

EDA AND STATISTICS
END TERM SAMPLE QUESTIONS
INT351

UNIT 1

1. Describe three common data visualization techniques used in EDA, and provide examples of when each is useful.

Exploratory Data Analysis (EDA) is a crucial step in the data analysis process, and data visualization is a key component of EDA. Here are three common data visualization techniques used in EDA and examples of when each is useful:

Histograms:

Histograms are useful for visualizing the distribution of a single numeric variable. They provide insights into the central tendency, spread, and shape of the data, helping you understand its underlying characteristics. In the case of a dataset of exam scores for a group of students, a histogram can be created to see how scores are distributed. This can reveal whether the scores follow a normal distribution, are skewed, or exhibit multiple modes, which can inform further analysis or decisionmaking.

Scatter Plots:

Scatter plots are used to visualize the relationship between two numeric variables. They help identify patterns, correlations, and outliers in the data. In the case of a dataset on the advertising spend and corresponding sales revenue for a set of products. By creating a scatter plot, it can be examined whether there is a linear relationship between advertising spend and sales. This can help in assessing the effectiveness of advertising campaigns.

Box Plots:

Box plots, also known as boxandwhisker plots, are useful for summarizing and comparing the distribution of a numeric variable across different categories or groups. They provide insights into the spread and central tendency of the data within each group. In the case of a dataset on the salaries of employees in a company, a box plot can be created to compare the salary distributions by department. This can help identify variations in salary ranges between different departments and potential outliers. Box plots are especially useful for detecting differences in the central tendency and spread of data in a categorical context.

2. What are some ethical considerations when working with private data in data analysis?

Working with private data in data analysis comes with significant ethical responsibilities to protect individuals' privacy and ensure that data is used responsibly. The following are some important ethical considerations when working with private data:

Informed Consent: Ensure that data subjects have given informed and voluntary consent for their data to be collected, stored, and used for the specific purposes outlined. It's important to be transparent about the data collection process, its purpose, and any potential risks involved.

Data Anonymization and Deidentification: Prioritize the anonymization or deidentification of data to protect the privacy of individuals. This involves removing or altering personally identifiable information (PII) such as names, addresses, and social security numbers to make it difficult to reidentify individuals.

Data Minimization: Collect and retain only the data that is necessary for the intended analysis or research. Avoid collecting excessive or irrelevant information, as this can increase privacy risks.

Data Security: Implement robust data security measures to protect the data from unauthorized access, breaches, or theft. This includes encryption, access controls, and secure storage practices.

Data Ownership and Control: Respect the ownership and control of data. Ensure that data subjects have the right to access, correct, or delete their data as per applicable data protection regulations (e.g., GDPR in Europe).

Data Sharing and Sharing Agreements: If data is shared with third parties, establish clear datasharing agreements that specify how the data will be used, protected, and for what purposes. Ensure that the data is used only in accordance with these agreements.

Purpose Limitation: Use the data only for the purpose for which it was collected and not for any other unrelated purposes without explicit consent.

Data Retention and Disposal: Set clear policies for data retention and disposal. Do not retain data longer than necessary, and securely dispose of data when it is no longer needed.

3. Provide an example of a Python library commonly used for web scraping and explain its key features.

One commonly used Python library for web scraping is BeautifulSoup. BeautifulSoup is a powerful and userfriendly library for parsing HTML and XML documents, making it easier to extract data from web pages. It provides a structured representation of the document, making it easy to navigate and extract data. It is possible to locate specific elements and their attributes using functions like `find()`, `find_all()`, and navigation methods like `parent`,

`next_sibling`, and `previous_sibling`. It can handle poorly formatted or broken HTML with a degree of tolerance, making it resilient to imperfect web pages. BeautifulSoup allows you to extract text, links, images, and other data from web pages effortlessly. The following is an example of how BeautifulSoup can be used to scrape data from a web page:

```
from bs4 import BeautifulSoup
import requests
url = 'https://google.com'
response = requests.get(url)
soup = BeautifulSoup(response.text, 'html.parser')
title = soup.title.string
print(f'Title of the page: {title}')
links = soup.find_all('a')
for link in links:
    print(link.get('href'))
```

4. What are some common methods for imputing missing values in a dataset, and under what circumstances should each method be used?

Imputing missing values is a crucial step in data preprocessing when dealing with datasets that have incomplete or missing data. The choice of imputation method depends on the nature of the data and the reasons behind the missing values. Some common methods for imputing missing values are given below:

Mean/Median/Mode Imputation:

This method is suitable for imputing missing values in numerical or categorical data where missing values are missing completely at random (MCAR). In MCAR scenarios, the missingness is unrelated to any other variables in the dataset. It can be solved by replacing missing values with the mean (for numerical data), median (for numerical data with outliers), or mode (for categorical data) of the observed values in the same column.

Forward Fill and Backward Fill:

It can be used in the case of timeseries data, sequential data, or datasets with a natural order are good candidates for forwardfill and backwardfill imputation. This method is appropriate when missing values occur intermittently in sequences. It can be solved by filling missing values with the most recent nonmissing value (forward fill) or the next nonmissing value (backward fill) in the sequence.

Linear Regression Imputation:

When missing data is related to other variables in the dataset, linear regression imputation can be effective. This method is suitable when the missingness is not completely random (e.g., missing values are related to some other observed variables). Use linear regression to predict missing values based on the relationship with other variables. Fit a regression model with the variable containing missing values as the target variable and other relevant variables as predictors.

KNearest Neighbors (KNN) Imputation:

KNN imputation is useful when missing values are related to nearby data points or neighbors in a multidimensional space. It's effective for datasets with complex, nonlinear relationships. Find the knearest data points with complete information and use their values to impute the missing values. The imputed value is typically a weighted average of the neighbors' values.

Multiple Imputation:

Multiple Imputation is a versatile technique suitable for datasets with missing data that may be missing not completely at random (MNAR). It's a probabilistic imputation method that provides multiple imputed datasets. It can be solved by generating multiple datasets with different imputed values for missing data, incorporating uncertainty into the imputation process. Each imputed dataset should be analysed separately and combine results to account for imputation uncertainty.

5. Explain two methods for identifying outliers in a dataset and discuss when each method is appropriate.

Identifying outliers in a dataset is essential for data quality assessment and anomaly detection. Outliers are data points that deviate significantly from the majority of the data and can have a substantial impact on statistical analysis. Here are two common methods for identifying outliers and when each method is appropriate:

Z-Score Method:

The Z-score method is suitable for datasets where the distribution of the data is approximately normal (Gaussian). It assumes that the data follows a bellshaped curve and calculates how many standard deviations a data point is away from the mean.

Steps for computing outliers are given below:

Calculate the mean (μ) and standard deviation (σ) of the dataset.

For each data point, compute the Z-score using the formula $Z = (X - \mu) / \sigma$, where X is the data point.

Set a threshold (e.g., $Z > 3$ or $Z < -3$) to identify extreme values as outliers. This threshold is often chosen based on the desired level of significance (e.g., 3 standard deviations represents a 99.7% confidence interval).

Advantages:

Easily interpretable and can be applied to normally distributed data.

Provides a quantifiable measure of how extreme a data point is.

Disadvantages:

May not perform well with nonnormal distributions.

Sensitivity to the chosen threshold.

IQR (Interquartile Range) Method:

The IQR method is robust to nonnormal distributions and is suitable for datasets with skewed or non-Gaussian distributions. It identifies outliers based on the spread of data within the interquartile range.

Steps for computing outliers are given below:

Calculate the first quartile (Q1) and third quartile (Q3) of the dataset.

Compute the IQR as $IQR = Q3 - Q1$.

Define lower and upper boundaries as $Q1 + 1.5 * IQR$ and $Q3 + 1.5 * IQR$, respectively.

Any data point below the lower boundary or above the upper boundary is considered an outlier.

Advantages:

Robust to nonnormal distributions and resistant to extreme values.

The $1.5 * IQR$ threshold is a common rule of thumb, but it can be adjusted to be more or less stringent depending on the application.

Disadvantages:

It may not be as informative about how extreme a data point is compared to the Zscore method.

The choice of the $1.5 * IQR$ multiplier can be somewhat arbitrary and may need adjustment based on the specific dataset.

6. What is the purpose of standardizing or normalizing data values in a dataset, and how is it done?

The purpose of standardizing or normalizing data values in a dataset is to transform the data in a way that makes it more suitable for certain types of analysis, particularly when the features (variables) have different scales or units. The key goals are:

Feature Scaling: Ensuring that all features are on a common scale. This is important for machine learning algorithms that use distancebased measures (e.g., Euclidean distance) or gradientbased optimization, as features with larger scales can dominate the analysis.

Improved Model Convergence: Scaling can help numerical optimization algorithms converge faster and improve the stability and performance of machine learning models.

MinMax Scaling (Normalization):

MinMax scaling transforms the data to a specific range, typically $[0, 1]$, making all features fall within this range. This method is particularly useful when you want to preserve the relationships and proportions between data points.

For each feature, use the following formula to normalize the data:

$$X_{\text{normalized}} = (X - X_{\text{min}}) / (X_{\text{max}} - X_{\text{min}})$$

Standardization (ZScore Scaling):

Standardization transforms data to have a mean of 0 and a standard deviation of 1. This method centers the data around the mean and scales it to have unit variance. It is particularly useful when the features have different units or scales, and when you want to remove the effect of the scale from the analysis.

For each feature, use the following formula to standardize the data:

$$X_{\text{standardized}} = (X - X_{\text{mean}}) / X_{\text{std}}$$

The following is an example in Python for standardization using the Zscore method with the `scikitlearn` library:

```
from sklearn.preprocessing import StandardScaler
import numpy as np
data = np.array([1.0, 2.0, 3.0, 4.0, 5.0]).reshape(1, 1)
scaler = StandardScaler()
scaled_data = scaler.fit_transform(data)
```

7. Please illustrate the essential stages of conducting exploratory data analysis (EDA) on a dataset by providing an example.

Exploratory Data Analysis (EDA) is a critical initial step in data analysis that helps you understand the main characteristics of your dataset, identify patterns, and gain insights before performing more complex analyses or modeling.

The following are the steps for performing EDA:

Consider a dataset containing information about customers' purchases in an online store. The dataset includes columns for customer ID, purchase amount, purchase date, customer age, and product category.

Data Collection and Loading: Start by obtaining the dataset, which may come from various sources like databases, spreadsheets, or APIs. Load the dataset into your data analysis environment (e.g., Python with pandas).

Data Inspection: Get an overview of the dataset by examining its structure, data types, and the first few rows. Check for any missing values or anomalies in the data.

Summary Statistics: Calculate summary statistics to understand the central tendency, spread, and distribution of numerical variables.

Data Visualization: Create visualizations to explore the data's distribution, relationships, and patterns. Use histograms, scatter plots, box plots, and other charts to visualize your data.

Data Transformation and Cleaning: Handle missing values, outliers, and data inconsistencies as needed. Perform feature engineering or transformation if required for further analysis.

Hypothesis Testing and Exploration: Formulate hypotheses and conduct statistical tests to explore relationships and patterns in the data. For example, you could test if there is a significant difference in purchase amounts between different product categories using ttests or ANOVA.

Correlation Analysis: Examine the relationships between numerical variables using correlation analysis (e.g., Pearson correlation). Visualize correlations with heatmaps. (Provide python code for the same in exam if asked)

UNIT 2

1. What are some common summary statistics used in univariate analysis for numerical data, and why are they important?

Summary statistics are essential for understanding the characteristics of numerical data in univariate analysis. They provide a concise and informative overview of the data distribution. Some common summary statistics used in univariate analysis for numerical data and their importance are:

Mean (Average):

The mean represents the center of the data distribution and provides insight into the central tendency. It's a measure of the average value and can help identify whether the data is skewed in one direction.

Formula: $\text{Mean} = (\text{Sum of all data points}) / (\text{Number of data points})$

Median

The median is the middle value when the data is sorted in ascending order. It is a robust measure of central tendency that is less affected by extreme outliers. Median is useful when the data is skewed or has outliers.

Calculation: Find the middle value when the data is sorted, or for an odd number of data points, it's the middle data point. For an even number of data points, it's the average of the two middle values.

Mode

The mode is the most frequently occurring value in the dataset. It helps identify the most common value in the data, which can be important for categorical or discrete numerical data.

Calculation: Count the frequency of each value and identify the one with the highest count.

Variance:

Variance measures the spread or dispersion of the data points. It quantifies how data points deviate from the mean. A high variance indicates that data points are spread out, while a low variance suggests data points are close to the mean.

Formula: $\text{Variance} = \Sigma (X - \text{Mean})^2 / (\text{Number of data points})$

Standard Deviation:

The standard deviation is the square root of the variance. It provides a measure of the average deviation from the mean. A higher standard deviation implies greater data variability.

Formula: $\text{Standard Deviation} = \sqrt{(\text{Variance})}$

Range:

The range is the difference between the maximum and minimum values in the dataset. It offers a simple way to understand the data's spread. However, it is sensitive to outliers.

Calculation: $\text{Range} = \text{Max} - \text{Min}$

Interquartile Range (IQR)

IQR measures the spread of the central 50% of the data, making it robust against outliers. It is used to identify the range within which most data points lie.

Calculation: $IQR = Q3 \text{ (Third Quartile)} - Q1 \text{ (First Quartile)}$

Skewness

Skewness quantifies the asymmetry of the data distribution. Positive skew indicates a tail on the right (data skewed to the right), while negative skew indicates a tail on the left.

Calculation: $Skewness = (3 * (\text{Mean} - \text{Median})) / \text{Standard Deviation}$

Kurtosis

Kurtosis measures the "tailedness" of the data distribution. High kurtosis indicates heavier tails, while low kurtosis indicates lighter tails compared to a normal distribution.

Calculation: $Kurtosis = \sum (X - \text{Mean})^4 / (\text{Standard Deviation})^4$

2. Differentiate between categorical ordered and categorical unordered variables in data analysis.

In data analysis, variables can be categorized as either categorical ordered or categorical unordered (nominal) based on the nature of the data and the relationships between categories.

Categorical Ordered Variables (Ordinal Variables):

Categorical ordered variables represent data with distinct categories or groups, but these categories have a meaningful order or hierarchy.

Ex: Educational levels (e.g., elementary school, high school, college, postgraduate), customer satisfaction ratings (e.g., very dissatisfied, dissatisfied, neutral, satisfied, very satisfied), income brackets (e.g., low income, middle income, high income).

Characteristics:

Categories have a specific, meaningful order or ranking.

The order indicates the relative position, level, or ranking of each category with respect to the others.

The intervals between categories may not be uniform or consistent.

Arithmetic operations like addition and subtraction are not meaningful for ordinal data.

Categorical Unordered Variables (Nominal Variables):

Categorical unordered variables represent data with distinct categories or groups, but there is no inherent order or hierarchy among these categories.

Ex: Colors (e.g., red, green, blue, yellow), animal types (e.g., cats, dogs, birds, fish), geographic regions (e.g., North, South, East, West).

Characteristics:

Categories have no natural order or ranking; they are merely labels.

The categories are mutually exclusive and distinct, with no inherent relationship between them.

Arithmetic operations like addition and subtraction are not meaningful for nominal data.

Frequency counts and proportions are typical summary statistics for nominal data.

3. How can the mean, median, and standard deviation help in understanding the distribution of numerical data?

The mean, median, and standard deviation are fundamental summary statistics that provide valuable insights into the distribution of numerical data. Each of these statistics reveals different aspects of the data distribution, and together, they offer a comprehensive view of the data's central tendency, variability, and shape.

Mean (Average):

The mean represents the central tendency of the data distribution. It is calculated by summing all data points and dividing by the total number of data points.

The mean helps you understand where the "center" of the data is located. When the data distribution is approximately symmetrical and normally distributed, the mean often coincides with the peak of the distribution (center of the bell curve). However, in the presence of outliers or skewness, the mean may be influenced by extreme values. Mean is useful for summarizing the location of the data.

Median:

The median is the middle value in the sorted dataset. It divides the data into two equal halves, with half of the data points falling below it and half above it.

The median is robust to extreme values (outliers) and is less affected by skewed data. It provides a better representation of central tendency when the distribution is not symmetric. It is particularly useful when you want to understand the typical value that is less influenced by extreme data points. The median provides a robust measure of central tendency that is less sensitive to outliers or skewed data.

Standard Deviation:

The standard deviation measures the spread or dispersion of the data points around the mean. A small standard deviation indicates that most data points are close to the mean, while a large standard deviation suggests that data points are more spread out.

The standard deviation helps in understanding the variability, consistency, or uncertainty associated with the data. In a normal distribution, about 68% of data points fall within one standard deviation of the mean, 95% within two standard deviations, and 99.7% within three standard deviations.

4. Differentiate between correlation and causation in the context of data analysis, and provide an example to illustrate the difference.

Correlation:

Correlation refers to a statistical relationship or association between two variables. It measures how changes in one variable are related to changes in another variable. Correlation does not imply causation, meaning that even if two variables are correlated, it does not necessarily mean that one variable causes the other. Correlation is quantified using correlation coefficients, such as the Pearson correlation coefficient (r), which ranges from -1 to 1. A positive value indicates a positive correlation (as one variable increases, the other tends to increase), a negative value indicates a negative correlation (as one variable increases, the other tends to decrease), and zero indicates no correlation.

Causation:

Causation implies a cause-and-effect relationship between two variables, where a change in one variable directly influences or causes a change in another variable. Causation asserts a directional link between variables, whereas correlation only reflects an association. Establishing causation typically requires conducting controlled experiments and demonstrating that a change in one variable leads to a predictable change in another. It involves establishing temporal precedence (cause precedes effect), demonstrating a correlation, and ruling out alternative explanations.

In a study, there's a positive correlation between ice cream sales and the number of drowning incidents at the beach. When ice cream sales increase, drowning incidents also tend to increase. This correlation is relatively strong (high positive r). However, it would be incorrect to conclude that ice cream sales cause drownings. This is an example of a correlation without causation. To establish causation in this scenario, there is a need to conduct controlled experiments, manipulate one variable while keeping others constant (e.g., temperature), and demonstrate that changes in ice cream sales directly lead to changes in drowning incidents, while ruling out alternative explanations.

5. Describe a common technique for comparing a numerical variable across different categories in a dataset.

A common technique for comparing a numerical variable across different categories in a dataset is to use a box plot, also known as a box-and-whisker plot. Box plots are effective for visualizing the distribution of numerical data within each category or group, making it easy to compare and identify differences in the central tendency, spread, and the presence of outliers.

In Python, libraries like Matplotlib or Seaborn can be used to create box plots. The box plot displays the distribution of the numerical variable within each category or group.

```
import seaborn as sns
import matplotlib.pyplot as plt
data = sns.load_dataset("iris")
plt.figure(figsize=(8, 6))
sns.boxplot(x="species", y="sepal_length", data=data)
plt.title("Comparison of Sepal Length Across Species")
```

```
plt.xlabel("Species")
plt.ylabel("Sepal Length")
plt.show()
```

The box in the plot represents the interquartile range (IQR), which contains the central 50% of the data. The horizontal line inside the box is the median, which represents the central value of the data distribution. The whiskers (lines extending from the box) typically indicate the data's range. They are often set to a specified distance (e.g., 1.5 times the IQR) from the quartiles, and data points beyond the whiskers are considered potential outliers. Individual data points outside the whiskers are often plotted as dots or small circles, indicating potential outliers. The box plot is organized by the categories or groups being compared, making it easy to compare the distribution of the numerical variable for each category.

By examining the box plots for different categories, it is possible to compare the central tendency (medians), spread (IQR), and the presence of outliers for the numerical variable within each category. Box plots are particularly useful for identifying variations between categories, assessing data symmetry, and spotting potential outliers.

6. Explain how scatter plots can be used to identify relationships between two numerical variables.

Scatter plots are a powerful visualization tool that can be used to identify relationships between two numerical variables. They provide a visual representation of the data points in a twodimensional space, making it easier to understand patterns, associations, and trends between the two variables.

Scatter plots display individual data points as dots on a Cartesian plane, with one variable on the xaxis and the other on the yaxis. Each dot represents a single data point, allowing us to see the values of both variables simultaneously. By examining the overall pattern of data points in the scatter plot, you can identify trends or relationships between the two variables. A clear pattern may suggest a relationship, while a random scattering of points may indicate no significant association.

Scatter plots can reveal different types of relationships, including:

Positive Linear Relationship: Data points tend to form a straight line from bottomleft to topright, indicating that as one variable increases, the other also tends to increase.

Negative Linear Relationship: Data points form a straight line from topleft to bottomright, showing that as one variable increases, the other tends to decrease.

No Relationship (Random Scatter): Data points are scattered randomly with no discernible pattern.

Curvilinear Relationship: Data points follow a curved or nonlinear pattern, suggesting a complex relationship.

Clustered Data Points: Groups or clusters of data points may suggest subpopulations or categories.

Scatter plots can also help identify outliers, which are data points that deviate significantly from the general trend. Outliers often appear as individual points far away from the main cluster of data. Outliers can have a substantial impact on the relationship between variables, and their identification is essential in understanding data patterns.

Scatter plots are often used in conjunction with correlation analysis to quantify the strength and direction of the relationship between the two variables. The Pearson correlation coefficient (r) is a common measure to determine the degree of linear association, with values ranging from -1 (perfect negative correlation) to 1 (perfect positive correlation).

UNIT 3

1. Explain the key properties that distinguish permutations from combinations.

Permutations and combinations are fundamental concepts that involve selecting or arranging elements from a set.

Permutations

Permutations take into account the order in which elements are arranged or selected from a set. Permutations can either allow repetition or not, leading to two types:

Permutations with Repetition (nPr): In this type, an element can be selected more than once. For example, when arranging the letters "A," "B," and "C" with repetition, you can have "ABC," "BCA," "CAB," and so on.

Permutations without Repetition (nPr): In this type, each element can only be selected once. For example, arranging the letters "A," "B," and "C" without repetition leads to "ABC," "ACB," "BAC," "BCA," "CAB," and "CBA."

Combinations:

Combinations do not consider the order in which elements are arranged or selected from a set. In other words, the arrangement does not matter.

Combinations can either allow repetition or not, leading to two types:

Combinations with Repetition (nCr): In this type, an element can be selected more than once, but the combinations are distinct based on the elements chosen. For example, selecting "A," "B," and "C" with repetition can result in combinations like "AAA," "AAB," "ABB," and so on.

Combinations without Repetition (nCr): In this type, each element can only be selected once, and the combinations are distinct based on the elements chosen. For example, selecting "A," "B," and "C" without repetition results in combinations like "AB," "AC," and "BC."

2. Define probability and explain its importance in the field of statistics and data analysis.

Probability is a branch of mathematics that quantifies uncertainty and randomness. It provides a way to measure the likelihood of different outcomes or events in various situations.

An event is an outcome or a set of outcomes.

A sample space is the set of all possible outcomes.

A probability distribution assigns a probability to each outcome or event in the sample space, with the sum of probabilities equaling 1.

Probability is essential in the field of statistics and data analysis for several reasons:

Quantifying Uncertainty: Probability provides a rigorous framework for dealing with uncertainty, allowing statisticians and data analysts to make informed decisions in the

presence of randomness. It helps assess the likelihood of different outcomes and understand the inherent variability in data.

Statistical Inference: In statistics, we often work with samples to make inferences about populations. Probability theory underlies statistical methods, such as hypothesis testing, confidence intervals, and regression analysis, which help us draw conclusions about populations based on sample data.

Modeling Data: Probability distributions, such as the normal distribution, binomial distribution, and Poisson distribution, are used to model and describe the behavior of random variables. These distributions provide a mathematical basis for understanding data patterns and making predictions.

Risk Assessment: Probability is crucial in risk assessment and decisionmaking. It allows organizations and individuals to evaluate and manage risks by estimating the likelihood of various outcomes and their associated consequences.

Machine Learning and Data Science: Probability plays a key role in machine learning and data science, particularly in areas like Bayesian statistics, where it is used to build probabilistic models for making predictions, clustering, classification, and anomaly detection.

Sampling and Survey Design: In survey sampling, probability sampling methods are used to ensure that the sample is representative of the population. This ensures that the results of a survey are generalizable to the larger population.

3. Differentiate between mutually exclusive and independent events in probability theory.

Mutually Exclusive Events:

Mutually exclusive events are events that cannot occur simultaneously. If one of these events happens, it excludes the possibility of the other event occurring at the same time.

If A and B are two mutually exclusive events, then the probability of both A and B occurring is $P(A \cap B) = 0$. In other words, the intersection of mutually exclusive events is an empty set.

Example: When rolling a six sided die, the events "getting an even number" and "getting an odd number" are mutually exclusive. You cannot roll a number that is both even and odd.

Mutually exclusive events cannot occur at the same time, and their intersection has a probability of 0.

Independent Events:

Independent events are events in which the occurrence or nonoccurrence of one event does not affect the probability of the other event occurring. In other words, the outcomes of independent events are unrelated.

If A and B are two independent events, then the probability of both A and B occurring is equal to the product of their individual probabilities: $P(A \cap B) = P(A) * P(B)$. When flipping a coin, the event "getting heads" and the event "rolling a six on a die" are independent. The outcome of the coin flip has no impact on the probability of rolling a six.

Independent events can occur simultaneously, but the probability of both events happening is the product of their individual probabilities.

4. Illustrate the addition rule and multiplication rule in probability with the help of an example.

Addition Rule:

The addition rule is used to calculate the probability of either of two mutually exclusive events occurring.

Consider a standard deck of 52 playing cards. We want to calculate the probability of drawing a red card (hearts or diamonds) from the deck. In this example, there are two mutually exclusive events: drawing a red card from the hearts suit (event A) and drawing a red card from the diamonds suit (event B).

Probability of drawing a red card from hearts (event A): There are 26 red hearts in the deck.

$$P(A) = \text{Number of favourable outcomes} / \text{Total possible outcomes} = 26 / 52 = 1/2.$$

Probability of drawing a red card from diamonds (event B): There are 26 red diamonds in the deck.

$$P(B) = \text{Number of favourable outcomes} / \text{Total possible outcomes} = 26 / 52 = 1/2.$$

Now, to find the probability of drawing a red card using the addition rule, we calculate the probability of either event A or event B occurring:

$$P(A \text{ or } B) = P(A) + P(B) = (1/2) + (1/2) = 1.$$

So, the probability of drawing a red card from the deck using the addition rule is 1 (100%).

Multiplication Rule:

The multiplication rule is used to calculate the probability of two or more independent events occurring.

In this example, there is a need to calculate the probability of drawing a red card from the hearts suit (event A) and then drawing another red card from the diamonds suit (event B).

Probability of drawing a red card from hearts (event A) is $1/2$ (as calculated earlier).

Now, to find the probability of drawing a red card from the diamonds suit after drawing a red card from hearts (using the multiplication rule for independent events), we calculate:

$$P(A \text{ and } B) = P(A) * P(B) = (1/2) * (1/2) = 1/4.$$

So, the probability of drawing a red card from the hearts suit and then drawing another red card from the diamonds suit using the multiplication rule is $1/4$.

5. Define random variable. Explain the concept of a probability distribution and why it is important in statistics.

Random Variable:

A random variable is a mathematical concept used in probability theory and statistics to quantify and represent uncertain or random events or outcomes. It assigns a numerical value to each possible outcome in a given probability space. Random variables can be discrete or continuous, depending on whether the set of possible outcomes is countable or uncountable.

For example, when rolling a six-sided die, you can define a random variable X that represents the outcome of the roll. X can take on values from 1 to 6, each corresponding to one of the possible outcomes.

Probability Distribution:

A probability distribution is a mathematical function or a description that specifies the probabilities associated with each possible outcome of a random variable. It provides a systematic way to describe the likelihood of each outcome occurring.

There are two main types of probability distributions:

1. **Discrete Probability Distribution:** This type is used when the random variable is discrete, meaning it can take on specific, separate values with gaps in between. Examples of discrete probability distributions include the binomial distribution, Poisson distribution, and geometric distribution.
2. **Continuous Probability Distribution:** This type is used when the random variable is continuous, meaning it can take on any value within a range. Examples of continuous probability distributions include the normal distribution, exponential distribution, and uniform distribution.

Probability distributions are fundamental in statistics for several reasons:

1. **Modeling and Analysis:** Probability distributions allow statisticians and data analysts to model and analyze the behavior of random variables. By understanding the underlying distribution, you can make inferences and predictions about real-world phenomena.
2. **Parameter Estimation:** Probability distributions often have parameters (e.g., mean and variance) that describe their characteristics. Estimating these parameters from data is a crucial part of statistical analysis.
3. **Hypothesis Testing:** Probability distributions play a central role in hypothesis testing. They help determine whether observed data is consistent with a particular hypothesis, providing a basis for statistical significance tests.
4. **Sampling Distributions:** The probability distribution of a statistic, called a sampling distribution, is essential for making inferences about a population from a sample. The Central Limit Theorem, for example, describes the properties of the sampling distribution of the sample mean.
5. **Prediction and Forecasting:** Probability distributions are used in forecasting and predictive modeling. For example, the normal distribution is commonly used to model and predict values in various fields, such as economics and finance.

6. **Define the expected value (mean) of a probability distribution and explain its significance in decision making.**

Expected Value (Mean), denoted as $E(X)$ or μ (mu) represents the long-term or average value of a random variable in a probability distribution. The expected value is calculated by summing the products of each possible value of the random variable and its corresponding probability of occurrence. For a discrete random variable, the formula for the expected value is:

$$E(X) = \sum [x * P(X=x)]$$

For a continuous random variable, the expected value is calculated using an integral instead of a summation.

The expected value has significant importance in decision making and risk analysis for several reasons:

1. **Measure of Central Tendency:** The expected value is a measure of central tendency, providing a single numerical summary that represents the "average" value or outcome. It helps in understanding the typical value you can expect from a random variable.
2. **Decision Criteria:** In decision-making scenarios with uncertainty, the expected value can serve as a criterion for choosing between different options or strategies. It allows decision-makers to compare options based on their expected outcomes.
3. **Risk Assessment:** When considering uncertain outcomes, it is not enough to look at the expected value alone. Decision-makers must also consider the associated risks. For this reason, the expected value is often used in conjunction with measures of variability, such as standard deviation or variance, to assess risk.
4. **Utility Theory:** In economics and decision theory, the concept of utility is used to measure an individual's preferences for different outcomes. The expected value can be incorporated into utility theory to make rational decisions based on personal preferences and risk tolerance.
5. **Insurance and Finance:** In the fields of insurance and finance, the expected value is used to assess the profitability or risk associated with investments, insurance policies, and financial decisions. It helps in pricing insurance premiums, valuing assets, and making investment choices.
6. **Game Theory:** In game theory and strategic decision-making, the expected value is used to evaluate the potential outcomes of different strategies and inform decision-makers about the best course of action.
7. **Project Management:** In project management and risk analysis, the expected value is used to estimate the likely cost or duration of a project. It helps project managers plan for contingencies and allocate resources efficiently.
8. **Statistical Inference:** The expected value is a key concept in statistical inference. When estimating population parameters from a sample, the sample mean is often used as an

estimator for the population mean. The expected value of an estimator indicates its bias and efficiency.

7. Describe the characteristics of the binomial distribution and provide an example of a real world situation that follows a binomial distribution.

The binomial distribution is a probability distribution that models the number of successes (usually denoted as "x") in a fixed number of independent and identical trials, where each trial has two possible outcomes: success (usually denoted as "S") or failure (usually denoted as "F").

The characteristics of the binomial distribution are as follows:

The binomial distribution considers a fixed number of trials (n), meaning you conduct the same experiment or trial a specified number of times.

Each trial is independent of the others. The outcome of one trial does not affect the outcomes of subsequent trials.

Each trial has only two possible outcomes: success or failure.

The probability of success (p) remains constant for each trial.

The random variable x, representing the number of successes, can only take on discrete values, typically non-negative integers (0, 1, 2, ...).

The probability mass function (PMF) for the binomial distribution is given by the binomial formula:

$$P(X = x) = \binom{n}{x} * p^x * (1 - p)^{(n - x)}$$

A classic real-world example of a situation that follows a binomial distribution is flipping a fair coin. The number of times the coin is flipped is decided beforehand. Each coin flip is independent. The outcome of one flip does not affect the outcome of subsequent flips. Each coin flip has two possible outcomes: heads (success) or tails (failure). Assuming the coin is fair, the probability of getting heads (a success) is 0.5 for each flip. The random variable x represents the number of heads you get in the 10 flips. It can only take on discrete values, such as 0, 1, 2, ..., 10. The binomial distribution formula can be used to calculate the probability of getting a specific number of heads in 10 coin flips.

For example, the binomial distribution can be used to find the probability of getting exactly 3 heads in 10 coin flips, given that the probability of getting a head in a single flip is 0.5.

The binomial distribution is applicable in various other scenarios, such as quality control in manufacturing, the success rate of marketing campaigns, and more, where there is a fixed number of trials with two possible outcomes and a constant probability of success.

8. How many different possible ways are there to create a 4-digit PIN code using the digits 0-9 such that no digit is repeated?

Given, $n = 10$, $r = 4$

Number of PIN codes = Permutations ($10P_4$)

$$= 10! / (10 - 4)!$$

$$= 10! / 6!$$

$$= (10 * 9 * 8 * 7 * 6!) / 6!$$

$$= 10 * 9 * 8 * 7$$

$$= 5,040$$

So, there are 5,040 different PIN codes possible when using 4 distinct digits from the set of digits 0-9.

9. The player is supposed to draw two cards successively (one after the other) without replacement from a deck of 52 cards. What is the probability of drawing at least one Ace in these two draws?

The deck has 52 cards initially, and there are 48 non-Aces.

$$P(\text{not drawing an ace in first draw}) = 48/52$$

Since one Ace has been removed from the deck in the first draw, there are now 51 cards in the deck, and 47 of them are non-Aces.

$$P(\text{not drawing an ace in second draw}) = 47/51$$

$$P(\text{not drawing an ace in both draws}) = (48/52) * (47/51) = (8/13) * (47/51)$$

$$\begin{aligned} P(\text{drawing at least one Ace}) &= 1 - [(8/13) * (47/51)] = 1 - (376/663) \\ &= 287/663 \end{aligned}$$

So, the probability of drawing at least one Ace is 287/663.

The addition rule states that the probability of either of two mutually exclusive events occurring is the sum of their individual probabilities.

There are 4 Aces in a deck of 52 cards.

$$P(\text{drawing an Ace in the first draw}) = 4/52$$

Since one Ace has been removed from the deck in the first draw, there are now 51 cards in the deck, and 3 of them are Aces.

$$P(\text{drawing an Ace in the second draw}) = 3/51$$

Using the addition rule, $P(\text{drawing at least one Ace}) = (4/52) + (3/51) = (1/13) + (1/17)$
 $= (17 + 13) / (13 * 17)$
 $= 30 / 221$

The result obtained using the addition rule matches the complement probability calculated earlier, which is 30/221.

Thus addition rule is verified.

UNIT-4

1. Explain Bayes' Theorem and provide an example of its application in medical diagnosis.

Answer: Bayes' Theorem is a mathematical formula used to update the probability for a hypothesis based on new evidence. In medical diagnosis, it's used to revise the probability of a disease given a positive or negative test result. For example, if a patient has a 95% probability of having a disease based on symptoms and a test with 90% accuracy is positive, Bayes' Theorem can help calculate the updated probability.

2. What is the significance of the Z-score in statistical analysis, and how is it calculated? Calculate the Z-score for a student who scored 120 on a standardized test with a mean of 100 and a standard deviation of 15.

Answer: The Z-score quantifies how many standard deviations a data point is from the mean in a normal distribution. It's calculated as $Z = (X - \mu) / \sigma$, where X is the data point, μ is the mean, and σ is the standard deviation. For a student who scored 120, the Z-score is $(120 - 100) / 15 = 1.33$.

3. Describe simple random sampling, systematic sampling, and stratified sampling in the context of survey design.

Answer:

- **Simple Random Sampling:** Involves selecting a random sample from the population where every member has an equal chance of being chosen.
- **Systematic Sampling:** Selects every "kth" member from a list after an initial random start.
- **Stratified Sampling:** Divides the population into subgroups or strata and then samples proportionally from each stratum.

4. Explain the advantages of stratified sampling in the context of conducting a customer satisfaction survey in a large shopping mall.

Answer: Stratified sampling is advantageous in this context because:

- It ensures representation from various customer segments (strata).
- It provides a more accurate estimate of overall customer satisfaction.
- It allows for comparisons between different customer groups.
- It can help identify specific areas for improvement based on stratum-specific feedback.

5. In a medical diagnosis scenario, describe the key components of Bayes' Theorem and how they relate to updating the probability of a disease.

Answer: The key components of Bayes' Theorem are:

- Prior Probability ($P(A)$): The initial probability of having the disease based on symptoms.
- Likelihood ($P(B|A)$): The probability of a positive or negative test result given the disease's presence.
- Evidence ($P(B)$): The overall probability of getting the test result.
- Posterior Probability ($P(A|B)$): The updated probability of having the disease after the test.

Bayes' Theorem relates these components to update the disease probability based on the test result, allowing for more accurate diagnoses.

6. What does a Z-score of -1.5 indicate in the context of a standardized test with a mean of 100 and a standard deviation of 15?

Answer: A Z-score of -1.5 indicates that the student's score is 1.5 standard deviations below the mean. It means the student's score is relatively lower compared to the average performance on the test.

7. What role does the significance level (α) play in hypothesis testing, and how is it typically chosen in practice?

Answer: The significance level (α) determines the probability threshold for making decisions in hypothesis testing. It represents the acceptable risk of Type I error (incorrectly rejecting the null hypothesis). Common choices for α are 0.05 or 0.01, and they indicate the level of confidence desired in the test results.

8. In stratified sampling, how are strata determined, and why is it important to have well-defined strata in the sampling process?

Answer: Strata are determined by categorizing the population into subgroups based on relevant characteristics (e.g., age, gender). Well-defined strata are important because they ensure that each subgroup is represented and that the sample reflects the population's diversity, leading to more accurate and meaningful results.

9. Can you explain the concept of the posterior probability in the context of Bayes' Theorem, and how is it useful in medical diagnosis?

Answer: The posterior probability, in the context of Bayes' Theorem, represents the updated probability of having a disease after considering a diagnostic test result. It's crucial in medical diagnosis as it provides a more accurate estimate of the disease's likelihood, incorporating both prior information (symptoms) and new evidence (test results).

10. When using systematic sampling in survey design, what is the significance of the "k" value, and how does it impact the representativeness of the sample?

Answer: The "k" value in systematic sampling represents the interval between selected sample points. The choice of "k" influences the representativeness of the sample. If "k" is too large, important data points might be missed; if it's too small, it may not differ significantly from simple random sampling. Properly choosing "k" is critical to maintain a representative sample.

UNIT- 5

1. What is the purpose of the critical value method in hypothesis testing?

Answer: The critical value method helps in making decisions about the null hypothesis by comparing a test statistic to a predetermined critical value. If the test statistic exceeds the critical value, the null hypothesis is rejected, indicating evidence for the alternative hypothesis. If it's less than or equal to the critical value, the null hypothesis is not rejected.

2. What is the significance level (α) in hypothesis testing, and how does it relate to the critical value?

Answer: The significance level (α) sets the threshold for making decisions about the null hypothesis. The critical value is chosen based on α and divides the distribution into the critical region (where you reject the null hypothesis) and the non-critical region (where you fail to reject the null hypothesis).

3. How does the choice of significance level (α) impact the decision in hypothesis testing?

Answer: The choice of α affects the trade-off between Type I and Type II errors. A lower α (e.g., 0.01) makes it harder to reject the null hypothesis, reducing Type I errors but increasing Type II errors. A higher α (e.g., 0.10) makes it easier to reject the null hypothesis, reducing Type II errors but increasing Type I errors.

4. Explain the null hypothesis (H_0) in the context of hypothesis testing.

Answer: The null hypothesis (H_0) is a statement that suggests there is no effect, no difference, or no relationship in the population being studied. It serves as the default assumption to be tested and potentially refuted based on data and evidence.

5. In the education sector, provide an example of how the null hypothesis (H_0) and alternative hypothesis (H_a) can be formulated when testing whether two groups of students are equal in academic performance.

Answer: H_0 : There is no significant difference in academic performance between Group A and Group B. H_a : There is a significant difference in academic performance between Group A and Group B.

6. How do teachers in the education sector use hypothesis testing to assess student performance differences?

Answer: Teachers can collect data on student performance and use hypothesis testing to evaluate whether there is evidence to reject the null hypothesis (no difference) in favor of the alternative hypothesis (a difference). This helps teachers make informed decisions about student performance and tailor their teaching methods accordingly.

7. If a teacher conducts a hypothesis test and obtains a p-value of 0.03 at a significance level of 0.05, what decision should the teacher make regarding the null hypothesis?

Answer: With a p-value of 0.03 less than the significance level of 0.05, the teacher should reject the null hypothesis, indicating evidence of a significant difference in the context being tested.

8. Why is it important to have a clear null hypothesis (H_0) and alternative hypothesis (H_a) when conducting hypothesis testing in education or any other field?

Answer: Clear hypotheses provide a framework for testing and help in drawing meaningful conclusions from data. They guide the direction of the analysis and communicate the research question and expected outcomes, ensuring that the test's purpose is well-defined and results are interpretable.

9. Explain the concept of hypothesis testing and its key components in statistical analysis.

Answer: Hypothesis testing is a fundamental statistical method used to make inferences about a population based on sample data. The key components of a hypothesis test include:

- **Null Hypothesis (H_0):** This is the default assumption or the statement of no effect, no difference, or no relationship in the population. It represents the status quo or what you're trying to test against.
- **Alternative Hypothesis (H_a):** This is the statement that contradicts the null hypothesis, suggesting there is an effect, difference, or relationship in the population. It represents what you're trying to find evidence for.
- **Significance Level (α):** This is the probability threshold that determines when to reject the null hypothesis. It represents the acceptable level of Type I error, which is the probability of incorrectly rejecting the null hypothesis when it is true.
- **Test Statistic:** This is a calculated value that measures how the sample data supports the alternative hypothesis. The test statistic is compared to a critical value or used to calculate a p-value to make the decision regarding the null hypothesis.

10. In a medical study, researchers want to test whether a new drug is more effective than the existing treatment. Describe how they can set up the null hypothesis (H_0) and the alternative hypothesis (H_a) for this study.

Answer: In the medical study comparing a new drug to the existing treatment, the null hypothesis (H_0) and the alternative hypothesis (H_a) can be formulated as follows:

- **Null Hypothesis (H_0):** The new drug is equally as effective as the existing treatment. In statistical terms, it can be represented as $H_0: \mu_{\text{new}} = \mu_{\text{existing}}$, where μ_{new}

represents the effectiveness of the new drug and μ_{existing} represents the effectiveness of the existing treatment.

- **Alternative Hypothesis (H_a):** The new drug is more effective than the existing treatment. In statistical terms, this can be represented as $H_a: \mu_{\text{new}} > \mu_{\text{existing}}$, where μ_{new} represents the effectiveness of the new drug and μ_{existing} represents the effectiveness of the existing treatment.

UNIT- 6

UNIT – 6

- 1. Write the significance of the p-value approach in hypothesis testing. How does it differ from the critical value approach, and what role does the significance level play in these methods?**

Answer: The p-value is a key concept in hypothesis testing that indicates the strength of evidence against the null hypothesis. It quantifies the probability of observing the test statistic (or something more extreme) if the null hypothesis were true. The smaller the p-value, the stronger the evidence against the null hypothesis. In hypothesis testing, we compare the p-value to a predetermined significance level (usually denoted as α , e.g., 0.05). If the p-value is smaller than α , we reject the null hypothesis; otherwise, we fail to reject it.

- 2. What does a small p-value indicate in hypothesis testing? If you conduct a hypothesis test and obtain a p-value of 0.02 at a significance level of 0.05, what decision should you make regarding the null hypothesis? Provide an explanation for your decision.**

Answer: A small p-value (e.g., 0.02) in hypothesis testing indicates that the observed data is unlikely to have occurred under the assumption that the null hypothesis is true. In this case, with a significance level of 0.05, the p-value of 0.02 is less than 0.05. Therefore, you should reject the null hypothesis. This means that you have found statistically significant evidence to support the alternative hypothesis.

- 3. Compare and contrast the p-value approach with the critical value approach in hypothesis testing. Discuss the advantages and limitations of each method, and when one might be preferred over the other.**

Answer: The p-value approach and the critical value approach are two methods used in hypothesis testing. The p-value approach provides a direct measure of the strength of evidence against the null hypothesis, while the critical value approach involves comparing the test statistic to a critical value derived from a predetermined significance level. The choice between these approaches often depends on the context and the researcher's preference. The p-value approach is more flexible and allows for a continuous assessment of evidence, while the critical value approach can be more straightforward in some cases.

- 4. In the context of hypothesis testing, explain the concept of Type I error and Type II error. How does the choice of significance level (α) influence the likelihood of these errors, and how are p-values related to these errors in hypothesis testing?**

Answer: In hypothesis testing, Type I error occurs when the null hypothesis is incorrectly rejected when it is true, while Type II error occurs when the null hypothesis is not rejected when it is false. The choice of significance level (α) in hypothesis testing influences the trade-off between these errors. A lower α (e.g., 0.01) reduces the chance of Type I error

but increases the chance of Type II error, and vice versa. P-values are directly related to Type I error; a smaller p-value makes it less likely to commit a Type I error.

5. What is the purpose of A/B testing, and why is it essential in web design and marketing?

Answer: A/B testing, also known as split testing, is a method used to compare two or more variations of a webpage, email, or marketing campaign to determine which one performs better in achieving specific goals, such as increasing user engagement, conversion rates, or revenue. It is essential in web design and marketing for several reasons:

- **Data-Driven Decision Making:** A/B testing allows businesses to make decisions based on empirical data rather than assumptions or intuition. It provides concrete evidence of what changes have a positive impact on user behavior.
- **Optimizing User Experience:** Web design and marketing campaigns can significantly impact user experience. A/B testing helps identify the design elements, content, or strategies that resonate most with the target audience, leading to better user engagement and satisfaction.
- **Efficiency and Cost-Effectiveness:** Instead of making sweeping changes based on guesswork, A/B testing enables incremental improvements. This iterative approach can save time and resources by focusing efforts on changes that are more likely to yield positive results.
- **Continuous Improvement:** A/B testing fosters a culture of continuous improvement. By regularly testing and refining designs and strategies, businesses can adapt to changing user preferences and market conditions.

6. Explain the role of the null hypothesis and alternative hypothesis in the A/B testing scenario described. Why is it important to have a clear hypothesis?

Answer: In the A/B testing scenario, the null hypothesis (H_0) and alternative hypothesis (H_a) play crucial roles:

1. **Null Hypothesis (H_0):** This is the default assumption that there is no significant difference between the two variations being tested (Design A and Design B). It serves as a baseline or reference point.
2. **Alternative Hypothesis (H_a):** This is the hypothesis that we are trying to establish, suggesting that there is a statistically significant difference between the two designs, indicating that Design B is more effective in terms of user engagement.

Having clear hypotheses is vital for several reasons:

3. **Statistical Testing:** Hypotheses provide the framework for statistical testing. They allow us to quantify and evaluate the differences observed between variations.
4. **Objective Decision Making:** Hypotheses guide the decision-making process by providing a clear framework for accepting or rejecting the null hypothesis based on statistical evidence.

5. **Avoiding Ambiguity:** Clear hypotheses help in avoiding ambiguity and subjectivity in interpreting the results. They ensure that the criteria for success are well-defined and measurable.
6. **Communication:** Clear hypotheses make it easier to communicate the goals and expectations of the A/B test to stakeholders and team members.

7. What is the significance level, and why is it set at 0.05 in this A/B test scenario?

What are the consequences of choosing a lower or higher significance level?

Answer: The significance level (α) is the predetermined threshold for statistical significance in an A/B test. It represents the acceptable level of risk for making a Type I error (incorrectly rejecting the null hypothesis). In this scenario, it is set at 0.05 (5%).

1. **Setting at 0.05:** Choosing a significance level of 0.05 means that you are willing to accept a 5% chance of making a Type I error. In other words, if there is no real difference between the designs, you would wrongly reject the null hypothesis 5% of the time.
2. **Consequences of a Lower Significance Level:** Setting a lower α , e.g., 0.01, reduces the risk of Type I error but increases the risk of Type II error (failing to detect a real difference). It makes it more conservative and requires stronger evidence to reject the null hypothesis.
3. **Consequences of a Higher Significance Level:** Setting a higher α , e.g., 0.10, increases the risk of Type I error but decreases the risk of Type II error. It is less conservative and requires weaker evidence to reject the null hypothesis.

The choice of significance level should be based on the specific goals and risks associated with the A/B test, as well as industry standards and best practices.

8. If the A/B test yields a p-value of 0.03 at a significance level of 0.05, what decision should be made regarding the null hypothesis in this scenario, and why?

Answer: If the A/B test yields a p-value of 0.03 at a significance level of 0.05, the decision should be to reject the null hypothesis. Here's why:

1. **P-value Interpretation:** The p-value is a measure of the strength of evidence against the null hypothesis. A smaller p-value indicates stronger evidence against the null hypothesis.
2. **Comparison with Significance Level:** In this scenario, the significance level (α) is set at 0.05. If the p-value is less than or equal to α ($p \leq 0.05$), it means that the evidence is significant enough to reject the null hypothesis.
3. **Decision:** With a p-value of 0.03 ($p < 0.05$), you have statistically significant evidence to support the alternative hypothesis. This suggests that Design B is more effective than Design A in terms of user engagement, based on the data collected during the A/B test.

Formulae:

Statistical Test	Formula for Test Statistic	Explanation of Variables
Z-Test (One-Sample Z-Test)	$Z = (\bar{X} - \mu) / (\sigma/\sqrt{n})$	- Z: The test statistic representing the number of standard deviations a sample mean is from the population mean. - \bar{X} : The sample mean. - μ : The population mean (under the null hypothesis). - σ : The population standard deviation. - n: The sample size.
T-Test (One-Sample T-Test)	$t = (\bar{X} - \mu) / (s/\sqrt{n})$	- t: The test statistic used to assess the difference between a sample mean and a population mean. - \bar{X} : The sample mean. - μ : The population mean (under the null hypothesis). - s: The sample standard deviation. - n: The sample size.
Paired T-Test	$t = (\bar{X}_d - \mu_d) / (sd/\sqrt{n})$	- t: The test statistic used to compare the means of paired data. - \bar{X}_d : The sample mean of paired differences. - μ_d : The population mean of paired differences (under the null hypothesis). - sd: The standard deviation of paired differences. - n: The number of pairs.
Independent Two-Sample T-Test	$t = (\bar{X}_1 - \bar{X}_2) / \sqrt{[(s_1^2/n_1) + (s_2^2/n_2)]}$	- t: The test statistic used to compare means of two independent samples. - \bar{X}_1 and \bar{X}_2 : Sample means of the two independent samples. - s_1 and s_2 : Sample standard deviations of the two samples. - n_1 and n_2 : Sample sizes of the two samples.
Analysis of Variance (ANOVA)	$F = (SSB / (k - 1)) / (SSW / (N - k))$	- F: The test statistic used to compare means of multiple groups. - SSB: The sum of squares between groups. - SSW: The sum of squares within groups. - k: The number of groups. - N: The total number of observations.
Pearson's Correlation Coefficient	$r = \Sigma((X - \bar{X})(Y - \bar{Y})) / \sqrt{[\Sigma(X - \bar{X})^2 * \Sigma(Y - \bar{Y})^2]}$	- r: The correlation coefficient indicating the strength and direction of a linear relationship between two variables. - X and Y: Paired data points. - \bar{X} and \bar{Y} : The means of X and Y, respectively.

