## Statistics - 1   Worksheet

By:
- Yatika Taneja

① Bernoulli random variables variables take only the values 1 or 0

Sol) True (a)

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

Sol) (a) Central limit theorem

③ Which of the following is incorrect with respect to use of poisson distribution

Sol) (b) Modeling bounded count data

④ Point out the correct statement

Sol) (d) All of the above mentioned

⑤ _____ Random variables are used to model rates

Sol) (c) Poisson distribution

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

Sol> (b) false

**7.** Which of the following testing is concerned with making decisions using data.

Sol> (b) Hypothesis

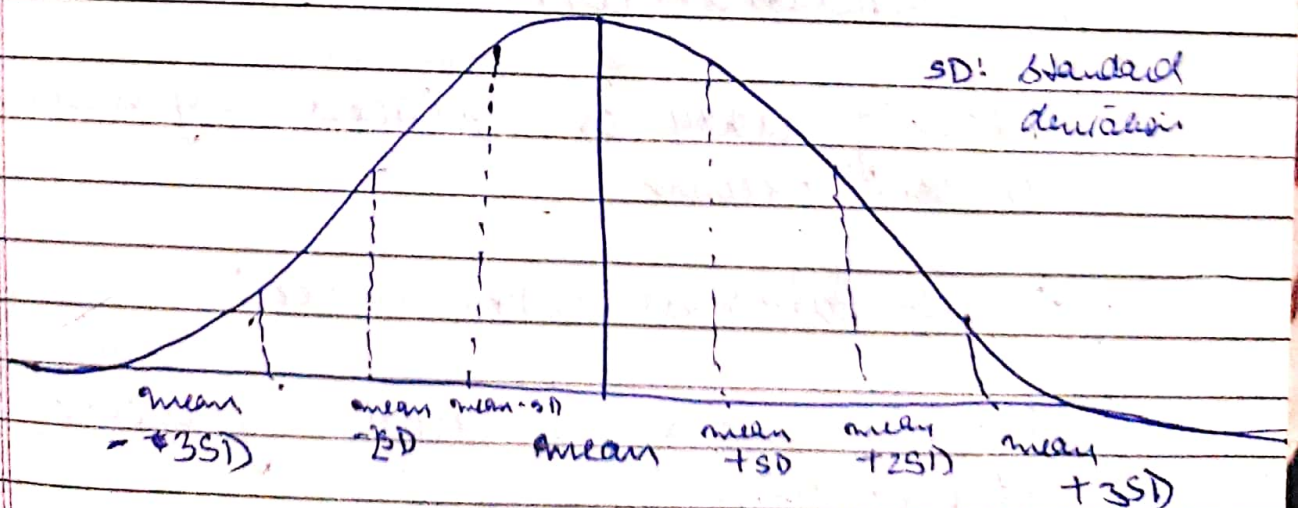**8.** Normalized data are centered at 0 and have units equal to standard deviations of the original data

Sol> (a) 0

**9.** Which of the following statement is incorrect with respect to outliers

Sol> (c) Outliers cannot conform to the regression relationship

**10.** What do you understand by the term Normal distribution?

- Normal distribution is a probability distribution that is symmetric about the mean

- The standard normal distribution has 2 parameters: mean & standard deviation

- The empirical rule states that:
  (i) 68% of values fall within 1 standard deviation of the mean
  (ii) 95% of values fall 2 standard deviation of the mean
  (iii) 99.7% of values are within 3 standard deviation of the mean

- In graph form, it appears as bell curve

SD: Standard deviation

mean
- #3SD

mean
-3D

mean-SD

mean

mean
+SD

mean
+2SD

mean
+3SD

(11) How do you handle missing data? What imputation ~~data~~ techniques do you recommend

① Remove all rows with missing data if there are not too many rows with missing data

② If more than 50% of rows of a specific column have missing data it is common to remove the particular column

③ Imputation Techniques :-

ⓐ Imputation with mean:

Missing data is replaced by mean of the column

ⓑ Imputation with median:

Missing data is replaced by median of the column

ⓒ Imputation with mode:

Missing data is replaced by mode of the column

(A) Imputation with Linear Regression:

The missing value is replaced by applying linear regression based on other feature values

(12) What is A/B testing?

- It is a hypothetical testing methodology for making decisions that estimate population parameters based on sample statistics

- A/B testing is a way to compare two versions of a variable to find out which performs better in a controlled environment

(13) Is mean imputation of missing data acceptable practice?

- Mean imputation should only be used for small numerical datasets

- For larger datasets, it is not suitable because:

(1) It doesn't factor the correlation between

features. It only works on the column
level

(ii) does Gives poor result on encoded
categorical features

(iii) Not accurate and does not account
for uncertainicity in the imputations

(14) What is linear regression in
Statistics

- Linear regression is a linear model
i.e a model that assumes linear
relationship between input variables
(n) and single output variable (y)

- Linear regression can be expressed
as:
$$y = B_0 + B_1^* n$$

where $B_0$ is the intercept & $B_1$ is
the coefficient

- In higher dimensions when we have
more than one input(n), the line
is called a plane or a higher-plane

(15) What are various branches of statistics?

Sol) The 2 branches of statistics are :-

I   Descriptive Statistics

- Analysis of data that helps discribe, show or summarize data in a meaningful way.

- It does not allow us to make conclusions beyond the data we have analyzed

- 2 types of descriptive statistics are :-

  ⇒ Measure of central tendency.
    · mean
    · median
    · mode

  ⇒ Measure of spread.
    · range
    · quartiles
    · variance
    · standard deviation

II   Inferential Statistics

- Inferential statistics are often used to compare the differences between treatment

groups.

Inferential statistics use measurement from sample of subjects in the experiment to compare the treatment groups and make generalizations about the layer population of subjects

## Statistics - 2 Worksheet

By:
Yatika Taneja

① What represent a Population parameter?

Sol) (B) mean

② What will be the median of following set of scores (18, 6, 12, 10, 15) ?

Sol) (C) 12

③ What is standard deviation?

Sol) (D) All of the above

(4) The intervals should be — in a grouped frequency distribution

Sol) c) Both of these

(5) what is the goal of descriptive of statistics?

Sol) (B) Summarizing & explaining specific set of data

(6) A set of data organized in a Participant by variables format is called

Sol) (B) Data Set

(7) In multiple regression, dependent variables are used

Sol) (c) 1

(8) which of the following is is used when you want to visually examine the relationship between 2 quantitative variables?

Sol) B) Scatter Plot

**⑨** 2 or more group means are compared by using

Sol) D) Analysis of variance

**⑩** ___ is a raw score which has been transformed into standard deviation units.

Sol) (a) z-score

**⑪** ___ is the value calculated when you want arithmetic average?

Sol) c) mean

**⑫** find the mean of these set of numbers (4, 6, 7, 9, 2000 000).

Sol) D) 4000 05 02

**⑬** ___ is a measure of central tendency that takes into account magnitude of scores

Sol) D) Mean

(14) ___ Is focus on describing or explaining data whereas ___ involves going beyond immediate & making inferences

Sol) (A) Descriptive & inferences

(15) What is the formula of range

Sol) (D) H-L