

1. Differentiate between a Pandas Series and DataFrame with examples.

◆ Pandas Series vs. DataFrame – Key Difference:

Feature	Pandas Series	Pandas DataFrame
Structure	1D (One-dimensional)	2D (Two-dimensional – rows and columns)
Data Type	Homogeneous (same data type)	Heterogeneous (can have different data types)
Index	Single index	Row and column indices
Like in Python	Similar to a list or 1D array	Similar to a table or Excel sheet
Usage	Store single column or labeled vector	Store complete datasets with multiple columns

◆ Pandas Series – Example:

```
import pandas as pd
```

```
# Creating a Series from a list
```

```
s = pd.Series([10, 20, 30, 40])
```

```
print(s)
```

✓ Output:

```
0    10
```

```
1    20
```

```
2    30
```

```
3    40
```

```
dtype: int64
```

📖 Explanation:

- Index on the left (0, 1, 2, 3)
- Values on the right (10, 20, 30, 40)

◆ Pandas DataFrame – Example:

```
import pandas as pd
```

Creating a DataFrame using a dictionary

```
data = {  
    'Name': ['Alice', 'Bob', 'Charlie'],  
    'Marks': [85, 90, 78]  
}  
  
df = pd.DataFrame(data)  
  
print(df)
```

✓ Output:

```
   Name  Marks  
0  Alice    85  
1   Bob    90  
2 Charlie    78
```

📖 Explanation:

- Two columns: Name and Marks
- Resembles a table or Excel sheet

🔹 Visual Analogy:

- **Series** → Like one column from an Excel file.
- **DataFrame** → The entire Excel file with multiple rows and columns.

🔹 When to Use Each:

Use Case	Use
Single column or vector	Series
Dataset with multiple fields	DataFrame

🔹 Conclusion:

✓ A **Series** is a **1D labeled array**, while a **DataFrame** is a **2D labeled data structure**. Both are core components of the **Pandas library**, widely used in **data science** for data analysis and manipulation.

2. Apply the `reindex()` method to a Pandas Series and handle missing values using `fillna()`

◆ Step 1: What is `reindex()` in Pandas?

- `reindex()` is used to **change the index** of a Series (or DataFrame).
 - If the new index has labels that **don't exist** in the original Series, Pandas inserts **NaN** (Not a Number → Missing Value).
-

◆ Step 2: What is `fillna()`?

- `fillna()` is used to **replace missing values (NaN)** with a specific value.
-

◆ Step 3: Full Example with Explanation

```
import pandas as pd
```

```
# Create original Series
```

```
s = pd.Series([100, 200, 300], index=[0, 1, 2])
```

```
print("Original Series:")
```

```
print(s)
```

```
# Reindex the Series
```

```
new_index = [0, 1, 2, 3, 4]
```

```
s_reindexed = s.reindex(new_index)
```

```
print("\nAfter Reindexing:")
```

```
print(s_reindexed)
```

```
# Fill missing values using fillna()
```

```
s_filled = s_reindexed.fillna(0)
```

```
print("\nAfter Filling NaN with 0:")
```

```
print(s_filled)
```

✓ Output:

```
Original Series:
```

```
0    100
```

```
1    200
```

```
2 300
dtype: int64
```

After Reindexing:

```
0 100.0
1 200.0
2 300.0
3  NaN
4  NaN
dtype: float64
```

After Filling NaN with 0:

```
0 100.0
1 200.0
2 300.0
3  0.0
4  0.0
dtype: float64
```

◆ Explanation:

Step	Result
------	--------

✓	Original Series Index: [0, 1, 2], Values: [100, 200, 300]
---	---

✓	After reindex() New Index: [0, 1, 2, 3, 4], Missing values = NaN
---	--

✓	After fillna(0) Missing values replaced with 0
---	--

◆ Visual Summary:

Index	Original	Reindexed	Filled
-------	----------	-----------	--------

0	100	100.0	100.0
---	-----	-------	-------

1	200	200.0	200.0
---	-----	-------	-------

Index Original Reindexed Filled

2	300	300.0	300.0
3	-	NaN	0.0
4	-	NaN	0.0

◆ Conclusion:

- `reindex()` lets you reshape or expand a Series.
- Missing values from new indices are handled using `fillna()` for cleaner data.

✓ This is essential for **data cleaning**, especially when aligning datasets or preparing data for machine learning.

3. Analyze the Role of `value_counts()` in Data Preprocessing

◆ What is `value_counts()` in Pandas?

- `value_counts()` is a method used in Pandas to **count the frequency** of unique values in a **Series** (column of a DataFrame).
 - It returns a **Series** with:
 - Unique values as **index**
 - Their count as **values**
-

◆ Syntax:

`Series.value_counts(normalize=False, sort=True, ascending=False, dropna=True)`

◆ Why is `value_counts()` important in data preprocessing?

Purpose	Explanation
✓ Data exploration	Helps understand how frequently values appear in a column.
✓ Detecting imbalance	In classification tasks (e.g., yes/no, spam/ham), helps check for class imbalance .
✓ Missing value analysis	Helps identify and count NaN values if <code>dropna=False</code> is set.
✓ Data cleaning	Helps find typos, duplicates, or rare values that need to be fixed or grouped.

Purpose	Explanation
✓ Feature engineering	Frequently used to convert categorical data into numerical counts or categories.

◆ Example:

```
import pandas as pd
```

```
data = pd.Series(['apple', 'banana', 'apple', 'mango', 'banana', 'banana', None])
```

```
# Count frequency
```

```
print(data.value_counts())
```

✓ Output:

```
banana 3
```

```
apple 2
```

```
mango 1
```

→ None is ignored by default. Use dropna=False to include missing values:

```
print(data.value_counts(dropna=False))
```

✓ Output:

```
banana 3
```

```
apple 2
```

```
mango 1
```

```
NaN 1
```

◆ Real-world Uses in Data Science:

Use Case	Example
✓ Check most popular product	df['Product'].value_counts()
✓ Count number of male/female	df['Gender'].value_counts()
✓ Detect rare classes	Useful for model balancing
✓ Identify missing values	dropna=False option

Use Case

Example

✓ Feature encoding ideas

Use counts as numerical features

◆ Conclusion:

✓ value_counts() plays a **key role in preprocessing**, helping data scientists:

- Understand distributions,
- Clean data,
- Detect imbalances,
- and prepare features for modeling.

🔍 It is one of the **first tools** used in **Exploratory Data Analysis (EDA)**.

4. Design a DataFrame and use describe() to summarize statistical properties

◆ What is describe() in Pandas?

- describe() is a built-in Pandas method used to **generate summary statistics** of **numeric columns** in a DataFrame.
 - It returns information like:
 - Count
 - Mean
 - Standard Deviation (std)
 - Minimum
 - 25%, 50%, 75% percentiles
 - Maximum
-

◆ Step-by-Step Example:

🔗 Step 1: Import Pandas and Create a DataFrame

```
import pandas as pd
```

```
# Sample student data
```

```
data = {  
    'Name': ['Alice', 'Bob', 'Charlie', 'David', 'Eva'],
```

```
'Math': [85, 78, 90, 72, 88],  
'Science': [91, 85, 89, 76, 95],  
'English': [78, 82, 85, 80, 90]  
}
```

```
df = pd.DataFrame(data)  
print("Original DataFrame:")  
print(df)
```

✓ **Output:**

	Name	Math	Science	English
0	Alice	85	91	78
1	Bob	78	85	82
2	Charlie	90	89	85
3	David	72	76	80
4	Eva	88	95	90

◆ **Step 2: Use describe()**

```
print("\nStatistical Summary:")  
print(df.describe())
```

✓ **Output of describe():**

	Math	Science	English
count	5.000000	5.000000	5.000000
mean	82.600000	87.200000	83.000000
std	6.720629	7.463632	4.690416
min	72.000000	76.000000	78.000000
25%	78.000000	85.000000	80.000000
50%	85.000000	89.000000	82.000000
75%	88.000000	91.000000	85.000000
max	90.000000	95.000000	90.000000

◆ Explanation of Output:

Statistic Meaning

count	Number of non-null entries
mean	Average value
std	Standard deviation (spread of data)
min	Smallest value
25%	1st quartile (25% of values are below this)
50%	Median (middle value)
75%	3rd quartile (75% of values are below this)
max	Largest value

◆ Conclusion:

✓ The describe() function is extremely useful in **data exploration**.

It provides a **quick statistical summary** that helps:

- Understand distributions,
- Detect outliers,
- Find missing data,
- and prepare the data for further analysis or machine learning.

🔍 This method is an essential part of **EDA (Exploratory Data Analysis)** in data science.

✓ **5. Use Pandas to write a DataFrame into text format and then read it back, maintaining data integrity**

◆ Overview:

- In real-world data science, we often need to **store data** and later **read it back** for analysis.
 - Pandas provides easy methods to:
 - Save data using to_csv() or to_txt()
 - Read it back using read_csv() or read_table()
 - These functions **preserve the structure and content** of the original DataFrame.
-

◆ Step-by-step Example:

📌 Step 1: Create a DataFrame

```
import pandas as pd
```

```
# Sample data
```

```
data = {  
    'Name': ['Alice', 'Bob', 'Charlie'],  
    'Age': [24, 27, 22],  
    'Score': [88.5, 92.0, 79.5]  
}
```

```
df = pd.DataFrame(data)  
print("Original DataFrame:")  
print(df)
```

✓ Output:

	Name	Age	Score
0	Alice	24	88.5
1	Bob	27	92.0
2	Charlie	22	79.5

📌 Step 2: Write DataFrame to a Text File (CSV format)

```
# Save as text file (comma-separated)
```

```
df.to_csv("student_data.txt", index=False)
```

📖 index=False prevents saving the row numbers as an extra column.

📌 Step 3: Read the Data Back into a New DataFrame

```
# Read the text file back
```

```
df_read = pd.read_csv("student_data.txt")
```

```
print("\nDataFrame After Reading Back:")
```

```
print(df_read)
```

✓ Output after reading:

	Name	Age	Score
0	Alice	24	88.5
1	Bob	27	92.0
2	Charlie	22	79.5

🔗 **Data Integrity Maintained** → All values are preserved exactly.

◆ How is Data Integrity Maintained?

- Column names, data types (int, float, string), and values are preserved.
 - No data is lost or changed in the process.
-

◆ Conclusion:

✓ Using Pandas, we can easily:

- Write a DataFrame to text (to_csv())
- Read it back using read_csv()
- And ensure **data integrity** is preserved

This is essential for **data storage**, **data sharing**, and **automated workflows** in data science projects.

6. Explain how missing data is handled in Pandas using built-in functions

◆ What is Missing Data?

- **Missing data** refers to **empty or null entries** in a DataFrame or Series.
 - In Pandas, missing values are represented as:
 - NaN (Not a Number) for numeric data
 - None for object/string data
-

◆ Why handle missing data?

Handling missing data is essential because:

- It can affect **data accuracy**

- Many algorithms **don't work** with missing values
 - It helps in **cleaning and preparing** the dataset for analysis or machine learning
-

◆ Built-in Pandas Functions to Handle Missing Data

Function	Description
isnull()	Detects missing values (returns True for NaN)
notnull()	Detects non-missing values
dropna()	Removes rows/columns with missing values
fillna()	Replaces missing values with a specified value or method
interpolate()	Fills missing values using interpolation

◆ Step-by-Step Example:

```
import pandas as pd

# Create DataFrame with missing values
data = {
    'Name': ['Alice', 'Bob', 'Charlie', None],
    'Age': [24, None, 22, 25],
    'Score': [88, 92, None, 85]
}

df = pd.DataFrame(data)
print("Original DataFrame:")
print(df)
```

✓ Output:

```
   Name  Age  Score
0  Alice  24.0   88.0
1   Bob   NaN   92.0
2 Charlie  22.0   NaN
3  None   25.0   85.0
```

◆ 1. Detect Missing Values

```
print(df.isnull())
```

- ✓ Returns True for each missing cell.

◆ 2. Drop Missing Values

```
df_drop = df.dropna()
```

```
print("\nAfter dropna():")
```

```
print(df_drop)
```

- ✓ Removes any row with even one missing value.

◆ 3. Fill Missing Values

```
df_filled = df.fillna({'Name': 'Unknown', 'Age': df['Age'].mean(), 'Score': 0})
```

```
print("\nAfter fillna():")
```

```
print(df_filled)
```

- ✓ Fills:
 - Missing names with "Unknown"
 - Missing age with **average age**
 - Missing score with 0

◆ 4. Interpolate Missing Values

```
df_interpolated = df.interpolate()
```

```
print("\nAfter interpolate():")
```

```
print(df_interpolated)
```

- ✓ Fills numeric gaps using **linear interpolation**

◆ Conclusion:

- ✓ Pandas provides powerful tools to **detect**, **drop**, **fill**, and **interpolate** missing values.

Clean data = Better models + More accurate analysis