# STATISTICS WORKSHEET-1-----Answers are mentioned in green

**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. Usually replacing the standard error by its estimated value does change the CLT.
a) True
**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

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

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

**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 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. It appears like a bell curve in its visual effect as where major distribution of data is around centre

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

Ans: We can handle missing data by replacement with Mode, Median, Forward or backward fill or any arbitrary values

Imputation techniques usually recommended are Regression, Interpolation and extrapolation, Substitution, Mean imputation among many others.

12. What is A/B testing?

Ans. A/B testing is a kind of test where two data variants are checked to see their performance based on a given metric or value. One group which performs better based on that metric is usually chosen and preferred over the other

13. Is mean imputation of missing data acceptable practice?

Ans: Due to its ignoring feature co-relation, it is not an acceptable practice in most of the cases

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?

Statistics have majorly categorised into two types:
1. Descriptive statistics
2. Inferential statistics

<u>Descriptive Statistics</u>
In this type of statistics, the data is summarised through the given observations. The summarisation is one from a sample of population using parameters such as the mean or standard deviation.
Descriptive statistics is a way to organise, represent and describe a collection of data using tables, graphs, and summary measures. For example, the collection of people in a city using the internet or using Television.
Descriptive statistics are also categorised into four different categories:
- Measure of frequency
- Measure of dispersion
- Measure of central tendency
- Measure of position

The frequency measurement displays the number of times a particular data occurs. Range, Variance, Standard Deviation are measures of dispersion. It identifies the spread of data. Central tendencies are the mean, median and mode of the data. And the measure of position describes the percentile and quartile ranks.

<u>Inferential Statistics</u>
This type of statistics is used to interpret the meaning of Descriptive statistics. That means once the data has been collected, analysed and summarised then we use these stats to describe the meaning of the collected data. Or we

can say, it is used to draw conclusions from the data that depends on random variations such as observational errors, sampling variation, etc.

Inferential Statistics is a method that allows us to use information collected from a sample to make decisions, predictions or inferences from a population. It grants us permission to give statements that goes beyond the available data or information. For example, deriving estimates from hypothetical research.