# Assignment 1: WIDS
# Answers to Questions 1–5

## Q1. Differentiate between Supervised and Unsupervised Learning

| Supervised Learning | Unsupervised Learning |
|---|---|
| Uses labeled data | Uses unlabeled data |
| Input–output pairs are available | Only input data is available |
| Used for prediction tasks | Used for pattern discovery |
| Examples: Regression, Classification | Examples: Clustering, Dimensionality Reduction |

## Q2. NumPy Array vs Pandas Series and DataFrame

A NumPy array is a homogeneous, multi-dimensional data structure used for efficient numerical computation.

| Feature | NumPy Array | Pandas Series | Pandas DataFrame |
|---|---|---|---|
| Data type | Homogeneous | Can be mixed | Can be mixed |
| Labels | No | Yes (index) | Yes (rows and columns) |
| Dimension | N-dimensional | 1D | 2D |
| Use case | Numerical computing | Labeled 1D data | Tabular data |

A Pandas DataFrame is conceptually similar to a dictionary of Pandas Series.

## Q3. Significance of Common Data Visualization Plots

- **Box Plot**: Displays data distribution using median, quartiles, and outliers.

- **Violin Plot**: Shows both summary statistics and the probability density of the data.

- **Histogram**: Represents frequency distribution of a numerical variable.

- **Scatter Plot**: Visualizes the relationship or correlation between two numerical variables.

# Q4. Output of the Given Code Snippets

## (a)

Reasoning followed

$$f(3) = \text{``B''} + \text{``A''} = \text{``BA''}$$
$$f(4) = \text{``BA''} + \text{``B''} = \text{``BAB''}$$
$$f(5) = \text{``BAB''} + \text{``BA''} = \text{``BABBA''}$$

**Output:**

```
BABBA
```

## (b)

The function prints values before and after the recursive call.
**Output:**

```
3-2-1-1-2-3-
```

# Q5. Hypothesis Testing and Statistical Distributions

Hypothesis testing is a statistical method used to make decisions about a population using sample data.

- **t-Distribution**: Used when sample size is small and population variance is unknown.

- **F-Distribution**: Used to compare variances and in Analysis of Variance (ANOVA).

- **Chi-square ($\chi^2$) Distribution**: Used for categorical data to test independence or goodness of fit.