

# Department of Engineering Sciences and Technology, Second Year Btech in Computer Science Project Based Learning-Python Assignment - 17

Name - Paritosh kolwadkar

SRN – 31231313

Roll no – 39

Batch – D2

Problem statement : **Write a program to handle missing values in a DataFrame using fillna() or dropna() methods. Perform aggregation operations (sum(), mean(), etc.) on columns of the DataFrame.**

Pre-requisites: Install the Pandas library:

```
pip install pandas
```

Knowledge of handling missing data in a DataFrame and performing aggregation operations.

Code:

```
# Import Pandas
import pandas as pd
import numpy as np

# Sample data with missing values
data = {
    "Name": ["Alice", "Bob", "Charlie", "David", "Eve"],
    "Age": [25, np.nan, 35, 40, np.nan],
    "Salary": [50000, 60000, np.nan, 70000, 65000],
```

```
    "Bonus": [5000, np.nan, 5500, 7000, 5200]
}

# Create a DataFrame
df = pd.DataFrame(data)

# Display the original DataFrame
print("Original DataFrame with Missing Values:")
print(df)

# Handle missing values:
# Option 1: Fill missing values with a specific value (e.g., 0 or mean)
df_filled = df.fillna({"Age": df["Age"].mean(), "Salary": df["Salary"].mean(),
"Bonus": 0})

# Option 2: Drop rows with any missing values
df_dropped = df.dropna()

# Perform aggregation operations on columns
age_mean = df_filled["Age"].mean()
salary_sum = df_filled["Salary"].sum()
bonus_mean = df_filled["Bonus"].mean()

# Display the results
print("\nDataFrame after Filling Missing Values:")
print(df_filled)

print("\nDataFrame after Dropping Rows with Missing Values:")
print(df_dropped)

print("\nAggregation Results:")
```

```
print(f"Mean Age: {age_mean}")

print(f"Total Salary: {salary_sum}")

print(f"Mean Bonus: {bonus_mean}")
```

## Explanation :

### Create a DataFrame:

- The **data** dictionary contains columns for **Name**, **Age**, **Salary**, and **Bonus**, with some missing values (represented as **np.nan**).
- The dictionary is converted into a Pandas DataFrame using **pd.DataFrame()**.

### Handling Missing Values:

- **Option 1: Fill Missing Values:**
  - The **fillna()** method is used to replace missing values:
    - The **Age** and **Salary** columns are filled with their respective mean values (calculated using **.mean()**).
    - The **Bonus** column is filled with **0** for missing values.
- **Option 2: Drop Rows with Missing Values:**
  - The **dropna()** method is used to remove rows that contain any missing values.

### Aggregation Operations:

- **mean()** is used to calculate the mean of the **Age** column after filling missing values.
- **sum()** is used to calculate the total sum of the **Salary** column after filling missing values.
- **mean()** is used again to calculate the mean of the **Bonus** column after filling missing values.

### Display Results:

- The original DataFrame with missing values is displayed.
- The DataFrame after filling missing values and the one after dropping rows are displayed.
- The aggregation results (mean of age, sum of salary, mean of bonus) are printed.

### Output:

Original DataFrame with Missing Values:

	Name	Age	Salary	Bonus
0	Alice	25.0	50000.0	5000.0
1	Bob	NaN	60000.0	NaN
2	Charlie	35.0	NaN	5500.0
3	David	40.0	70000.0	7000.0
4	Eve	NaN	65000.0	5200.0

DataFrame after Filling Missing Values:

	Name	Age	Salary	Bonus
0	Alice	25.0	50000.0	5000.0
1	Bob	31.25	60000.0	0.0
2	Charlie	35.0	61250.0	5500.0
3	David	40.0	70000.0	7000.0
4	Eve	31.25	65000.0	5200.0

DataFrame after Dropping Rows with Missing Values:

	Name	Age	Salary	Bonus
0	Alice	25.0	50000.0	5000.0
3	David	40.0	70000.0	7000.0

Aggregation Results:

Mean Age: 31.25

Total Salary: 316250.0

Mean Bonus: 4450.0

## Output Explained:

- **Handling Missing Data:**
  - `fillna()` replaces missing values with specified values (like the mean or zero).
  - `dropna()` removes rows with any missing values.
- **Aggregation Methods:**
  - `mean()` calculates the average of a column.
  - `sum()` calculates the sum of a column.
- **Use Cases:**
  - Filling missing values is useful when retaining all rows is important, while dropping rows is effective when data loss is acceptable.
  - Aggregation functions are used to summarize data and gain insights.

This program demonstrates how to handle missing data and perform basic aggregation tasks, both common operations in data preprocessing.