

## Unit - 2

### 1) Data Aggregation :

\* Aggregation is a process of implementing any mathematical operation on a dataset.

\* Aggregation is a one of technique used in pandas to manipulate the data in dataframe.

\* `Dataframe.agggregation()` function is used to apply aggregation.

Sum : It returns the sum of the values

Min : Returns the minimum of values.

Max : Returns maximum of values.

or Since aggregation only works with numeric type columns.

\* Let us consider a numeric column dataset.

	length	width	height	price
0		0.89027	48.8	13450
1	0.8118			
2	0.8114	0.82709	48.8	16500
3	0.8124	0.909722	52.4	16500
4	0.824	0.901944	54.3	13500
5				
6				
7				
8				
9				

Now let's apply aggregation functions to the table by

function to the fable by

## Applying aggregation

dataset.agg("mean", axis="rows").

## The output.

length 0.837102

width 0.915126.

height: 53.7667.

Price: 1320/-

## 2) Data manipulation using Pandas:

\* Pandas is a powerful data manipulation library for Python.

- \* It provides data structures like Dataframes and series.
- \* It also contains variety functions to manipulate and analyze data.

\* Here is a brief overview of some common data manipulation tasks using Pandas.

1. Reading and writing Data

2. Exploring Data

3. Selecting Data

4. Cleaning Data

5. Manipulating columns

6. Grouping and Aggregating

7. Merging Dataframes

8. Time Series Operations.

## \* Reading and writing Data:

```
import pandas as pd.  
df = pd.read_csv('your data.csv')  
df.to_csv('output.csv', index=False)
```

## \* Exploring Data:

df.head() // displays first five rows of Dataframe.

df.describe() // displays basic statistics of Dataframe

df.info() // display information about dataframe

## \* Selecting Data:

Column = df['column-name']

// Selecting a single column

columns = df[[column1, column2]]

// Selecting multiple columns.

Cleaning Data:

`df.dropna()` // drops rows with missing values.

`df.fillna(value)` // fill missing values with a specific value.

`df.drop_duplicates()` // removes duplicate values.

`df.replace(old_value, new_value)`

Manipulating Columns

# Adding new columns.

`df['new_column'] = df['column1'] + df['column2']`

# drop columns.

`df.drop(['column1', 'column2'])`

Grouping and aggregation

`df.groupby('columnname')`

// aggregation .agg(['sum', 'mean', 'count']).

## Merging Dataframes

result = pd.concat([df1, df2])

// concatenation

// merging

pd.merge(df1, df2).

## 3. pandas Object

\* In pandas, the primary data structures are DataFrame and Series.

\* which are built on top of numpy.

These structures are designed to handle data manipulations

### Dataframe:

\* A two dimensional labeled data structure with columns that can be of different types (numeric, string, boolean, etc...)

```
import pandas as pd  
data = {'column1': [1, 2, 3], 'column2': ['A', 'B', 'C']}  
df = pd.DataFrame(data)
```

## Attributes of DataFrame:

df.shape: Returns the number of rows and columns.

df.columns: Returns column labels

df.index: Returns row labels.

## Series:

\* A one-dimensional labeled array capable of holding any data type.

```
import pandas as pd
```

data = [1, 2, 3, 4]

series = pd.Series(data)

attribute:

- Series. values : Returns the data as a Numpy array.

- Series. index : Returns the index labels

#### 4) Indexing :

\* In pandas, data indexing refers to the process of selecting and retrieving specific subsets of data from dataframes or series.

\* Indexing in pandas can be done using flexibility for accessing rows, columns, or individual elements.

#### \* Selection by Label:

- loc [ ]

- uses 'labels' to select data.

- Allows row selecting and column by label

## of Selection by position:

### \* `iloc[]`:

\* uses integer position to select data.

\* Allows row selecting and column selecting.

### \* conditional indexing:

\* USES conditions to select data.

`df['Age'] > 30]`

### \* Multiindexing

\* It is also known as hierarchical indexing.

\* It allows you to have multiple index levels on an axis.

## 5) Combining datasets

\* Merge

\* concat

\* append

\* join

## 6. Vectorized String Operations:

### 1) Accessing Elements

```
import pandas as pd
```

```
str = "Sherbin"
```

```
printf(str[0])
```

Output

```
s.
```

### 2) Changing Case

```
import pandas as pd
```

```
str = "Sherbin"
```

```
str2 = "SHERBIN"
```

```
str2.lower()
```

```
str.upper()
```

Output: Sherbin, SHERBIN

### 3. String length:

```
import pandas as pd
```

```
str = "Shekhar"
```

```
str.length()
```

Output

7

### 4. String Concatenation:

```
import pandas as pd
```

```
str1 = "Hello"
```

```
str2 = "world"
```

```
str3 = str1 + str2
```

Output

Hello world.

### 5. Substring Matching:

```
str = "working"
```

```
str.contains("or")
```

Output:

True.

## \* String replacement:

```
str = "Sherbin";
```

```
str.replace("b", "l");
```

Output

Sherlin.