. **tail()** – Returns the last few rows of a DataFrame (default is 5 rows).
7. **loc[]** – Accesses a group of rows and columns by labels or a boolean array in a DataFrame.

Q1. Create Pandas series [10,20,30,40,50] without index

Solution:

```
1 import pandas as pd # Importing the pandas library
2
3 # Creating a simple list of data
4 data = [10, 20, 30, 40, 50]
5
6 # Converting the list into a Pandas Series (with default indexing 0, 1, 2, ...)
7 series = pd.Series(data)
8
9 # Display the Series with default indexing
10 print("Series with indexing:")
11 print(data)
12
13 # Converting the Series to a string, but without the index
14 unindexed = series.to_string(index=False)
15
16 # Display the Series without the index
17 print("\nSeries without indexing:")
18 print(unindexed)
19
```

Output:

```
Series with indexing:
[10, 20, 30, 40, 50]

Series without indexing:
10
20
30
40
50
```

## Q2. Create Pandas DataFrame example.

Solution:

```
1 import pandas as pd # Importing the pandas library
2
3 # Creating a dictionary with two keys: 'Batsman' and 'Averages'
4 cricket = {
5     "Batsman": ["Rohit Sharma", "Suryakumar Yadav", "Virat Kohli"], # List of batsmen
6     "Averages": [46, 89, 56] # Corresponding batting averages
7 }
8
9 # Converting the dictionary into a Pandas DataFrame
10 df = pd.DataFrame(cricket)
11
12 # Converting the DataFrame to a string representation without the index
13 result = df.to_string(index=False)
14
15 # Printing the DataFrame in string format, without index
16 print(result)
17
18 # Printing the type of the 'df' variable, which is a Pandas DataFrame
19 print(type(df))
20
```

Output:

```
      Batsman  Averages
Rohit Sharma      46
Suryakumar Yadav  89
Virat Kohli       56
<class 'pandas.core.frame.DataFrame'>
```

Q3. Example: Monthly Sales Data Imagine you are a sales manager for a retail company, and you want to analyze the monthly sales performance of a particular product in a given year. You have recorded the monthly sales figures for that product, and you want to represent this data using a Pandas Series.

Solution:

```
1 import pandas as pd # Importing the pandas library
2
3 # List of months
4 months = ['January', 'February', 'March', 'April', 'May', 'June', 'July', 'August',
5           'September', 'October', 'November', 'December']
6
7 # Corresponding sales data for each month
8 sales_data = [12000, 13500, 14200, 12800, 14000, 15500, 16200, 15800, 16500,
9               17800, 18500, 17200]
10
11 # Combining the two lists into a dictionary where each month is mapped to sales
12 data = dict(zip(months, sales_data))
13
14 # Creating a DataFrame from the dictionary, where each dictionary item (month-sales
15 # pair) becomes a row
16 df = pd.DataFrame(data.items(), columns=["Month", "Sales"])
17
18 # Converting the DataFrame to a string, excluding the index, for a cleaner output
19 result = df.to_string(index=False)
20
21 # Printing the DataFrame in string format
22 print(result)
23
24 # Printing the type of 'df' to confirm it's a DataFrame
25 print(type(df))
```

Output:

```
Month Sales
January 12000
February 13500
March 14200
April 12800
May 14000
June 15500
July 16200
August 15800
September 16500
October 17800
November 18500
December 17200
<class 'pandas.core.frame.DataFrame'>
```

Example: Read an excel file and demonstrate the use of loc slicing, head and tail methods.

Solution:

```
1 import pandas as pd # Importing pandas for data handling
2 import numpy as np  # Importing numpy for numerical operations
3
4 # Reading data from an Excel file into a Pandas DataFrame
5 data = pd.read_excel("Excelfile.xlsx")
6
7 # Converting the 'Marks' column into a NumPy array for numerical operations
8 marks = np.array(data['Marks'])
9
10 # Calculating the average of the 'Marks' column using NumPy's average function
11 average = np.average(marks)
12
13 # Extracting the first 2 rows (top of the DataFrame)
14 top = data.head(2)
15
16 # Extracting the last 2 rows (bottom of the DataFrame)
17 bottom = data.tail(2)
18
19 # Extracting the middle rows based on specific row indices (rows from index 2 to 7)
20 middle = data.loc[2:7]
21
22 # Setting 'Student Name' as the index for the 'top' DataFrame
23 top.set_index("Student Name", inplace=True)
24
25 # Setting 'Student Name' as the index for the 'bottom' DataFrame
26 bottom.set_index("Student Name", inplace=True)
27
28 # Setting 'Student Name' as the index for the 'middle' DataFrame
29 middle.set_index("Student Name", inplace=True)
30
31 # Printing the top 2 rows
32 print(top, "\n")
33
34 # Printing the bottom 2 rows
35 print(bottom, "\n")
36
37 # Printing the middle rows (rows from index 2 to 7)
38 print(middle)
39
40 # Printing the average of the 'Marks' column
41 print("\nAverage Marks: ", average)
42
```

Output:

Marks	
Student Name	
Roy	70
Jason	89
Marks	
Student Name	
Jake	88
Hilary	80
Marks	
Student Name	
Smith	91
Kyle	91
Robin	96
Spike	71
Matt	67
Robinson	69
Average Marks: 81.2	

Student Name	Marks
Roy	70
Jason	89
Smith	91
Kyle	91
Robin	96
Spike	71
Matt	67
Robinson	69
Jake	88
Hilary	80