

## 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
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
3. Which of the following is incorrect with respect to use of Poisson distribution?
  - b) Modeling bounded count data
4. Point out the correct statement.
  - d) All of the mentioned
5. \_\_\_\_\_ random variables are used to model rates
  - c) Poisson
6. 10. Usually replacing the standard error by its estimated value does change the CLT.
  - b) False
7. 1. Which of the following testing is concerned with making decisions using data?
  - b) Hypothesis
8. 4. Normalized data are centered at \_\_\_\_\_ and have units equal to standard deviations of the original data.
  - a) 0
9. Which of the following statement is incorrect with respect to outliers?
  - c) Outliers cannot conform to the regression relationship

**Q10and 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, 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?

**Use deletion methods to eliminate missing data**

The deletion methods only work for certain datasets where participants have missing fields. There are several deleting methods – two common ones include Listwise Deletion and Pairwise Deletion. It means deleting any participants or data entries with missing values. This method is particularly advantageous

to samples where there is a large volume of data because values can be deleted without significantly distorting readings.

### **Use regression analysis to systematically eliminate data**

Regression is useful for handling missing data because it can be used to predict the null value using other information from the dataset. There are several methods of regression analysis, like Stochastic regression. Regression methods can be successful in finding the missing data, but this largely depends on how well connected the remaining data is. Of course, the one drawback with regression analysis is that it requires significant computing power, which could be a problem if data scientists are dealing with a large dataset.

### **Data scientists can use data imputation techniques**

Data scientists use two data imputation techniques to handle missing data: Average imputation and common-point imputation. Average imputation uses the average value of the responses from other data entries to fill out missing values. However, a word of caution when using this method – it can artificially reduce the variability of the dataset.

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

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.

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

True, imputing the mean preserves the mean of the observed data. So if the data are missing completely at random, the estimate of the mean remains unbiased.

#### **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?**

There are three real branches of statistics: data collection, descriptive statistics and inferential statistics.