

## **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

Answer – 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 mentioned

Answer - A

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

Answer – 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 mentioned

Answer – C

5. \_\_\_\_\_ random variables are used to model rates.
- a) Empirical
  - b) Binomial
  - c) Poisson
  - d) All of the mentioned

Answer – C

6. 10. Usually replacing the standard error by its estimated value does change the CLT.
- a) True
  - b) False

Answer – 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 mentioned

Answer - 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) 10

Answer – 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

Answer - 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?

Normal Distribution is basically a probability distribution which is symmetric about the mean and shows data near the mean are more frequent in occurrence than data away from mean. It is also known as bell curve in graph.

The standard normal distribution has two parameters: the mean and the standard deviation. For a normal distribution, 68% of the observations are within  $\pm$  one standard deviation of the mean, 95% are within  $\pm$  two standard deviations, and 99.7% are within  $\pm$  three standard deviations.

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

The best way to handle missing data is Imputation technique which helps to fill the NAN data and remove the skewedness of the graph. There are many techniques

1. Mean or Median Imputation
2. Random Forest

12. What is A/B testing?

A/B testing or bucket testing, is a user experience research methodology. A/B testing is a shorthand for a simple randomized controlled experiment, in which two samples (A and B) of a single vector-variable are compared. These values are similar except for one variation which might affect a user's behavior. A/B tests are widely considered the simplest form of controlled experiment. However, by adding more variants to the test, its complexity grows.

13. Is mean imputation of missing data acceptable practice?

Mean imputation is not acceptable practice and is taken as last resource because of the reasons below:

- Mean imputation reduces the variance of the imputed variables.
- Mean imputation shrinks standard errors, which invalidates most hypothesis tests and the calculation of confidence interval.
- Mean imputation does not preserve relationships between variables such as correlations.

14. What is linear regression in statistics?

Linear regression is used to estimate the relationship between two or more quantitative variables. You can use linear regression when you want to know:

1. How strong the relationship is between two or more variables
2. The value of the dependent variable at a certain value of the independent variable (e.g. the amount of soil erosion at a certain level of rainfall).

15. What are the various branches of statistics?

There are 3 branches of Statistics:

1. Data Collection

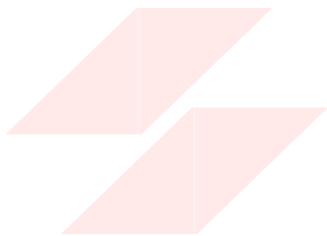
Data collection is all about how the actual data is collected. For the most part, this needn't concern us too much in terms of the mathematics (we just work with what we are given), but there are significant issues to consider when actually collecting data.

2. Descriptive Statistics

Descriptive statistics is the part of statistics that deals with presenting the data we have. This can take two basic forms – presenting aspects of the data either visually (via graphs, charts, etc.) or numerically (via averages and so on).

3. Inferential Statistics

Inferential statistics is the aspect that deals with making conclusions about the data. This is quite a wide area; essentially you are asking 'What is this data telling us, and what should we do?'



**FLIP ROBO**

---