

**STATISTICS WORKSHEET-1**

**Q1 to Q9 have only one correct answer. Choose the correct option to answer your question.**

1. Bernoulli random variables take (only) the values 1 and 0.
  - a) True
  - b) FalseAns a
2. Which of the following theorem states that the distribution of averages of iid variables, properly normalized, becomes that of a standard normal as the sample size increases?
  - a) Central Limit Theorem
  - b) Central Mean Theorem
  - c) Centroid Limit Theorem
  - d) All of the mentionedAns d
3. Which of the following is incorrect with respect to use of Poisson distribution?
  - a) Modeling event/time data
  - b) Modeling bounded count data
  - c) Modeling contingency tables
  - d) All of the mentionedAns B
4. Point out the correct statement.
  - a) The exponent of a normally distributed random variables follows what is called the log- normal distribution
  - b) Sums of normally distributed random variables are again normally distributed even if the variables are dependent
  - c) The square of a standard normal random variable follows what is called chi-squared distribution
  - d) All of the mentionedAns D
5. \_\_\_\_\_ random variables are used to model rates.
  - a) Empirical
  - b) Binomial
  - c) Poisson
  - d) All of the mentionedAns C
6. 10. Usually replacing the standard error by its estimated value does change the CLT.
  - a) True
  - b) FalseAns B
7. 1. Which of the following testing is concerned with making decisions using data?
  - a) Probability
  - b) Hypothesis
  - c) Causal
  - d) None of the mentionedANS b
8. 4. Normalized data are centered at \_\_\_\_\_ and have units equal to standard deviations of the original data.
  - a) 0
  - b) 5
  - c) 1
  - d) 10Ans a
9. Which of the following statement is incorrect with respect to outliers?
  - a) Outliers can have varying degrees of influence

- b) Outliers can be the result of spurious or real processes
- c) Outliers cannot conform to the regression relationship
- d) None of the mentioned

Ans c

**Q10 and Q15 are subjective answer type questions, Answer them in your own words briefly.**

10. What do you understand by the term Normal Distribution?

- Ans Normal distribution, also known as the Gaussian distribution, is a probability distribution that is symmetric about the mean, showing that data near the mean are more frequent in occurrence than data far from the mean. The normal distribution is the proper term for a probability bell curve. In a normal distribution the mean is zero and the standard deviation is 1. It has zero skew and a kurtosis of 3

11. How do you handle missing data? What imputation techniques do you recommend?

In real life, many datasets will have many missing values, so dealing with them is an important step.

Why do you need to fill in the missing data? Because most of the machine learning models that you want to use will provide an error if you pass NaN values into it. The easiest way is to just fill them up with 0, but this can reduce your model accuracy significantly.

1. Deleting the columns with missing data
2. Deleting the rows with missing data
3. Filling the missing data with a value – Imputation
4. Imputation with an additional column
5. Filling with a Regression Model

12. What is A/B testing?

A/B testing is a type of experiment in which you split your web traffic or user base into two groups, and show two different versions of a web page, app, email, and so on, with the goal of comparing the results to find the more successful version.

13. Is mean imputation of missing data acceptable practice?

the dataset we want to use for Machine Learning contains missing data. The quick and easy workaround is to substitute a mean for numerical features and use a mode for categorical ones. Even better, someone might just insert 0's or discard the data and proceed to the training of the model. In the following article, I will explain why using a mean or mode can significantly reduce the model's accuracy and bias the results. I will also point you to few alternative imputation algorithms which have their respective Python libraries that you can use out-of-the-box

We can use some of other algorithms in place of mean and median Like KNN, Random forest, fuzzy kmean clustering, python imputation library.

14. What is linear regression in statistics?

Linear regression analysis is used to predict the value of a variable based on the value of another variable. The variable you want to predict is called the dependent variable. The variable you are using to predict the other variable's value is called the independent variable.

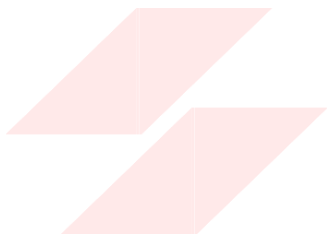
15. What are the various branches of statistics?

The two major areas of statistics are descriptive and inferential statistics.

Descriptive statistics mostly focus on the central tendency, variability, and distribution of sample data. Central tendency means the estimate of the characteristics, a typical element of a sample or population, and includes descriptive statistics such as mean, median, and mode. Variability refers to a set of statistics that show how much difference there is among the elements of a sample or population along the characteristics measured, and includes metrics such as range, variance, and standard deviation.

Inferential statistics are tools that statisticians use to draw conclusions about the characteristics of a population, drawn from the characteristics of a sample, and to decide how certain they can be of the reliability of those conclusions. Based on the sample size and distribution statisticians can calculate the probability that statistics, which measure the central tendency, variability, distribution, and relationships between characteristics within a data sample, provide an accurate picture of the corresponding parameters of the whole population from which the sample is drawn.

---



# FLIP ROBO