# STATISTICS WORKSHEET-1 Sayed Raza Ali Hussain

**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, 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. Imputing the Missing Value

- Replacing With Arbitrary Value.
- Replacing With Mode.
- Replacing With Median.
- Replacing with previous value – Forward fill.
- Replacing with next value – Backward fill.
- Impute the Most Frequent Value.

  Common Imputation techniques are
- Mean imputation. Simply calculate the mean of the observed values for that variable for all individuals who are non-missing
- Substitution
- Hot deck imputation
- Cold deck imputation.
- Regression imputation.
- Stochastic regression imputation.
- Interpolation and extrapolation.

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

Ans. A/B testing in its simplest sense is an experiment on two variants to see which performs better based on a given metric. Typically, two consumer groups are exposed to two different versions of the same thing to see if there is a significant difference in metrics like sessions, click-through rate, and/or conversions

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

Ans. Mean imputation is typically considered terrible practice since it ignores feature correlation

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

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

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).

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.