

Answer File - Fliprobo || STATISTICS WORKSHEET- 1 ||

Task 4

Batch - DS2406

Name - Sakshi Jha

Ans-1. a. True - A Bernoulli random variable is a type of discrete random variable representing the outcome of a single Bernoulli trial, an experiment, or a process resulting in a binary outcome.

Ans-2. a. Central Limit Theorem

The theorem that states the distribution of averages of iid (independent and identically distributed) variables, properly normalized, becomes that of a standard normal as the sample size increases.

Ans-3. b) Modeling bounded count data

The Poisson distribution is used for modeling count data that can take any non-negative integer value and is typically unbounded, meaning there is no upper limit on the count. Bounded count data, where the counts are limited to a certain range, is not appropriately modeled by the Poisson distribution.

Ans-4. The correct statements are:

a) The exponent of a normally distributed random variable follows what is called the log-normal distribution.

c) The square of a standard normal random variable follows what is called the chi-squared distribution.

However, statement (b) is incorrect because the sum of normally distributed random variables is only guaranteed to be normally distributed if the variables are independent. If the variables are dependent, additional considerations are required to determine the distribution of the sum.

Ans-5 c. Poisson

Ans-6. b. False

Replacing the standard error by its estimated value does not fundamentally change the Central Limit Theorem (CLT). The CLT still holds, meaning that the distribution of the sample mean will approach a normal distribution as the sample size increases, regardless of whether the actual standard error or an estimated standard error is used.

Ans-7. b. Hypothesis

Ans-8. a. 0

Normalized data units are equal to standard deviations of the original data, as they are centered at zero.

Ans-9. b. The incorrect statement concerning outliers is:

c) Outliers cannot conform to the regression relationship

This statement is incorrect because outliers can sometimes conform to the regression relationship, even though they are extreme values compared to the rest of the data. Outliers may or may not fit the general pattern or trend represented by the regression model. The other statements (a and b) are correct as outliers can indeed have varying degrees of influence and can result from both spurious and real processes.

Ans-10. The Normal Distribution, or Gaussian distribution, is a fundamental concept in statistics and probability theory. It is characterized by several key features that make it widely applicable and important in various fields:

Shape and Symmetry: The distribution is symmetric around its mean (average). This means that if you were to plot a Normal Distribution on a graph, it would form a bell-shaped curve with the highest point directly above the mean.

Parameters: It is defined by two main parameters:

Mean (μ): This is the center of the distribution. It tells you where the peak of the bell curve is located.

- **Standard Deviation (σ):** This measures how spread out the data points are from the mean. A large standard deviation means the data points are more spread out.
- **Central Limit Theorem:** One of the most important properties of the Normal Distribution is its relationship to the Central Limit Theorem. This theorem states that if you take many samples from any population (regardless of its original distribution), and you calculate the average of each sample, those averages will be normally distributed around the population mean.

Empirical Rule: The Normal Distribution follows a rule of thumb known as the 68-95-99.7 rule:

About 68% of the data falls within one standard deviation of the mean ($\mu \pm \sigma$).

About 95% falls within two standard deviations ($\mu \pm 2\sigma$).

About 99.7% falls within three standard deviations ($\mu \pm 3\sigma$).

Applications: Normal distributions are incredibly useful in statistics and everyday life. They are used to model things like heights, weights, test scores, and many other measurements

that tend to cluster around an average with most values near the mean and fewer values further away.

In essence, the Normal Distribution is a fundamental concept in statistics that helps us understand and predict how data is spread out around an average value, making it a cornerstone of statistical analysis and inference.

Ans-11. Handling missing data is crucial in statistical analysis to ensure that the results are accurate and unbiased. Here are some common approaches and imputation techniques used to handle missing data:

1. Identify Missing Data

- **Missing Completely at Random (MCAR):** The missingness is unrelated to any variables, observed or unobserved.
- **Missing at Random (MAR):** The missingness is related to observed data but not to the missing values themselves.
- **Missing Not at Random (MNAR):** The missingness is related to the missing values themselves, i.e., the missingness depends on information that was not collected.

2. Handling Missing Data

- **Complete Case Analysis:** Exclude records with any missing values. This approach can lead to loss of data and potential bias if the missingness is related to the outcome.
- **Imputation:** Replace missing values with estimated values.

3. Imputation technique based on the above condition.

- **Consider the nature of missing data:** Choose an imputation method that aligns with whether the data is MCAR, MAR, or MNAR.
- **Multiple Imputation:** Generally preferred as it accounts for uncertainty in imputed values and provides more accurate statistical estimates.
- **Domain knowledge:** Use domain knowledge to guide the choice of imputation method and interpret results.
- **Evaluate assumptions:** Assess the assumptions underlying the chosen imputation method.
- **Sensitivity analysis:** Test the robustness of results to different imputation methods.

In practice, the choice of imputation method depends on the specific dataset, the missing data mechanism, and the goals of the analysis. Each method has its strengths and limitations, and it's essential to consider these factors when handling missing data in statistical analysis.

Ans-12. A/B testing, also referred to as split testing, is a statistical technique used to compare two versions of a webpage, app, marketing campaign, or other digital asset to determine which one performs better. It involves a controlled experiment in which variant A

(the control) is compared against variant B (the treatment) by randomly assigning users or subjects to one of the two groups.

Key Components of A/B Testing:

- 1. Randomization:** Users are randomly assigned to either group A or group B to ensure an unbiased assignment process and comparable groups.
- 2. Variant Implementation:** Each group is exposed to a different version of the asset being tested (e.g., webpage design, advertisement).
- 3. Outcome Measurement:** The performance of each variant is measured using a predefined metric, such as click-through rate, conversion rate, revenue generated, or any other relevant KPI (Key Performance Indicator).
- 4. Statistical Analysis:** Statistical methods are applied to determine if there is a statistically significant difference in performance between the two variants. This analysis helps decide whether to adopt the new variant (if it outperforms the control) or stick with the existing one.

Ans-13. Mean imputation, where missing values are replaced with the mean of observed values for that variable, is a straightforward and commonly used method to handle missing data. However, its acceptability and appropriateness depend on several factors and considerations:

Advantages of Mean Imputation:

- 1. Simple and Quick:** Mean imputation is easy to implement and computationally efficient, making it attractive for handling missing data in large datasets.
- 2. Preserve Sample Size:** It retains all observations in the dataset, avoiding data loss that can occur with other methods like complete case analysis.
- 3. Maintains Variable Distribution:** Mean imputation preserves the mean and variance of the variable, which can be important for maintaining statistical properties in subsequent analyses.

Disadvantages of Mean Imputation:

- 1. Bias Introduction:** Mean imputation can introduce bias, particularly if the missing data mechanism is not completely at random (MCAR). Imputing with the mean assumes that the missing values have the same mean as the observed values, which may not be true.
- 2. Underestimation of Variance:** Imputing missing values with the mean tends to underestimate the variance of the variable, potentially affecting statistical tests and confidence intervals.
- 3. Impact on Relationships:** Mean imputation ignores relationships between variables. If variables are correlated, imputing with the mean can distort these relationships.

4. **Not Suitable for Categorical Data:** Mean imputation is suitable for numerical data but not for categorical or ordinal data where the mean may not be meaningful.

Alternatives to Mean Imputation:

1. **Multiple Imputation:** Generates multiple imputed datasets and averages results to account for uncertainty. It provides more reliable estimates than single imputation methods like mean imputation.
2. **Regression Imputation:** Predict missing values using other variables as predictors in a regression model. This method accounts for relationships between variables.
3. **Model-Based Imputation:** Impute missing values using more complex models (e.g., decision trees, neural networks) that capture nonlinear relationships and interactions.
4. **Hot-Deck Imputation:** Assign missing values to the value of a similar record based on other variables.

Ans-14. Linear regression is a statistical method used to understand and model the relationship between two or more variables. It assumes that this relationship is linear, meaning changes in one variable are associated with proportional changes in another. Here are the key points:

Key Concepts of Linear Regression:

1. **Dependent Variable (Y):** This is the main variable we want to predict or explain based on other variables. It's typically continuous and numeric.
2. **Independent Variables (X):** These are the variables used to predict or explain the dependent variable. They can be numeric or categorical.
3. **Linear Relationship:** Linear regression assumes that the relationship between the independent and dependent variables can be approximated by a straight line.
4. **Model Equation:** The model tries to fit a line that best describes the relationship between Y and X.
5. **Objective:** The goal of linear regression is to find the best-fitting line that minimizes the difference between predicted values and actual values of Y. This helps in making predictions and understanding how changes in X affect Y.

Types and Applications:

- **Simple Linear Regression:** Involves one independent variable predicting a dependent variable.
- **Multiple Linear Regression:** Involves multiple independent variables predicting a dependent variable.

Linear regression is used widely across various fields, including economics, social sciences, health sciences, and business. Examples include predicting sales based on advertising spending, analyzing the impact of education and experience on income, or understanding how factors like temperature affect energy consumption.

Linear regression assumes that the relationship between variables is linear, observations are independent, residuals (errors) are normally distributed, and there's no significant multicollinearity (high correlation among independent variables).

Hence, Linear regression is a fundamental statistical tool for analyzing relationships between variables and making predictions. It provides insights into how changes in one variable can be expected to affect another, making it invaluable in data analysis, research, and decision-making processes.

Ans-15. In the context of data science, statistics plays a foundational role, providing essential tools and techniques for analyzing, interpreting, and making sense of data. Here are some of the key branches of statistics that are particularly relevant in data science:

- **Descriptive Statistics:** Descriptive statistics involves methods for summarizing and describing data. In data science, descriptive statistics are used to understand the basic properties of data sets, such as measures of central tendency (mean, median, mode), measures of dispersion (variance, standard deviation), and graphical representations (histograms, box plots).
- **Inferential Statistics:** Inferential statistics is crucial in data science for drawing conclusions and making predictions about populations based on sample data. Techniques such as hypothesis testing, confidence intervals, and regression analysis are used to infer relationships and patterns in data.
- **Probability Theory:** Probability theory provides the mathematical foundation for understanding uncertainty and randomness in data. In data science, probability theory is used to model and analyze probabilistic events, and distributions of data, and to quantify uncertainty in predictions and decisions.
- **Statistical Learning (Machine Learning):** Statistical learning, often synonymous with machine learning in data science, involves developing algorithms and models that can learn from data and make predictions or decisions. Techniques such as regression, classification, clustering, and dimensionality reduction fall under this branch.
- **Bayesian Statistics:** Bayesian statistics is increasingly relevant in data science for its approach to updating beliefs and making predictions based on prior knowledge and observed data. Bayesian methods are used in machine learning, decision-making under uncertainty, and modeling complex systems.

- **Time Series Analysis:** Time series analysis is critical in data science for analyzing data collected over time, such as stock prices, weather patterns, or sensor data. Techniques like autoregressive models, moving averages, and spectral analysis are used to uncover trends, and seasonality, and forecast future values.
- **Spatial Statistics:** Spatial statistics deals with data that have spatial or geographical components, such as maps, GPS data, or satellite imagery. In data science, spatial statistics is used for spatial interpolation, clustering, and modeling spatial relationships.
- **Experimental Design and Analysis:** Experimental design involves planning and conducting experiments to gather data and test hypotheses. Data scientists use techniques such as factorial designs, randomized controlled trials (RCTs), and analysis of variance (ANOVA) to design experiments effectively and draw valid conclusions.