

Statistics Worksheet - 1 Answers.

Q.1.Ans:- A) True

Q.2.Ans:- A) Central limit theorem

Q.3.Ans:- B) Modeling bounded count data

Q.4.Ans:- C) The square of a standard normal random variable follows what is called chi-squared distribution

Q.5.Ans:- C) Poisson

Q.6.Ans:- B) False

Q.7.Ans:- B) Hypothesis

Q.8.Ans:- A) 0

Q.9.Ans:- C) Outliers cannot confirm to the regression relationship

Q.10.Ans:- The Normal Distribution, also known as the Gaussian distribution or bell curve, is a fundamental concept in statistics and probability theory. It is characterized by a symmetric, bell-shaped curve that describes the probability distribution of a continuous random variable.

In a normal distribution:

1. The mean (μ) represents the central tendency of the distribution and is the point around which the curve is symmetrically centered.

2. The standard deviation (σ) determines the spread or dispersion of the data points around the mean. A larger standard deviation results in a wider and flatter curve, while a smaller standard deviation results in a narrower and taller curve.

3. The curve is asymptotic, meaning it approaches but never touches the horizontal axis.

4. The total area under the curve is equal to 1, representing the total probability of all possible outcomes.

The normal distribution is characterized by several important properties:

- It is unimodal, having a single peak at the mean.
- It is symmetric around the mean, with equal probabilities of values occurring on either side of the mean.
- It is continuous, meaning that the variable can take any value within a range.
- It is defined over the entire real number line, from negative infinity to positive infinity.

The normal distribution is widely used in various fields, including statistics, natural sciences, social sciences, engineering, finance, and more. It serves as a fundamental model for many natural and human-made phenomena due to its prevalence in nature and its mathematical tractability. Additionally, many statistical methods and inferential techniques are based on the assumption of normality or utilize properties of the normal distribution for their validity.

Q.11.Ans:- Handling missing data is a crucial aspect of data pre-processing in machine learning and statistical analysis. There are several strategies to deal with missing data, and the choice of technique depends on the nature of the data and the specific

requirements of the analysis. Some common approaches include:

1. Deletion:

- **Listwise deletion (complete-case analysis):** Remove entire rows of data that contain missing values. This approach is straightforward but may lead to loss of valuable information, especially if missing values are not completely random.

- **Pairwise deletion:** Analyze available data for each variable pair-wise, ignoring missing values in other variables. This approach retains more data but can lead to biased estimates if missingness is related to the variables being analyzed.

2. Imputation:

- **Mean/Median imputation:** Replace missing values with the mean or median of the observed values for that variable. This method is simple and effective but may lead to underestimation of variance and biased estimates if missingness is related to the variable being imputed.

- **Mode imputation:** Replace missing categorical values with the mode (most frequent) value of the observed values for that variable.

- **Hot-deck imputation:** Replace missing values with values from similar cases in the dataset.

- **K-nearest neighbors (KNN) imputation:** Replace missing values with the average of the nearest neighbors' values. This method considers the relationships between variables and can produce more accurate imputations.

- **Multiple imputation:** Generate multiple imputed datasets by estimating missing values multiple times, each time incorporating randomness. This approach captures uncertainty in the imputation process and produces more accurate parameter estimates.

3. Predictive modeling:

- Use machine learning algorithms to predict missing values based on observed data. This approach can capture complex relationships in the data but may be computationally expensive and require careful model selection and validation.

The choice of imputation technique depends on factors such as the proportion of missing data, the mechanism of missingness, the distribution of the data, and the analysis goals. It is often recommended to compare the performance of different imputation methods and evaluate their impact on the results of subsequent analyses. Additionally, it's essential to assess the assumptions underlying each imputation method and consider potential biases introduced by the imputation process.

Q.12.Ans:- A/B testing, also known as split testing, is a statistical method used in the field of marketing, product development, and web analytics to compare two or more versions of a webpage, app, marketing campaign, or product feature in order to determine which one performs better.

In A/B testing, the variations being tested are referred to as "A" and "B". The A version is typically the current version or the baseline, while the B version includes one or more changes or modifications that are being tested against the control.

The basic steps involved in A/B testing are as follows:

1. Hypothesis formulation: Define a hypothesis about the change or modification being tested and its

expected impact on the desired outcome (e.g., click-through rate, conversion rate, user engagement).

2. Design experiment: Randomly divide the audience or users into two or more groups, with each group exposed to a different version of the webpage, app, or feature.

3. Collect data: Track relevant metrics or key performance indicators (KPIs) for each group, such as clicks, conversions, or engagement.

4. Statistical analysis: Use statistical methods to analyze the data and determine whether any observed differences in performance between the variations are statistically significant or are likely due to chance.

5. Draw conclusions: Based on the results of the analysis, determine whether the treatment version outperforms the control version and whether the observed differences are practically significant enough to warrant implementation of the changes.

A/B testing allows businesses and organizations to make data-driven decisions and optimize their products, websites, and marketing campaigns based on empirical evidence rather than subjective opinions or intuition. It is a powerful tool for improving user experience, increasing conversion rates, and driving business growth.

Q.13.Ans:- Mean imputation, where missing values are replaced with the mean of the observed values for that variable, is a simple and commonly used method to handle missing data. However, its acceptability depends on the specific context and assumptions of

the data being analyzed. Here are some considerations regarding mean imputation:

Pros:

- 1. Simple and easy to implement:** Mean imputation is straightforward and requires minimal computational resources.
- 2. Preserves sample size:** Mean imputation allows for the retention of all observations in the dataset, unlike deletion methods that remove missing values.
- 3. Maintains variable distribution:** Mean imputation preserves the distribution of the variable in the dataset, as it only replaces missing values with a central tendency measure.

Cons:

- 1. May introduce bias:** Mean imputation assumes that missing values are missing completely at random (MCAR), meaning that the probability of missingness is unrelated to any observed or unobserved variables. If data are missing not at random (MNAR) or missing at random (MAR), mean imputation can introduce bias into the analysis.
- 2. Underestimates variability:** Mean imputation tends to underestimate the variability of the data, as it artificially reduces the variance of the imputed variable.
- 3. Distorts relationships:** Mean imputation can distort relationships between variables, particularly if missingness is related to the variable being imputed or to other variables in the dataset.
- 4. Does not reflect uncertainty:** Mean imputation provides a single imputed value for missing data, ignoring the uncertainty associated with the imputation process.

Q.14.Ans:- Linear regression is a statistical method used to model the relationship between a dependent variable and one or more independent variables. It is one of the most widely used techniques in statistics and machine learning for predictive modeling and inference.

In linear regression, the relationship between variable and the independent variables is assumed to be linear:

$$Y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \varepsilon$$

Where:

- Y is the dependent variable
- x_1, x_2, \dots, x_p are the independent variable
- $\beta_0, \beta_1, \dots, \beta_p$ are the coefficients or parameters of the model that represent the effect of each independent variable on the dependent variable.
- ε is the error term, which represents the variability in the dependent variable that is not explained by the independent variables. It is assumed to be normally distributed with mean 0 and constant variance.

The goal of linear regression is to estimate the coefficients $\beta_0, \beta_1, \dots, \beta_p$ that best fit the observed data. This is typically done by minimizing the sum of squared differences between the observed and predicted values of the dependent variable, a method known as ordinary least squares regression.

Overall, linear regression is a versatile and powerful tool for analyzing relationships between variables and making predictions based on observed data.

Q.15.Ans:- Statistics is a broad field that encompasses various branches and subfields, each focusing on different aspects of data analysis, inference, and modeling. Some of the major branches of statistics include:

1. Descriptive Statistics: Descriptive statistics involves methods for summarizing and describing the features of a dataset. This includes measures of central tendency (e.g., mean, median, mode), measures of dispersion (e.g., variance, standard deviation), and graphical techniques (e.g., histograms, box plots) for visualizing data.

2. Inferential Statistics: Inferential statistics involves making inferences and drawing conclusions about a population based on a sample of data. This includes hypothesis testing, confidence intervals, and estimation of population parameters.

3. Probability Theory: Probability theory is the mathematical framework for quantifying uncertainty and randomness. It includes concepts such as probability distributions, random variables, independence, and conditional probability.

4. Biostatistics: Biostatistics is the application of statistical methods to biological and health-related data. It includes areas such as clinical trials, epidemiology, medical research, and public health.

5. Econometrics: Econometrics is the application of statistical methods to economic data. It includes areas such as regression analysis, time series analysis, and causal inference in economics.

6. Bayesian Statistics: Bayesian statistics is an approach to statistical inference that uses Bayesian probability to update beliefs about parameters and make predictions based on data and prior knowledge. It includes methods such as Bayesian inference, Bayesian modeling, and Markov Chain Monte Carlo (MCMC) sampling.

7. Multivariate Statistics: Multivariate statistics involves the analysis of datasets with multiple variables. It includes techniques such as multivariate regression, factor analysis, principal component analysis, and cluster analysis.

8. Nonparametric Statistics: Nonparametric statistics involves methods that do not rely on specific assumptions about the underlying distribution of the data. This includes rank-based tests, kernel density estimation, and resampling methods such as bootstrapping.

9. Spatial Statistics: Spatial statistics is the analysis of spatial data, including methods for modeling spatial patterns, spatial autocorrelation, and spatial interpolation.

10. Statistical Learning: Statistical learning is the field that combines statistics and machine learning techniques to develop algorithms for prediction and classification. It includes methods such as linear regression, logistic regression, decision trees, support vector machines, and neural networks.

These are just a few examples of the many branches and subfields within the field of statistics. Each branch has its own techniques, methodologies, and applications, but they are all interconnected and

contribute to our understanding of data and uncertainty.