

Statistics Worksheet-1

1.The correct answer is: a) True

Bernoulli random variables are a type of discrete random variable that can only take on two values: 0 and 1. They are often used to model binary outcomes, such as success or failure, yes or no, or 0 or 1.

2.The correct answer is: a) Central Limit Theorem

The Central Limit Theorem (CLT) states that, given certain conditions, the distribution of the average of a large number of independent and identically distributed (iid) random variables will be approximately normal, even if the underlying distribution is not normal. The CLT is a fundamental concept in probability theory and is widely used in statistical analysis.

3.The correct answer is: b) Modeling bounded count data

4.The correct answer is: d) All of the mentioned

Here's why:

a) The exponent of a normally distributed random variable follows what is called the log-normal distribution: This is true. If X is a normally distributed random variable, then e^X follows a log-normal distribution.

b) Sums of normally distributed random variables are again normally distributed even if the variables are dependent: This is also true, but with a caveat. If the dependent normally distributed random variables are jointly normally distributed, then their sum is also normally distributed. However, if they are not jointly normally distributed, then the sum may not be normally distributed.

c) The square of a standard normal random variable follows what is called chi-squared distribution: This is true. If X is a standard normal random variable (i.e., $X \sim N(0, 1)$), then X^2 follows a chi-squared distribution with 1 degree of freedom.

5.The correct answer is: c) Poisson

Poisson random variables are commonly used to model rates, such as:

The number of defects per unit of a product

The number of phone calls received per minute

The number of accidents per year

The number of customers arriving per hour

6.The correct answer is: b) False

7.The correct answer is: b) Hypothesis

Hypothesis testing is a statistical technique used to make decisions based on data. It involves formulating a hypothesis about a population parameter, collecting sample data, and using statistical methods to determine whether the data provide sufficient evidence to support or reject the hypothesis.

8.The correct answer is: a) 0

9.The correct answer is d) None of the mentioned.

Here's why:

a) Outliers can have varying degrees of influence: This is correct.

Outliers can have a significant impact on statistical models, and their influence can vary depending on the type of model, the size of the dataset, and the nature of the outlier.

b) Outliers can be the result of spurious or real processes: This is also correct. Outliers can be the result of errors in data collection, measurement errors, or other spurious processes. On the other hand, they can also be the result of real processes or phenomena that are not well-represented by the majority of the data.

c) Outliers cannot conform to the regression relationship: This is correct as well. By definition, outliers do not conform to the expected pattern or relationship in the data, including regression relationships. Therefore, all the options are correct, and the correct answer is d) None of the mentioned.

10.The Normal Distribution, also known as the Gaussian Distribution or Bell Curve, is a continuous probability distribution that is widely used in statistics and data analysis.

The normal distribution is commonly used to model real-valued random variables that are expected to be distributed symmetrically around the mean, such as:

Human heights and weights

IQ scores

Errors in measurement

Stock prices

and many more!

The normal distribution has many important properties and is used extensively in statistical inference, hypothesis testing, and confidence intervals.

11.Data Dropping: One approach is to simply drop the rows with missing data, also known as listwise deletion. However, this method can lead to biased parameter estimates and underestimated standard errors, especially if the data is not Missing Completely at Random (MCAR).

Single Imputation: Another approach is single imputation, where a single estimate of the missing value is used. This can be done using various methods, including:

Mean/Median Imputation: Replace the missing value with the mean or median of the observed values for that variable.

Random Sample Imputation: Replace the missing value with a random sample from the observed values for that variable.

Hot Deck Imputation: Replace the missing value with a randomly chosen value from an individual in the sample who has similar values on other variables.

Substitution: Impute the value from a new individual who was not selected to be in the sample.

12.A/B testing, also known as split testing, is a method of comparing two or more versions of a product, web page, or application to determine which one performs better in terms of achieving a specific goal or

metric. The goal of A/B testing is to identify changes that can improve user engagement, conversion rates, revenue, or other desired outcomes.

13. Mean imputation of missing data is a common practice, but it's not always the most acceptable or recommended approach. Here's why:

Mean imputation involves replacing missing values with the mean value of the respective feature or variable. For example, if you have a dataset with missing values in a column, you would replace those missing values with the average value of that column.

In summary, while mean imputation is a common practice, it's not always the best approach. It's essential to consider the limitations and potential biases of mean imputation and explore alternative methods that better account for the complexity of the data.

14. Linear regression is a fundamental concept in statistics that models the relationship between a dependent variable (also called the outcome or response variable) and one or more independent variables (also called predictor or feature variables). The goal of linear regression is to create a linear equation that best predicts the value of the dependent variable based on the values of the independent variables.

15. Statistics is a diverse field that encompasses various branches, each focusing on a specific aspect of data analysis and interpretation. Here are some of the main branches of statistics:

1. Descriptive Statistics

Descriptive statistics deals with summarizing and describing the basic features of a dataset, such as measures of central tendency (mean, median, mode) and variability (range, variance, standard deviation).

2. Inferential Statistics

Inferential statistics involves making conclusions or inferences about a population based on a sample of data. This branch includes hypothesis testing, confidence intervals, and significance testing.

3. Exploratory Data Analysis (EDA)

EDA is an approach to analyzing data to understand its underlying structure and patterns. It involves using various techniques, such as data visualization, to identify relationships and anomalies in the data.

4. Machine Learning

Machine learning is a subfield of artificial intelligence that involves developing algorithms and models that can learn from data and make predictions or decisions. It includes supervised, unsupervised, and reinforcement learning.

5. Bayesian Statistics

Bayesian statistics is a branch of statistics that uses Bayes' theorem to update the probability of a hypothesis based on new data. It provides a framework for quantifying uncertainty and making decisions under uncertainty.

6. Time Series Analysis

Time series analysis deals with analyzing and forecasting data that is collected over time. It involves techniques such as autoregressive integrated moving average (ARIMA) models and exponential smoothing.

7. Spatial Statistics

Spatial statistics is concerned with analyzing and modeling data that is associated with geographic locations. It involves techniques such as spatial autocorrelation and geostatistics.

8. Biostatistics

Biostatistics is the application of statistical principles to medical and health-related data. It involves designing experiments, analyzing data, and interpreting results in fields such as epidemiology, clinical trials, and public health.

9. Computational Statistics

Computational statistics involves developing algorithms and software for statistical analysis and modeling. It includes areas such as numerical analysis, optimization, and simulation.

10. Mathematical Statistics

Mathematical statistics is concerned with the theoretical foundations of statistics, including probability theory, stochastic processes, and mathematical modeling.

These branches are not mutually exclusive, and many statistical techniques and methods overlap across multiple branches.