# 😊 Series Data Type

**Pandas in Python** is a powerful and widely-used open-source data manipulation and analysis library. It is built on top of other fundamental Python packages, such as **NumPy**, and provides two main data structures: **Series** and **DataFrame**, which are designed to handle and manipulate structured data efficiently.

## Key Features of Pandas:

1. **DataFrames**: This is a two-dimensional labeled data structure, similar to a table in a database, Excel spreadsheet, or a data frame in R. It consists of rows and columns, where each column can hold data of different types (e.g., integers, floats, strings).

2. **Series**: A one-dimensional array-like object that holds data of any type (integers, floats, strings, etc.). You can think of a **Series** as a single column in a DataFrame.

3. **Data Handling**: Pandas makes it easy to read and write data from various file formats such as CSV, Excel, SQL databases, JSON, and more.

4. **Data Cleaning**: Pandas includes tools for handling missing data, transforming data types, removing duplicates, and applying transformations to clean and prepare data for analysis.

5. **Indexing**: Pandas allows for flexible and efficient indexing and selection of data. You can use labels or positional indexing for easy data retrieval.

6. **Data Aggregation**: With built-in aggregation and grouping capabilities, Pandas allows users to perform complex operations, such as summing, averaging, or counting values across groups of data.

7. **Time Series Support**: Pandas is designed to handle time series data effectively, including functionalities for resampling, shifting, and calculating rolling windows.

8. **Vectorized Operations**: Pandas supports fast, element-wise operations over large datasets, making it highly efficient for large-scale data processing.

**Benefits of Using Pandas:**

1. **Ease of Use**: Pandas simplifies data manipulation and cleaning through an intuitive, easy-to-use syntax, which makes working with data more accessible to both beginners and experts.

2. **High Performance**: It provides optimized performance for handling large datasets, often reducing the time complexity of tasks like filtering, sorting, and transforming data.

3. **Integration**: Pandas can be easily integrated with other Python libraries (e.g., Matplotlib, Seaborn for plotting; Scikit-learn for machine learning) to create a powerful ecosystem for data analysis.

4. **Rich Functionality**: Pandas provides a wide array of functions for data manipulation, such as handling missing data, reshaping, merging datasets, and performing time series analysis.

5. **Data Cleaning**: Built-in functions allow for easy handling of data irregularities such as missing values, duplicated entries, and inconsistent data formats.

## Use Cases of Pandas:

1. **Data Cleaning and Preprocessing**:

   o **Missing data**: Filling, dropping, or interpolating missing values.

   o **Data transformations**: Changing data types or applying functions to columns.

   o **Removing duplicates**: Identifying and removing duplicate rows in a dataset.

2. **Exploratory Data Analysis (EDA)**:

   o Summarizing data using functions like .describe(), .mean(), .median(), .std().

   o Visualizing the distribution of data using histograms, box plots, etc.

   o Aggregating data to explore patterns or trends.

3. **Data Manipulation**:

   o Filtering and selecting data based on specific conditions.

   o Merging or joining datasets from multiple sources (e.g., SQL databases, CSV files).

   o Pivoting, reshaping, or stacking data for better analysis or presentation.

4. **Time Series Analysis**:

   o Handling time-based data, such as stock market data or weather data.

   o Resampling data to different time intervals.

   o Analyzing trends and patterns over time.

5. **Data Analysis for Machine Learning**:

   o Feature engineering: Creating new features or transforming existing ones to improve machine learning models.

   o Preparing data for machine learning models by splitting, scaling, or normalizing the dataset.

**Conclusion:**

Pandas is an essential tool for anyone working with data in Python. Its flexibility and high-level functions make it a key player in tasks like data cleaning, exploratory analysis, and preparing data for machine learning. Whether you are a beginner or an experienced data scientist, mastering Pandas will help streamline your data analysis workflow and improve efficiency.

# 😄 To begin with the Lab

1. The first main data type we will learn about for pandas is the Series data type.
2. A Series is very similar to a NumPy array (in fact it is built on top of the NumPy array object).
3. What differentiates the NumPy array from a Series, is that a Series can have axis labels, meaning it can be indexed by a label, instead of just a number location.
4. It also doesn't need to hold numeric data, it can hold any arbitrary Python Object.
5. We start by importing NumPy and Pandas as np and pd. Using the help function, we can read about the series data type in the Jupyter Notebook.

```
[12]:  import numpy as np
       import pandas as pd
```

```
[13]:  help(pd.Series)
```

```
Help on class Series in module pandas.core.series:

class Series(pandas.core.base.IndexOpsMixin, pandas.core.generic.NDFrame)
 |  One-dimensional ndarray with axis labels (including time series).
 |
 |  Labels need not be unique but must be a hashable type. The object
 |  supports both integer- and label-based indexing and provides a host of
 |  methods for performing operations involving the index. Statistical
 |  methods from ndarray have been overridden to automatically exclude
 |  missing data (currently represented as NaN).
 |
 |  Operations between Series (+, -, /, *, **) align values based on their
 |  associated index values-- they need not be the same length. The result
 |  index will be the sorted union of the two indexes.
 |
 |  Parameters
 |  ----------
 |  data : array-like, Iterable, dict, or scalar value
```

6. In this code, a list of country names (`myindex`) and corresponding years (`mydata`) are used to create a pandas Series.
7. The `pd.Series(data=mydata, index=myindex)` explicitly assigns labels to the data.
8. Then, a random list of integers is generated to simulate ages, and another pandas Series is created using names as the index.
9. This demonstrates how Series in pandas can map labeled indices to numerical data, allowing for labeled data analysis.

```
[14]: myindex = ['USA','Canada','Mexico']
      •••
```

```
[15]: mydata = [1776,1867,1821]
      •••
```

```
[16]: myser = pd.Series(data=mydata)
```

```
[17]: myser
```

```
[17]: 0    1776
      1    1867
      2    1821
      dtype: int64
```

```
[18]: pd.Series(data=mydata,index=myindex)
```

```
[18]: USA       1776
      Canada    1867
      Mexico    1821
      dtype: int64
```

```
[23]: ran_data = np.random.randint(0,100,4)
```

```
[24]: ran_data
```

```
[24]: array([39, 35, 37, 23])
```

```
[26]: names = ['Andrew','Bobo','Claire','David']
      •••
```

```
[27]: ages = pd.Series(ran_data,names)
```

```
[28]: ages
```

```
[28]: Andrew    39
      Bobo      35
      Claire    37
      David     23
      dtype: int32
```

10. In this code, a dictionary named ages is created, mapping names to their corresponding ages.
11. When this dictionary is passed to pd.Series(), it creates a pandas Series where the dictionary keys become the index (labels) and the values become the data.
12. This is a simple way to convert structured dictionary data into a pandas Series for easier analysis and manipulation.

```
[29]:  ages = {'Sammy':5,'Frank':10,'Spike':7}

[30]:  ages

[30]:  {'Frank': 10, 'Sammy': 5, 'Spike': 7}

[31]:  pd.Series(ages)

[31]:  Sammy      5
       Frank     10
       Spike      7
       dtype: int64
```

13. In this code, dictionaries representing first and second quarter sales data for various countries are converted into pandas Series.
14. When sales_Q1 is accessed using 'Japan', it returns the sales value for Japan using the label-based index.
15. Accessing sales_Q1[0] returns the first item (Japan's sales) using integer-based indexing, which is also supported by pandas Series alongside label-based access.

```
[32]:  # Imaginary Sales Data for 1st and 2nd Quarters for Global Company
       q1 = {'Japan': 80, 'China': 450, 'India': 200, 'USA': 250}
       q2 = {'Brazil': 100,'China': 500, 'India': 210,'USA': 260}

[33]:  # Convert into Pandas Series
       sales_Q1 = pd.Series(q1)
       sales_Q2 = pd.Series(q2)

[34]:  sales_Q1

[34]:  Japan      80
       China     450
       India     200
       USA       250
       dtype: int64

[35]:  # Call values based on Named Index
       sales_Q1['Japan']

[35]:  80

[36]:  # Integer Based Location information also retained!
       sales_Q1[0]

[36]:  80
```

16. This code demonstrates basic operations on a pandas Series. `sales_Q1.keys()` retrieves the index labels (country names).
17. Performing `sales_Q1 * 2` multiplies every value in the series by 2, and `sales_Q2 / 100` divides each value in the second quarter sales series by 100.
18. These operations are automatically applied element-wise due to pandas' broadcasting capabilities.

```
[40]:   # Grab just the index keys
        sales_Q1.keys()
```

```
[40]:   Index(['Japan', 'China', 'India', 'USA'], dtype='object')
```

```
[41]:   # Can Perform Operations Broadcasted across entire Series
        sales_Q1 * 2
```

```
[41]:   Japan     160
        China     900
        India     400
        USA       500
        dtype: int64
```

```
[42]:   sales_Q2 / 100
```

```
[42]:   Brazil    1.0
        China     5.0
        India     2.1
        USA       2.6
        dtype: float64
```

19. This code highlights how pandas handles arithmetic operations between two Series with different indices.
20. When adding sales_Q1 + sales_Q2, pandas aligns by index and assigns NaN to mismatched entries. Using sales_Q1.add(sales_Q2, fill_value=0) fills missing values with 0 before performing the addition, preventing NaN and producing a complete result.

```
[43]:   # Notice how Pandas informs you of mismatch with NaN
        sales_Q1 + sales_Q2
```

```
[43]:   Brazil      NaN
        China     950.0
        India     410.0
        Japan       NaN
        USA       510.0
        dtype: float64
```

```
[44]:   # You can fill these with any value you want
        sales_Q1.add(sales_Q2,fill_value=0)
```

```
[44]:   Brazil    100.0
        China     950.0
        India     410.0
        Japan      80.0
        USA       510.0
        dtype: float64
```