

# Statistics Worksheet 1 Solutions

1. Option a -> True
2. Option a-> Central Limit Theorem
3. Option b -> Modelling Bounded Count Data
4. Option d -> All of the mentioned
5. Option c -> Poisson
6. Option b -> False
7. Option b -> Hypothesis
8. Option a -> 0
9. Option c -> Outliers Cannot conform to the regression relationship
10. Normal distribution is a distribution in which the mean is zero and the standard deviation is 1. It forms a bell shaped curve. Normal distributions are symmetric. For a normal distribution, 68% of the observations are within  $\pm$  one standard deviation of the mean, 95% are within  $\pm$  two standard deviations of the mean, and 99.7% are within three standard deviations of the mean. Real life data rarely follows a perfect normal distribution. Normal Distribution has zero skewness. The normal distribution has a kurtosis of three which indicates it has neither fat tails nor thin tails.
11. Ways of handling missing data->
  - a. Deleting the entire row in which there is a cell with missing data
  - b. Imputing The value->
    - i. Replace all the missing data with an arbitrary value
    - ii. Replace missing data with the mean of the column
    - iii. Replace missing data with the mode of the column
    - iv. Replace missing data with the median of the column

- v. Replace missing data with the previous value in the column(Forward Fill)
  - vi. Replace missing data with the next value in the column(backward fill)
  - vii. Interpolation
12. A/B testing, also known as split testing, refers to a randomized process wherein two or more versions of a variable are shown to different segments of website visitors at the same time to determine which version leaves the maximum impact and derives the business metric.
13. Mean imputaion is typically considered a terrible practice since it ignores feature correlation. Consider the following scenario: we have a table with age and income, now suppose 21 year old has a missing value of the income, if we average the income of all people aged between 20 and 60, the 20 year old has a higher income than what his actual income is.
- Mean imputation also decreases the variance of our data while increasing our bias.
14. Linear regression is used to predict the value of a variable on the value of another variable. The variable we want to predict is called the dependent variable and the variable we use to predict is called the independent variable. The goal of linear regression is to predict the slope and the Y-intercept of a line that predicts the relation between predictor variable and prediction variable.

$$Y = B_0 + B_1 * X$$

Where  $B_0$  is the Y-intercept of the line and  $B_1$  is the slope

The above equation is the one we get from linear regression. The actual relation also has an error term which is independent of  $X$  and thus linear regression is unable to predict that error term

Actual Relation =

$$y = b_0 + b_1x + \text{error term}$$

The error term is generally irreducible error.

15. There are two branches of statistics ->
- a. Descriptive Statistics -> It deals with the presentation and collection of data. This is usually the first part of the statistical analysis. In descriptive statistics we simply state what the data tells us.
  - b. Inferential Statistics - > It involves drawing the right conclusion from the statistical analysis that has been performed using descriptive statistics. It attempts to apply the conclusions obtained from the experiments to more general populations. This tries to answer questions about populations and samples that have not been tested in the given experiment.