# EXPLORATORY DATA ANALYSIS

Exploratory Data Analysis (EDA) is the process of thoroughly examining and characterizing data in order to find its underlying **characteristics**, possible **anomalies**, and hidden **patterns** and **relationships** and test hypotheses, and check assumptions using visual and quantitative techniques.

**Purpose:**

To gain insights and understanding of the data before building predictive or statistical models.

☐ **Why EDA is Important:**

- Identifies patterns and trends in data.
- Detects outliers and anomalies.
- Validates assumptions for statistical modeling.
- Guides feature engineering and selection.
- Helps decide suitable modeling techniques.

## Goals of EDA

1. **Understand the Dataset Structure**
   - Determine the number of rows, columns, data types, and the nature of variables.
   - Gain a clear picture of what the dataset contains and its context.
2. **Identify Data Quality Issues**
   - Detect missing, inconsistent, or erroneous values.
   - Identify outliers or anomalies that can distort model performance.
3. **Discover Patterns and Relationships**
   - Find correlations between features.
   - Spot trends and dependencies in data that can guide modeling.
4. **Check Statistical Assumptions**

- o Verify assumptions such as normality, independence, and homoscedasticity (required by many ML algorithms and statistical models).
5. **Guide Feature Engineering and Selection**
   - o Help decide which variables to keep, modify, or create.
   - o Suggest transformations (e.g., log-scaling, binning).
6. **Choose Appropriate Modeling Approaches**
   - o By understanding the data distribution and patterns, you can determine which algorithms or statistical methods might perform best.
7. **Generate Hypotheses**
   - o EDA helps formulate initial hypotheses about the data that can be tested later.

---

## Why is EDA Important?

+ Informs Subsequent Analysis:

EDA provides a solid foundation for building more complex models by revealing the underlying structure and characteristics of the data.

+ Reduces Errors:

Identifying and addressing data quality issues early on can prevent errors in later stages of analysis.

+ Improves Model Performance:

Understanding data distributions and relationships can help in selecting appropriate models and feature engineering.

+ Facilitates Decision Making:

EDA helps in extracting meaningful insights from data, enabling informed decisions and actions.

**Common Techniques Used in EDA:**

- Descriptive Statistics: Calculating measures like mean, median, standard deviation, etc., to summarize data.
- Data Visualization: Creating charts and plots (histograms, scatter plots, box plots, etc.) to visually explore data and identify patterns.
- Data Cleaning: Identifying and handling missing values, outliers, and inconsistencies in the data.
- Feature Engineering: Creating new variables from existing ones to improve model performance.
- Multivariate Analysis: Exploring relationships between multiple variables simultaneously, often using techniques like correlation matrices or scatterplot matrices.
- Dimensionality Reduction: Reducing the number of variables while retaining important information.

## Steps in EDA

1. **Data Understanding and Preparation:**

   - Understanding the Data:

     Gain a comprehensive understanding of the dataset's structure, variables, and potential issues.

   - Data Cleaning:

     Address missing values (imputation or removal), handle duplicates, correct data types, and fix errors or inconsistencies.

   - Data Transformation:

     Normalize or standardize data, create new features (feature engineering), and aggregate or disaggregate data as needed.

2. **Data Exploration:**

   - Univariate Analysis:

Analyze individual variables using descriptive statistics (mean, median, standard deviation, etc.) and visualizations like histograms and box plots.

+ Bivariate Analysis:

Examine relationships between pairs of variables using scatter plots, correlation coefficients, and cross-tabulations.

+ Multivariate Analysis:

Investigate interactions between multiple variables using techniques like pair plots, correlation matrices, and potentially more complex visualizations.

+ Outlier Detection:

Identify and handle outliers that may be data errors or represent interesting cases.

+ Pattern Recognition:

Look for trends, clusters, and other patterns within the data.

3. **Data Visualization:**

+ Creating Visualizations:

Use various plots like histograms, scatter plots, box plots, heatmaps, and more to gain insights and communicate findings.

+ Interactive Visualizations:

Employ tools for interactive exploration, allowing for dynamic filtering and zooming into data points.

4. **Hypothesis Generation and Refinement:**

+ Formulating Questions: Generate questions about the data based on initial observations and patterns.
+ Testing Hypotheses: Use the data to test potential explanations and refine research questions.

# Considerations

- Iterative Process:

  EDA is an iterative process; findings from one step often lead to new questions and further analysis.

- Domain Knowledge:

  Understanding the context of the data (domain knowledge) is crucial for interpreting results and making informed decisions.

- Choosing the Right Tools:

  Select appropriate tools and techniques based on the type of data, analysis goals, and available resources.

# Key EDA Questions

- What is the data shape and quality?
- Which features are most correlated?
- Are there missing values?
- Are there outliers?
- What is the distribution of target variable (if supervised learning)?
- Are there class imbalances?

## Univariate data analysis

Univariate data analysis involves examining a **single variable** in a dataset to understand its distribution, central tendency, and variability.

**Goals:**

- Describe the basic properties of the data.
- Identify patterns (e.g., skewness, outliers).
- Visualize the frequency or proportion of data.
- Univariate analysis can be applied to both categorical and numerical data.

## Characterizing Data with Descriptive Statistics

Descriptive statistics summarize the key characteristics of a single variable.

Common Measures:

1. **Central Tendency:** Mean, Median, Mode.
2. **Spread:** Variance, Standard Deviation, Range, Interquartile Range (IQR).
3. **Shape:** Skewness (asymmetry) and Kurtosis (peakedness).
4. **Frequency:** Counts and proportions for categorical data.

   df['Age'].describe() – numerical variable

   df['Gender'].value_counts()  - categorical variable

# Univariate Distribution Plots

These plots show **how values of a single variable are distributed**.

- **Histogram:** For numeric data, shows the frequency distribution.
- **Density Plot (KDE):** Smooth probability density estimate.

- **Boxplot:** Shows median, quartiles, and outliers.
- **Violin Plot:** Combines boxplot with KDE.

## Univariate Comparison Plots

These compare **the values of a single variable across categories or groups**.

- Boxplot across groups:

sns.boxplot(x='Loan_type', y='Income', data=df) - Compares the distribution of Income across different Loan_type categories.

- Bar Plot (for numeric mean comparison across categories):

## Univariate Composition Plots

These show the composition of categories within a single variable (mainly for categorical data).

- Bar Chart:
- Pie Chart (less preferred in analytics but widely used):
- Stacked Bar (when combined with other variables):
  Shows composition within a categorical variable.

## Univariate hypothesis tests

Univariate hypothesis tests analyze **a single variable** to determine whether its distribution or properties significantly differ from expectations.

Common Univariate Tests:

- **One-Sample t-test:** Checks if the mean of a numeric variable is equal to a hypothesized value.
- **Z-test:** Similar to t-test, used when population variance is known.

- **Chi-Square Goodness of Fit Test:** Checks if a categorical variable follows a specified distribution.

## Hypothesis Testing

Hypothesis testing is a **statistical process** used to make decisions or inferences about a population parameter based on sample data.

### Steps of Hypothesis Testing:

1. State the hypotheses:
   - Null Hypothesis ($H_0$): No effect or difference (e.g., mean income = ₹50,000).
   - Alternative Hypothesis ($H_1$): There is an effect/difference (e.g., mean income ≠ ₹50,000).
2. Choose significance level ($\alpha$):
   - Usually 0.05 (5% risk of Type I error).
3. Select the test statistic:
   - Example: t-statistic for mean tests, $x^2$ for categorical data.
4. Compute the test statistic and p-value.
5. Make a decision:

   If p-value < $\alpha$, reject $H_0$ (evidence suggests $H_1$).

   Otherwise, fail to reject $H_0$.

## Error Types in Hypothesis Testing

### Type I Error ($\alpha$):

- Rejecting a true null hypothesis (false positive).
- Example: Concluding a drug works when it actually doesn't.

### Type II Error ($\beta$):

- Failing to reject a false null hypothesis (false negative).

- Example: Concluding a drug doesn't work when it actually does

## Test Statistic

A **test statistic** is a standardized value used to compare sample data with the null hypothesis.

- **Formula for t-test (one-sample):**
  - Formula for t-test (one-sample):

$$t = \frac{\bar{x} - \mu_0}{s/\sqrt{n}}$$

where:

- $\bar{x}$ = sample mean
- $\mu_0$ = hypothesized mean
- $s$ = sample standard deviation
- $n$ = sample size

- 
  - **Interpretation:** The larger the absolute value of the test statistic, the stronger the evidence against $H_0$.

## Understanding the p-value

- **Definition:** The probability of obtaining a result **at least as extreme** as the observed result, assuming $H_0$ is true.
- **Low p-value (< 0.05):** Evidence against $H_0$ (statistically significant).
- **High p-value (> 0.05):** Not enough evidence to reject $H_0$.

**Example:** If p-value = 0.02, it means that **if $H_0$ were true**, there is a 2% chance of observing the given sample data or something more extreme.

## Multivariate Analysis

Multivariate analysis involves examining **three or more variables simultaneously** to find **patterns, relationships, and interactions** between them.

**Goal:**

- To understand how multiple features relate to each other.
- To detect **redundancies, dependencies, or patterns** that might affect data modeling.
- It includes tools like **correlation matrices, covariance, PCA, and multivariate regression**.

# Finding Relationships in Data

To explore relationships between numerical variables, we often use **covariance** and **correlation**.

## Covariance

Covariance measures **how two variables vary together**.

- If both variables increase together → **Positive Covariance**.
- If one increases while the other decreases → **Negative Covariance**.
- If unrelated → **Covariance close to 0**.

**Units:** Covariance is **not standardized**, so its value depends on the scale of variables

# Correlation

Correlation measures the **strength and direction of a linear relationship** between two variables, standardized to the range **[-1, 1]**.

## Types of Correlation (By Direction)

### a) Positive Correlation

- When one variable increases, the other variable also increases.
- Example: **Height vs. Weight** – taller people often weigh more.
- Correlation coefficient $r > 0$.
- **Plot:** Points trend upward on a scatterplot.

### b) Negative Correlation

- When one variable increases, the other decreases.
- Example: **Price vs. Demand** – as price increases, demand decreases.
- Correlation coefficient $r < 0$.
- **Plot:** Points trend downward on a scatterplot.

### c) Zero Correlation

- No relationship between the two variables.
- Example: **Height vs. Favorite Color** – unrelated variables.
- Correlation coefficient $r \approx 0$.
- **Plot:** Points are randomly scattered.

---

# Types of Correlation (By Linearity)

### a) Linear Correlation

- Variables change at a constant rate.
- Example: **Study Time vs. Marks** – more study hours generally yield better scores.
- Pearson correlation is typically used.

b) Non-linear (Curvilinear) Correlation

- Variables are related but not linearly (e.g., parabolic or U-shaped relationship).
- Example: **Stress vs. Performance** – performance increases with stress up to a point, then decreases.

# Types of Correlation (By Method of Measurement)

## a) Pearson Correlation (r)

- Measures **linear relationship** between two continuous variables.
- Values range from **-1 to +1**.
- Assumes data is normally distributed.
- 

## Spearman Rank Correlation (ρ)

- **Non-parametric** measure (uses ranks instead of actual values).
- Used for ordinal data or when variables are not normally distributed.
- Detects **monotonic** relationships (not just linear).

# Strength of Correlation (Magnitude)

| Value of r | Strength |
|---|---|
| 0.9 – 1.0 | Very strong |
| 0.7 – 0.9 | Strong |
| 0.5 – 0.7 | Moderate |
| 0.3 – 0.5 | Weak |
| 0 – 0.3 | Very weak/none |

Visual Representation of Types

- **Positive:** ↗ (upward trend).
- **Negative:** ↘ (downward trend).
- **Zero:** scattered without pattern.
- **Non-linear:** U-shape or other curves.

| Type | Use Case | Examples |
| --- | --- | --- |
| Multivariate Distribution | Show joint distributions of 2+ variables | Pairplot, KDE, Jointplot |
| Multivariate Comparison | Compare multiple groups across variables | Grouped Boxplot, Violin, FacetGrid |
| Multivariate Relationship | Visualize relationships among multiple numeric features | Pairplot, Bubble plot, 3D Scatter, Heatmap |
| Multivariate Composition | Show proportions across multiple categories | Stacked Bar, Mosaic, Treemap |

1. **What is the main purpose of Exploratory Data Analysis (EDA)?**

   a) To directly build machine learning models

   b) To explore and understand data characteristics

   c) To store data securely

   d) To clean data for production use

   **Answer:** b) To explore and understand data characteristics

2. **Which of the following is NOT a goal of EDA?**

   a) Identifying data patterns

   b) Checking statistical assumptions

   c) Optimizing hyperparameters

   d) Guiding feature engineering

   **Answer:** c) Optimizing hyperparameters

3. **Which plot is most suitable for visualizing the distribution of a single numerical variable?**

a) Scatter plot

b) Histogram

c) Heatmap

d) Mosaic plot

**Answer:** b) Histogram

4. **Which correlation method in pandas is used for ordinal data?**

a) Pearson

b) Spearman

c) Kendall

d) Both b and c

**Answer:** d) Both b and c

5. **What does a correlation coefficient close to 0 imply?**

a) Strong positive relationship

b) Strong negative relationship

c) No linear relationship

d) The variables are dependent

**Answer:** c) No linear relationship

6. **Which of the following is NOT a step in EDA?**

a) Data Cleaning

b) Model Deployment

c) Outlier Detection

d) Visualization

**Answer:** b) Model Deployment

7. **Which Python function provides a quick statistical summary of a DataFrame column?**

a) df.summary()

b) df.describe()

c) df.overview()

d) df.stats()

**Answer:** b) df.describe()

8. **Which plot is best suited for detecting outliers in a dataset?**

a) Box Plot

b) Histogram

c) Bar Plot

d) Pie Chart

**Answer:** a) Box Plot

9. **Which of the following is a measure of data spread?**

a) Median

b) Mode

c) Variance

d) Mean

**Answer:** c) Variance

10.     **What is the purpose of feature engineering during EDA?**

a) To remove all columns from the dataset

b) To create new variables from existing ones

c) To shuffle rows in the dataset

d) To convert a dataset into text format

**Answer:** b) To create new variables from existing ones

---

# MCQs – Univariate & Bivariate Analysis

6. **Which plot is best for visualizing the frequency of categories in a categorical variable?**

a) Histogram

b) Pie Chart

c) Bar Chart

d) Scatter Plot

**Answer:** c) Bar Chart

7. **Which plot is commonly used to visualize the relationship between two numerical variables?**

a) Histogram

b) Scatter Plot

c) Pie Chart

d) Count Plot

**Answer:** b) Scatter Plot

8. **What does a cross-tabulation (crosstab) table show?**

a) The average values of a numeric column

b) Relationships between two categorical variables

c) Variance of data

d) Skewness of data

**Answer:** b) Relationships between two categorical variables

9. **What does a heatmap of correlation values visualize?**

a) Frequency of categorical values

b) Pairwise correlation between numerical features

c) Boxplot distribution

d) Error metrics of models

**Answer:** b) Pairwise correlation between numerical features

10. **Which of the following methods is most suitable to compare numerical values across groups?**

a) Violin Plot

b) Pie Chart

c) Line Chart

d) Treemap

**Answer:** a) Violin Plot

# MCQs – Hypothesis Testing

11.      **In hypothesis testing, the null hypothesis ($H_0$) usually states:**

a) There is a significant effect

b) There is no effect or no difference

c) The data is always skewed

d) The mean equals zero

**Answer:** b) There is no effect or no difference

12.      **What does a p-value less than 0.05 generally indicate?**

a) Accept the null hypothesis

b) Reject the null hypothesis

c) There is no evidence against $H_0$

d) The test is invalid

**Answer:** b) Reject the null hypothesis

13.      **Which error occurs when we reject a true null hypothesis?**

a) Type I Error

b) Type II Error

c) Type III Error

d) Type IV Error

**Answer:** a) Type I Error

14.      **Which test is used to check if a sample mean is equal to a known population mean?**

a) Chi-square test

b) One-sample t-test

c) ANOVA test

d) Z-test

**Answer:** b) One-sample t-test

15.      **Which hypothesis test is used for categorical data frequency distribution?**

a) Chi-square test

b) t-test

c) Z-test

d) F-test

**Answer:** a) Chi-square test

---

# MCQs – Correlation and Covariance

16. **Which of the following statements about correlation is true?**

a) Correlation values range from -10 to +10

b) A correlation of 0 means no linear relationship

c) A correlation of -1 means a strong positive relationship

d) Correlation and causation are the same

**Answer:** b) A correlation of 0 means no linear relationship

17. **Which correlation method is best for non-linear monotonic relationships?**

a) Pearson

b) Spearman

c) Kendall

d) Both b and c

**Answer:** d) Both b and c

18. **If two variables have a covariance of 0, this means:**

a) They are independent

b) They have no linear relationship

c) They are strongly correlated

d) Their mean is 0

**Answer:** b) They have no linear relationship

19. **What is the main disadvantage of covariance compared to correlation?**

a) It is more accurate

b) It is unit-dependent and not normalized

c) It only works for categorical data

d) It cannot be calculated in Python

**Answer:** b) It is unit-dependent and not normalized

20. **Which of these represents a very strong negative correlation?**

a) -0.95

b) 0.8

c) -0.3

d) 0.1

**Answer:** a) -0.95

# B. 5-Mark Questions

1. Define **Exploratory Data Analysis (EDA)** and explain why it is important in data science projects.
2. Explain the difference between **univariate, bivariate, and multivariate analysis** with examples.
3. What are **descriptive statistics** in univariate analysis? Describe measures of central tendency and spread.
4. List and explain **any four common visualization techniques used in EDA** with their use cases.
5. What is **hypothesis testing**? Explain the concept of **Type I and Type II errors** with examples.
6. Differentiate between **covariance** and **correlation**. Why is correlation preferred over covariance?
7. Explain **positive, negative, and zero correlation** with suitable examples.

# C. 10-Mark Questions

1. **Discuss the goals and benefits of EDA** in detail.
2. **Explain the different steps in EDA**, starting from data understanding to hypothesis generation, with real-life examples.
3. **Describe univariate data analysis techniques** (descriptive statistics, distribution plots, comparison plots, composition plots) with examples and Python commands.
4. **What is multivariate analysis?** Discuss the types of multivariate plots (distribution, comparison, relationship, composition) with examples.
5. **Explain hypothesis testing** in detail. Write steps of a typical hypothesis test, including the interpretation of p-values and test statistics.
6. Explain **Pearson, Spearman, and Kendall correlation methods** with use cases and examples.