# 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**

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

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

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**

5. _____ random variables are used to model rates.

a) Empirical

b) Binomial

**c) Poisson**

d) All of the mentioned

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

a) True

**b) False**

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

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

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 WORKSHEET

*Q10and 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 graphical form, the normal distribution appears as a "bell curve".

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

**Ans**- Missing data appear when no value is available in one or more variables of an individual. Due to Missing data, the statistical power of the analysis can reduce, which can impact the validity of the results.

**Zero Replacement**: Here, you replace the missing value with zero irrespective of everything.

**Min or Max Replacement**: Replace the missing value with the minimum or maximum value of a feature.

**Mean/ Median/ Mode Replacement**: Replace missing value with mean or median or most frequent feature value.

12. What is A/B testing?

**Ans** - A/B testing, also known as split testing, refers to a randomized experimentation process wherein two or more versions of a variable (web page, page element, etc.) are shown to different segments of website visitors at the same time to determine which version leaves the maximum impact and drives business metrics.

13. Is mean imputation of missing data acceptable practice?

**Ans** - Mean imputation is typically considered terrible practice since it ignores feature correlation. Consider the following scenario: we have a table with age and fitness scores, and an eight-year-old has a missing fitness score. If we average the fitness scores of people between the ages of 15 and 80, the eighty-year-old will appear to have a significantly greater fitness level than he actually does.

Second, mean imputation decreases the variance of our data while increasing bias. As a result of the reduced variance, the model is less accurate and the confidence interval is narrower.

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.

15. What are the various branches of statistics?

**Ans** - Statistics is a study of presentation, analysis, collection, interpretation and organization of data

There are two main branches of statistics
- Inferential Statistic.
- Descriptive Statistic.

Inferential Statistics:
Inferential statistics used to make inference and describe about the population. These stats are more useful when its not easy or possible to examine each member of the population.

Descriptive Statistics:
Descriptive statistics are use to get a brief summary of data. You can have the summary of data in numerical or graphycal form.