

## 7. Compare between Data Analytics and Data Science

Data Analytics focuses primarily on analyzing existing data to find actionable insights. It often involves descriptive and diagnostic analytics to understand past performance and trends.

Data Science encompasses a broader scope, including data analytics, but also involves predictive modeling, machine learning, and the development of algorithms to forecast future outcomes. Data science often requires a deeper understanding of programming and statistical modeling.

14. Break down the concept of Exploratory Data Analysis (EDA) and illustrate its significance in statistics. Highlight how EDA facilitates the identification of patterns, detection of outliers, and validation of assumptions using both graphical and statistical methods. (Long Answer)

Exploratory Data Analysis (EDA) is the process of summarizing and visualizing data to understand its main characteristics before formal modeling.

Significance in Statistics.

Helps uncover patterns and relationships between variables.

Detects outliers and anomalies that may affect analysis.

Validates assumptions like normality and equal variance needed for statistical tests.

How EDA helps:

Identification of patterns: Uses scatter plots, histograms, and correlation analysis.

Detection of outliers: Uses box plots and Z-score calculations.

Validation of assumptions: Uses Q-Q plots, histograms, and tests like Shapiro-Wilk.

Overall, EDA ensures data quality and guides appropriate statistical methods for accurate results.

15. Analyze how outliers can be identified in a dataset by examining their impact on data distribution. Compare at least two different detection methods, such as the Z-score method and the Interquartile Range (IQR) method, and explain how each method determines outliers. (Long Answer)

Comparison of Two Outlier Detection Methods:

### 1. Z-Score Method:

- Calculates how many standard deviations a data point is from the mean.
- Formula:  $Z = \frac{(X - \mu)}{\sigma}$ , where  $X$  is the data point,  $\mu$  is the mean, and  $\sigma$  is the standard deviation.
- Typically, points with  $|Z| > 3$  are considered outliers.
- Works well for normally distributed data but sensitive to mean and standard deviation distortion caused by extreme values.

### 2. Interquartile Range (IQR) Method:

- Uses the middle 50% of data between the 1st quartile (Q1) and 3rd quartile (Q3).
- Calculates  $IQR = Q3 - Q1$ .
- Outliers are values below  $Q1 - 1.5 \times IQR$  or above  $Q3 + 1.5 \times IQR$ .
- More robust to skewed data and less affected by extreme values.

19. What is machine learning? Justify its role and importance in data science with a brief explanation. (Long Answer)

Machine Learning (ML):

Machine learning is a branch of artificial intelligence that enables computers to learn from data and improve their performance on tasks without being explicitly programmed.

Role and Importance in Data Science:

- Automates decision-making by building predictive models from data.
- Finds complex patterns and relationships that traditional analysis might miss.
- Enables scalable analysis on large and complex datasets.
- Supports tasks like classification, regression, clustering, and recommendation essential for extracting insights and making data-driven decisions.

## 16. Compare quantitative data and qualitative data?

### Comparison of Quantitative Data and Qualitative Data:

Aspect	Quantitative Data	Qualitative Data
Definition	Numerical data that can be measured and counted	Non-numerical data describing qualities or categories
Nature	Expressed in numbers (e.g., height, age)	Expressed in words or labels (e.g., color, gender)
Types	Discrete (countable) and continuous (measurable)	Nominal (categories without order) and ordinal (categories with order)
Analysis	Uses statistical methods (mean, median, correlation)	Uses thematic analysis, coding, or categorization
Examples	Temperature, weight, test scores	Eye color, occupation, opinions

## 17. Explain different stages of data Science?(life cycle of Data Science)

### Different Stages of Data Science:

1. **Problem Understanding:** Define the business or research problem clearly.
2. **Data Collection:** Gather relevant data from various sources like databases, APIs, or files.
3. **Data Cleaning:** Handle missing values, remove errors, and preprocess data for analysis.
4. **Exploratory Data Analysis (EDA):** Explore and visualize data to find patterns and insights.
5. **Modeling:** Build and train machine learning or statistical models on the prepared data.
6. **Evaluation:** Assess model performance using metrics like accuracy, precision, or RMSE.
7. **Deployment:** Implement the model into production for real-world use.
8. **Monitoring & Maintenance:** Continuously monitor and update the model as needed.

## 18. Define one-sample t-test. Explain when it is used in statistical analysis. (Long Answer)

### One-Sample t-Test:

A one-sample t-test is a statistical method used to determine whether the mean of a single sample differs significantly from a known or hypothesized population mean.

### When it is used:

- When the population mean is known or specified.
- When the sample size is small (typically less than 30).
- When the population standard deviation is unknown.
- To test hypotheses about the average value of a population based on sample data.

## 8. Evaluate the effectiveness of the median over the mean in specific data scenarios. Justify your reasoning with appropriate examples.

### Effectiveness of the Median Over the Mean:

The median is often more effective than the mean in datasets that are skewed or contain outliers. While the mean considers all values and can be easily affected by extreme values, the median represents the middle value and remains stable regardless of outliers.

### Justification with Examples:

#### i) Skewed Data:

In income data, if most people earn between \$30,000–\$50,000 but a few earn millions, the mean might rise to \$100,000, which misrepresents the typical income. The median (e.g., \$40,000) gives a better picture of what most people earn.

#### ii) Presence of Outliers:

For test scores: 55, 58, 60, 62, 95.

$$\text{Mean} = (55 + 58 + 60 + 62 + 95) / 5 = 66$$

$$\text{Median} = 60$$

The high score of 95 pulls the mean up, while the median remains unaffected, giving a more accurate center for the typical score.

9. Analyze the relationship between a population and a sample in inferential statistics, and examine their respective roles through examples.

#### Relationship Between Population and Sample in Inferential Statistics:

In inferential statistics, we study a sample to make conclusions or predictions about a population. A population includes all elements of interest, while a sample is a smaller, manageable subset selected from the population.

Roles:

i) **Population:** Represents the entire group we want to understand.

*Example:* All college students in a country.

ii) **Sample:** Represents the portion of the population used to collect data.

*Example:* 500 students randomly selected from various colleges.

Purpose in Inferential Statistics:

We use statistical methods (like confidence intervals, hypothesis testing) on the sample data to make inferences about the population. This approach is efficient when studying the entire population is too costly, time-consuming, or impossible.

Example:

If a researcher wants to know the average screen time of teenagers in a country:

**Population:** All teenagers in the country.

**Sample:** A group of 1,000 teenagers surveyed.

The sample mean is used to estimate the population mean.

10. Analyze the structure and components of a confusion matrix in classification tasks. Interpret its elements (TP, FP, TN, FN) with an example to assess model performance. (Long Answer)

#### Confusion Matrix in Classification Tasks:

A confusion matrix is a performance measurement tool for classification models. It shows the number of correct and incorrect predictions made by the model compared to the actual outcomes.

Structure of a Confusion Matrix:

	Predicted Positive	Predicted Negative
Actual Positive	True Positive (TP)	False Negative (FN)
Actual Negative	False Positive (FP)	True Negative (TN)

Components Explained:

- **True Positive (TP):** Model correctly predicts the positive class.
- **False Positive (FP):** Model wrongly predicts positive for a negative case (Type I error).
- **True Negative (TN):** Model correctly predicts the negative class.
- **False Negative (FN):** Model wrongly predicts negative for a positive case (Type II error).

Example – Disease Detection Model:

Suppose we test 100 patients, where 40 actually have a disease and 60 do not.

- Model predicts 35 correctly as having the disease → TP = 35
- Predicts 5 patients with the disease as healthy → FN = 5
- Predicts 10 healthy patients as having the disease → FP = 10
- Predicts 50 healthy patients correctly → TN = 50

Interpretation:

- A high TP and TN means the model is accurate.
- Low FP and FN means fewer misclassifications.
- Helps calculate metrics like accuracy, precision, recall, and F1-score for deeper performance evaluation.

12. Draw and explain the bias-variance tradeoff in machine learning. Interpret its impact on model performance. (Long Answer) Bias-Variance Tradeoff:

a) Bias refers to error due to overly simplistic assumptions in the learning algorithm.

i) High bias means the model is too simple and underfits the data, missing important patterns.

a) Variance refers to error due to the model being too sensitive to small fluctuations in the training data.

i) High variance means the model is too complex and overfits the data, capturing noise as if it were a pattern.

Tradeoff Explanation:

i) Increasing model complexity usually reduces bias but increases variance.

ii) Decreasing complexity increases bias but reduces variance.

iii) The goal is to find a balance where both bias and variance are minimized to reduce the total error.

Impact on Model Performance:

Model Type	Bias	Variance	Result
High Bias	High	Low	Underfitting
High Variance	Low	High	Overfitting
Balanced Model	Moderate	Moderate	Good generalization

13. Identify underfitting and overfitting in machine learning through model behavior and apply your understanding to interpret their impact on prediction accuracy, using diagrams where appropriate. (Long Answer)

a) Underfitting:

i) Model Behavior: The model is too simple to capture the underlying pattern of the data.

ii) Impact: Poor performance on both training and test data — low accuracy and high error.

iii) Cause: High bias, low variance.

b) Overfitting :

i) Model Behavior: The model is too complex and fits noise or random fluctuations in the training data.

ii) Impact: Excellent performance on training data but poor generalization to new, unseen data — high accuracy on training but low on test data.

iii) Cause: Low bias, high variance.

Summary:

Underfitting → model too simple → poor accuracy on train & test data.

Overfitting → model too complex → great accuracy on train but poor on test data.

Goal: Find a balance for good generalization and high prediction accuracy.

11. Analyze the impact of different imputation techniques on a given dataset with missing values by comparing their outcomes based on relevant criteria. (Long Answer)

To analyze the impact of different imputation techniques on a dataset with missing values, compare their outcomes based on the following criteria:

1. Accuracy and Model Performance:

Evaluate how each imputation method affects the predictive accuracy or other performance metrics of models trained on the imputed data. Some methods preserve data patterns better, leading to improved model results.

2. Bias and Variance:

- Simple methods like mean/median imputation can introduce bias by underestimating variance.
- More advanced methods (e.g., K-Nearest Neighbors, Multiple Imputation) better capture the data distribution, reducing bias and preserving variance.

3. Computational Complexity:

- Simple techniques (mean, median) are fast and easy to implement.
- Complex methods (e.g., multiple imputation, model-based) require more computation time and resources.

4. Suitability to Data Type and Missingness Pattern:

- Mean/median imputation suits numerical data with random missingness (MCAR).
- Model-based or KNN imputation works better for data with patterns or mixed types (numerical + categorical).

5. Preservation of Relationships:

Advanced methods maintain correlations between features better than simple imputations, leading to more realistic datasets and reliable analysis.

12. Analyze the effect of different imputation techniques (such as mean, median, and KNN imputation) on model performance when applied to a dataset containing missing values. Compare their outcomes using appropriate performance metrics (Long Answer)

When applying different imputation techniques like mean, median, and KNN imputation to a dataset with missing values, their effects on model performance can be analyzed as follows:

1. Mean Imputation

- Replaces missing values with the mean of the feature.
- Simple and fast but can reduce variance and may introduce bias if data is skewed.
- May lead to underestimation of model error and slightly lower accuracy, especially with non-symmetric distributions.

2. Median Imputation

- Uses the median value to replace missing data.
- More robust to outliers and skewed data compared to mean imputation.
- Often results in better model performance than mean imputation when data is not normally distributed.

3. KNN Imputation

- Imputes missing values based on the closest k neighbors' values.
- Preserves data relationships and feature correlations better.
- Usually leads to higher model accuracy and F1-score due to more realistic imputations but is computationally expensive.

Performance Comparison Using Metrics

Imputation Method	Accuracy	Precision	Recall	F1-Score	Computation Time
Mean	Moderate	Moderate	Moderate	Moderate	Low
Median	Moderate-High	Moderate-High	Moderate-High	Moderate-High	Low
KNN	High	High	High	High	High

1. Compare and contrast the advantages and disadvantages of ensemble methods like Bagging, Boosting, and Stacking. (for 5 marks)

Bagging

*Advantages:* Reduces variance and overfitting by training multiple models on different random subsets of data (e.g., Random Forest). It improves stability and accuracy.

*Disadvantages:* Can be less effective if the base models are biased; it mainly reduces variance, not bias.

Boosting

*Advantages:* Focuses on correcting errors from previous models by sequentially training weak learners, which reduces both bias and variance and often achieves high accuracy (e.g., AdaBoost, XGBoost).

*Disadvantages:* More prone to overfitting if not properly regularized and is computationally more expensive.

Stacking

*Advantages:* Combines different types of models (heterogeneous learners) to leverage their strengths, often leading to better predictive performance than individual models.

*Disadvantages:* More complex to implement and interpret, requires careful tuning of meta-model, and can be computationally intensive.

2. Analyze the impact of hyperparameter tuning in Random Search in deep learning models. (for 5 marks)

⇒ Hyperparameter tuning using Random Search significantly impacts the performance of deep learning models by efficiently exploring a wide range of hyperparameter values. Unlike grid search, which exhaustively tests all combinations, Random Search samples random combinations, allowing it to find better-performing hyperparameters faster and with fewer iterations.

This method increases the chance of discovering optimal or near-optimal settings for parameters like learning rate, batch size, number of layers, and dropout rates, which directly affect model accuracy and generalization. Random Search also reduces computational cost by focusing on promising areas in the hyperparameter space instead of wasting resources on less important combinations.

3. A regression analysis between apples (y) and oranges (x) resulted in the following least squares line:  $y = 100 + 2x$ . Predict the implication if oranges are increased by 1 (for 5 marks)

⇒ The regression equation is  $y = 100 + 2x$ , where y represents apples and x represents oranges. If oranges (x) increase by 1 unit, the predicted number of apples (y) will increase by 2 units, because the slope coefficient is 2. This implies a positive relationship between oranges and apples, meaning that for every additional orange, the number of apples is expected to increase by 2. The intercept (100) indicates the predicted number of apples when the number of oranges is zero. Overall, increasing oranges by 1 leads to a predicted increase of 2 apples according to the model.

9. Analyze how Logistic Regression is applied to classify binary outcomes. Examine how its approach, underlying assumptions, and output differ from those of Linear Regression. (Long Answer)

Logistic Regression is used to classify binary outcomes (e.g., 0 or 1) by modeling the probability that an input belongs to the positive class. It applies the sigmoid function to a linear combination of input features to ensure the output lies between 0 and 1:

#### Differences from Linear Regression:

1. Purpose:
  - *Logistic Regression*: Classification (binary outcomes).
  - *Linear Regression*: Prediction of continuous values.
2. Output:
  - Logistic Regression outputs probabilities, which are mapped to class labels (e.g., using a 0.5 threshold).
  - Linear Regression outputs continuous values without bounds.
3. Function Used:
  - Logistic Regression uses the sigmoid function.
  - Linear Regression uses a linear function.
4. Assumptions:
  - Logistic Regression assumes a log-odds linear relationship between features and the probability.
  - Linear Regression assumes a linear relationship between features and output.
5. Error/Cost Function:
  - Logistic Regression uses log loss (cross-entropy).
  - Linear Regression uses mean squared error (MSE).

13. Analyze how Support Vector Machines (SVM) separate data points into distinct classes by identifying optimal hyperplanes. Examine the role of kernel functions in transforming non-linearly separable data into a higher-dimensional space for effective classification.

Support Vector Machines (SVM) classify data by finding the optimal hyperplane that best separates data points of different classes. The optimal hyperplane maximizes the margin, which is the distance between the hyperplane and the nearest data points from each class (called support vectors). This helps improve generalization and robustness.

For linearly separable data, SVM finds a straight-line (or flat-plane in higher dimensions) boundary. However, when data is not linearly separable, SVM uses kernel functions to map the data into a higher-dimensional space where a linear separation becomes possible.

#### Common Kernel Functions:

- Linear kernel: No transformation; used when data is already linearly separable.
- Polynomial kernel: Captures curved boundaries.
- Radial Basis Function (RBF) kernel: Projects data into infinite dimensions to handle complex, non-linear patterns.

By using kernel functions, SVM avoids explicit transformation, applying the "kernel trick" to compute dot products efficiently in the new space. This enables effective classification of complex datasets