

## **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) False

Ans: a) True

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 mentioned

Ans: a) Central Limit Theorem

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 mentioned

Ans: b) Modeling bounded count data

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 mentioned

Ans: d) All of the mentioned

5. \_\_\_\_\_ random variables are used to model rates.

- a) Empirical  
b) Binomial  
c) Poisson  
d) All of the mentioned

Ans: c) Poisson

6. Usually replacing the standard error by its estimated value does change the CLT.

- a) True                      b) False

Ans: b) False

7. Which of the following testing is concerned with making decisions using data?

- a) Probability  
b) Hypothesis  
c) Causal  
d) None of the mentioned

Ans: b) Hypothesis

8. Normalized data are centered at \_\_\_\_\_ and have units equal to standard deviations of the original data.

- a) 0    b) 5    c) 1    d) 10

Ans: a) 0

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) Outliers can have varying degrees of influence

**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. In graph form, normal distribution will appear as a bell curve.

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

Ans: There are 7 ways to handle missing values in the dataset:

- Deleting Rows with missing values
- Impute missing values for continuous variable
- Impute missing values for categorical variable
- Other Imputation Methods
- Using Algorithms that support missing values
- Prediction of missing values
- Imputation using Deep Learning Library

We have to see the feature of missing value first then depending on the dependency of it on the outcome and on the other features it should be handled.

I will recommend the simplest imputation method is **replacing missing values** with the mean or median values of the dataset at large, or some similar summary statistic. This has the advantage of being the simplest possible approach, and one that doesn't introduce any undue bias into the dataset.

**12. What is A/B testing?**

Ans: A/B testing is a basic randomized control experiment. It is a way to compare the two versions of a variable to find out which performs better in a controlled environment.

For instance, let's say you own a company and want to increase the sales of your product. Here, either you can use random experiments, or you can apply scientific and statistical methods. A/B testing is one of the most prominent and widely used statistical tools.

In the above scenario, you may divide the products into two parts – A and B. Here A will remain unchanged while you make significant changes in B's packaging. Now, on the basis of the response from customer groups who used A and B respectively, you try to decide which is performing better.

It is a hypothetical testing methodology for making decisions that estimate population parameters based on sample statistics. The **population** refers to all the customers buying your product, while the **sample** refers to the number of customers that participated in the test.

**13. Is mean imputation of missing data acceptable practice?**

Ans: Not always, but it totally depends on the type of dataset and its dependency on the null values that how it will be affecting the outcomes.

**14. What is linear regression in statistics?**

Ans: 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.

This form of analysis estimates the coefficients of the linear equation, involving one or more independent variables that best predict the value of the dependent variable. Linear regression fits a straight line or surface that minimizes the discrepancies between predicted and actual output values. There are simple linear regression calculators that use a "least squares" method to discover the best-fit line for a set of paired data. You then estimate the value of X (dependent variable) from Y (independent variable).

**15. What are the various branches of statistics?**

Ans: There are mainly two types/branches of statistics: 1. Descriptive & 2. Inferential

**1.Descriptive statistics:** if data can be described without any statistical tools, then it is called descriptive statistics. ex, marks in class, height of student.

**2.Inferential statistics:** if data is too big then then we use inferential statistics.

