

Q.1)

Write a short note on:-

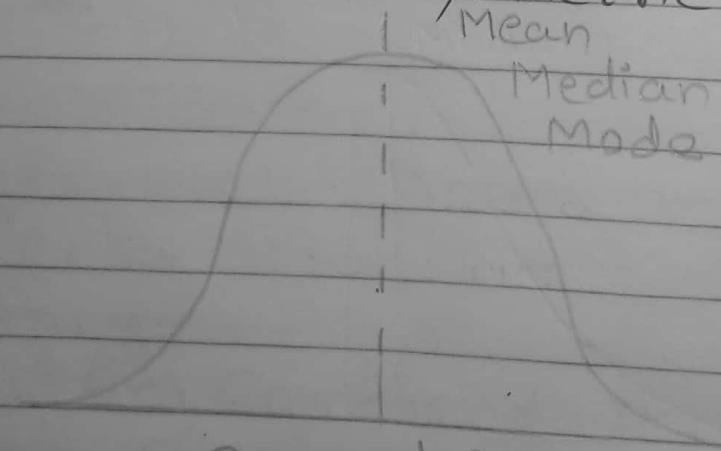
Q.1) a) Skewness and Kurtosis

→ i) The extent to which a distribution of data points is concentrated at one end or the other is known as 'Skewness'.

ii) The degree of peakedness of a distribution of points is called as 'Kurtosis'.

iii) Based on the three Central Tendency measures we can decide the Skewness of the distribution.

a) If $\text{mean} = \text{median} = \text{mode}$, then we can ensure that the shape of the distribution is Symmetric.



Symmetrical
Distribution.

b) If $\text{mode} < \text{median} < \text{mean}$, then we can ensure that the shape of the distribution trails to the right, is positively skewed.

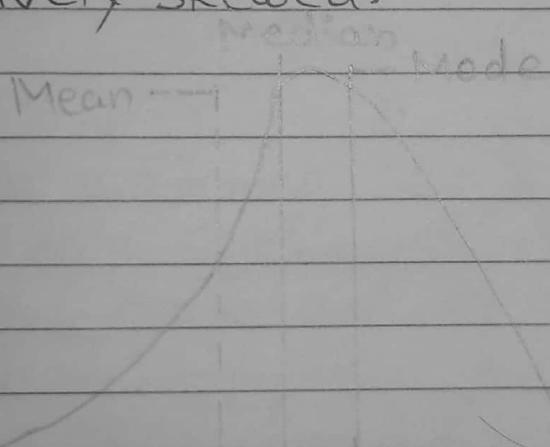
Median
Mode

Mean



Positive
Skew

- c) If $\text{mean} < \text{median} < \text{mode}$, the shape of the distribution trails to the left, is negatively skewed.



Negative
Skew

- iv) Kurtosis is measured Based on the Moments

Formula's are as follows:-

$$\beta_1 = \frac{\mu_3^2}{\mu_2^3}$$

$$\beta_2 = \frac{\mu_4}{\mu_2^2}$$

$\gamma_1 = \sqrt{\beta_1}$
(Gamma₁)
Coefficient
of Skewness

$\gamma_2 = \beta_2 - 3$
(Gamma₂)
Coefficient of
Kurtosis

If,

Decision

If,

a) $B_2 = 3$

Mesokurtic

$\gamma_2 = 0$

b) $B_2 < 3$

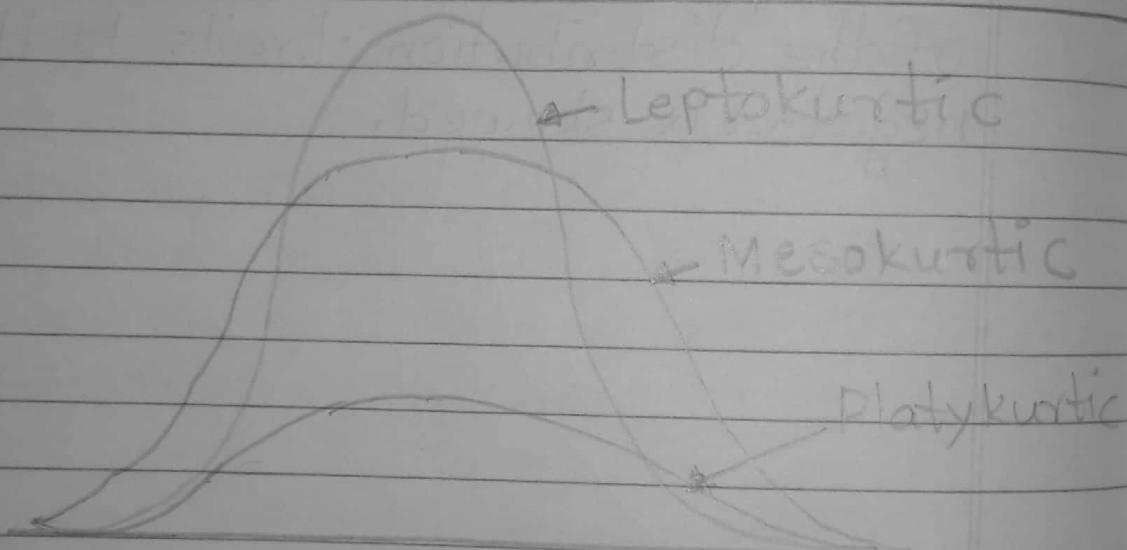
Platykurtic

$\gamma_2 < 0$

c) $B_2 > 3$

Leptokurtic

$\gamma_2 > 0$



Kurtosis

- v) Third Central Moment gives us skewness.
vi) Whereas, Fourth central Moment gives us Kurtosis.

vii) Two more popular Measures of Skewness are

- a) Karl Pearson's Coefficient of skewness
- b) Bowley's Coefficient of Skewness
- a) Karl Pearson's Coefficient of Skewness

Formula:-

$$S = \frac{\text{mean-mode}}{\text{S.D}} = \frac{\bar{x} - Z}{\sigma}$$

(Standard Deviation)

Note:- If mode is ill-defined, we can go for following formula;

$$S = \frac{3(\text{mean-median})}{\text{S.D}}$$

(Standard Deviation)

$$= \frac{3(\bar{x} - M)}{\sigma}$$

Since we know, $Z = 3M - 2\bar{x}$

b) Bowley's Coefficient of Skewness:
 Bowley's Coefficient of skewness is based on quartiles. Also called Galton skewness.

Formula:-

$$S = \frac{(Q_3 - Q_2) - (Q_2 - Q_1)}{(Q_3 - Q_2) + (Q_2 - Q_1)}$$

Q.1 b) Type I and Type II Error

→ Statisticians use specific definitions and symbols for signifying the error in the hypothesis.

→ Rejecting a null hypothesis when it is true is called a Type I error, and its probability (which, as we have seen, is also the significance level of the test) is symbolized α (Alpha.)

- Alternatively, accepting a null hypothesis, when it is false is called a Type II Error and its probability is symbolized as β (Beta)
- There is a trade off between these two Errors:-
 - The Probability of making one type of error can be reduced only if we are willing to increase the probability of making the other type of error.
- To deal with this trade off, in personal and professional situations, decision makers decide the appropriate level of significance by examining the costs or penalties attached to both type of errors.
- If we do not want to occur Type I-Error as well as Type -II Error we need to increase the number of Samples, -

Q.1) c) P-Value

- The Computed probability of getting the observed result or any result at least as extreme in its difference from what the null hypothesis would imply-is called the p-value.
- A p-value of 0.05 is the de-facto standard cut off between significant

- and non-significant results.
- If this de facto value is used as the critical value, it will result in wrong result 5% of the time.
 - p value tells us the probability that the results occurred by chance. It is the probability of wrongly rejecting null hypothesis. If p-value less than 0.05, we say that the test result is significant at 5% level of significance.
 - The observed significance level is a criterion that can be computed from a sample data. It is a probability
 - i) If $p\text{-value} < \alpha$; reject the null hypothesis. There is statistical significant difference between two treatments.
 - ii) If $p\text{-value} > \alpha$; accept the Null Hypothesis.

Q-1) d) ANOVA vs T-test

T-test

- i) T-test is a hypothesis test that is used to compare the means of two samples.

ANOVA

- i) ANOVA (Analysis of Variance) is a statistical technique that is used to compare means of more than two samples.

ii) The probable distribution is F-distribution.

iii) Null Hypothesis always differs based on the tailedness of the test.

iv) Test statistic used for T-test

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}}$$

(for one Sample)

Hypothesis & It will change as we moved from using one Sample to two Sample.

ii) The Probable distribution is F-distribution.

iii) Null Hypothesis i.e., $H_0: \mu_1 = \mu_2 = \mu_3$ or $H_1: \text{At least two population means are different.}$

iv) Test statistic for ANOVA is

$$F = \frac{S_b}{S_w}$$

Where S_b = between Column Variance,
 S_w = Within Column Variance,

Q.1) e) Nature of Normal Curve
→ A normal curve is a bell-shaped curve which shows the probability distribution of a continuous random variable.

→ Moreover, the normal curve represents a normal distribution... Also

the standard normal curve represents a normal curve with mean 0 and standard deviation 1.

- Here we have four characteristics of a Normal Distribution.

i) Normal distributions are symmetric, unimodal, and asymptotic.

ii) and the mean, median and mode are all equal.

iii) A normal distribution is perfectly symmetrical around its center.

iv) That is, the right side of the center is a mirror image of the left side.

- In normal distribution curve, The curve is symmetrical and bell shaped, showing that trials will usually give a result near the average, but will occasionally deviate by large amounts.

Q.1) f) Probability

→ In everything, we guess the chances of successful outcomes, from business to medicine to medicine to the weather.

→ A probability provides a quantitative description of the chances or likelihoods associated with various outcomes.

- It provides a bridge between descriptive and inferential Statistics
- It is a numerical measure which indicates the chance of occurrence.
- If an experiment has n equally likely simple events and if m be the favorable number of outcome to an event A . Then the probability of A , $P(A)$ is,

$$P(A) = \frac{\text{Number of favourable outcomes}}{\text{Total number of outcomes}}$$

$$= \frac{m}{n}$$

Results from Probability / Rules :-

- $0 \leq P(A) \leq 1$
- $P(A) = 1 - P(A')$
- $P(\emptyset) = 0$ where \emptyset is null event.

- Probability of an Event must have following Properties
- i) $P(A)$ must be between 0 and 1.
- ii) If event A can never occur, $P(A) = 0$. If event A always occurs when the experiment is performed, $P(A) = 1$
- iii) The sum of the probabilities for

all simple events in S equals 1.

Q.1) g) Cluster vs stratified Sampling	cluster	stratified
i) Sampling cluster	Sampling refers to a Sampling method where in the members of the population are selected at random, from naturally occurring groups called 'clusters.'	i) Stratified Sampling is one in which the population is divided into homogeneous segments, and then, the sample is randomly taken from the segments.
ii) Randomly selected clusters are taken from all the individuals.		ii) Randomly selected individuals are taken from all the strata.
iii) Selection of population elements are performed Collectively.		iii) Selection of population elements are performed individually.
iv) Between groups their might		iv) Homogeneous elements

Contain similarity of the data samples.

are found within group of each strata.

v) within groups there might be different categories of data.

vi) Heterogeneous elements are found within groups.

vi) Bifurcation of data has been done through naturally occurring groups.

vi) Bifurcation of data has been imposed by the researcher.

vii) Objective of performing cluster sampling is to reduce cost and improve efficiency.

vii) objective of performing stratified sampling is to increase precision and representation.

a.) b) Sampling error vs non-Sampling error

	Sampling error	Non-Sampling error
Meaning	Sampling error is a type of error, occurs due to the sample selected does not perfectly represents the population of interest.	An error occurs due to Sources other than Sampling, While Conducting Survey activities is known as non-Sampling error.

Cause:- Deviation between Sample mean and population mean

Deficiency and analysis of data.

Type Random

Random

Occurs only when Sample is selected.

or
Non-random
Both in Sample and census.

Sample Possibility of error reduced with the increase in Sample size

It has nothing to do with the Sample Size.

Q.1) i) Binomial vs Poisson Distribution

Binomial
Distribution

- i) Binomial distribution is biparametric.
i.e. It is featured by two parameters 'n' and 'p'.

Poisson

- i) where as Poisson distribution is unparametric i.e. characterized by a single parameter λ .

- 2) In a Binomial distribution, there are only two possible outcomes. i.e., Success or Failure.

- 2) Conversely, there are an unlimited number of possible outcomes in the case of Poisson Distribution.

- 3) In binomial distribution mean(np) $>$ Variance(npg)

- 3) while in Poisson distribution mean = Variance

- 4) For Discrete Variables

- 4) For Continuous Variables

Q.1) j) Sampling Distribution of the Mean

- The mean of the Sampling distribution of the mean is the mean of the population from which the scores were sampled.
- Therefore, if a population has a mean μ , then the mean of the Sampling distribution of the mean is also μ .
- The symbol μ_M is used to refer to the mean of the Sampling distribution of the mean.
- Therefore the formula for the mean of the Sampling distribution of the mean can be written as:

$$\mu_M = \mu$$

→ VARIANCE:-

The variance of the Sampling distribution of the mean is computed as follows:-

$$\sigma_M^2 = \frac{\sigma^2}{N}$$

That is the Variance of the Sampling distribution of the mean is the population Variance divided by N , the

Sample size (the number of scores used to compute a mean.) Thus, the larger the sample size the smaller the variance of the sampling distribution of the mean.

- The standard error of the mean is the standard deviation of the sampling distribution of the mean. It is therefore the square root of the variance of the sampling distribution of the mean and can be written as:

$$\sigma_M = \frac{\sigma}{\sqrt{N}}$$

- The standard error is represented by a σ because it is a standard deviation.
- The subscript (M) indicates that the standard error in question is the standard error of the mean.

Q.1) K) Central Limit Theorem

The Central Limit Theorem states that the Sampling distribution of the Sample means approaches a normal distribution as the sample size gets larger.

- no matter what the shape of the population distribution.

→ Central Limit Theorem (CLT) is an important because all the samples will follow a normal (approximately) distribution pattern with a standard deviation approximately equal to the standard deviation of the population, divided by each sample's size.

a. i) Correlation Coefficient:-

→ A correlation coefficient is a numerical measure of same type of correlation (meaning a statistical relation between two variables).

→ The two columns of a given data set can be the two variables.

→ Several types of correlation coefficient exist, each with their own definition and own range of usability and characteristics.

→ They all assume the values in the range of -1 to 1, where

± 1 indicates the strongest possible agreement and 0 indicating the strongest possible disagreement.

→ Certain problems include the possibility of incorrectly being used to infer causal relationship between the variables.

Properties:-

- 1) The absolute values of both sample and population (Pearson Correlation) are less than or equal to 1.
- 2) Correlation Coefficients equal to -1 or 1 correspond to data points lying exactly on a line.
- 3) The Correlation Coefficient (Pearson) is Symmetric than $\text{Corr}(x, y) = \text{Corr}(y, x)$

Q. 1 M) Outlier detection:-

→ Outliers are extreme values that deviate from other observations on data, they may indicate a variability in measurements, experimental error.

→ They are of 2 kinds: univariate and multivariate.

→ Most Common Causes of outliers on a data set:

- 1) Data entry errors (human)
- 2) Measurement errors (instrument)
- 3) Experiment errors (data extraction)

- 4) Data Processing errors (data manipulation)
 - 5) Sampling errors (extracting or mixing data.)
 - 6) Natural (novelties in data)
 - 7) Intentional (dummy to test detection methods)
- Some of the most popular methods of outlier detection:-
- Z-Score / Extreme Value analysis
 - Statistical modelling
 - Linear regression models
 - Proximity based models.

- Q.1) N) Properties of normal distribution
- The normal distribution is also referred to as Gaussian or Gauss distribution.
- It is made relevant by the Central limit theorem, which states that the averages obtained from independent, identically distributed random variables with almost the same normal distributions regardless of the type of distribution they are sampled from.
- Normal distribution follows the following characteristics.

1) It is Symmetric

- The distribution curve can be divided in the middle to produce two equal halves.

2) The mean, median, mode

- The middle point of a normal distribution is the point with the maximum frequency, which means it possess the most number of variables/ observations.

3) Empirical rule

- There is a constant proportion of distance lying under the curve between the mean and specific number of standard deviation from the mean.

Example: 68.25% of all cases falls within 1 standard deviation.
95% of all cases fall within 2 standard deviation.
99% of all cases fall within 3 standard deviation.

4) Skewness and Kurtosis

- They are Coefficients that measure how different a distribution is from a normal distribution.

- Skewness measures the symmetry which kurtosis measures the thickness of the tail ends.

Q.1) Q) Assumptions of Linear regression:-

- These are 5 basic assumptions of linear Regression:-
- Linear relationship between the features and target :-
- we assume that there is a linear relationship and thus capture that.
- Little or no multicollinearity between the features:
 - Multicollinearity is a state of very high inter-correlations among the independent variables.
 - It weakens the statistical power of the regression model.
- Homoscedasticity:
 - It describes a situation in which the errors term (that is the noise or random distribution in the relationship) is the same across all values of the independent variable.
- Normal distribution of error terms:
 - The assumption is that the error (residuals) follows a

normal curves

→ little or no auto correlation in the residuals.

- Auto correlation occurs when the residual errors are dependent on each other.
- The presence of correlation in error terms drastically reduces model's accuracy.

Q.1) p) Assumptions of Chi-Square goodness of fit test

Assumption 1:- one categorical variable.

→ dichotomous, nominal or ordinal

→ Examples:-

Dichotomous:- Gender (male, female)
Treatment (medication, no medication)
Education (under grad, Post Grad)

Nominal :- Ethnicity (Caucasian, Indian, Hispanic)
Profession (Doctor, teacher, nurse, dentist)

ordinal :- Scale (point scale)
Physical activity (sedentary, low, moderate, high)

Assumption 2: You should have independence of observation
→ which means you should have there is no relationship between any of the cases.

Assumption 3:- The group of the Categorical Variable must be mutually exclusive.

Assumption 4:- There must be at least 5 expected frequencies in each group of your categorical variable.

Q. 1) a) Correlation vs Covariance

Covariance:-

- It is the relationship between a pair of random variables where change in one variable causes change in another variable.
- It can take any value between $-\infty$ to $+\infty$, where the negative value represents the negative relationship whereas the positive value represent the positive relationship.
- It is used for Linear relationships.

→ It gives the direction of relationship between variables

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n}$$

Correlation:-

→ It shows whether and how strongly pairs of variable are related to each other.

→ Correlation takes values between +1 and -1, where in value close to +1 is strongly positive Correlation and value close to -1 is strongly negative Correlation.

→ It gives direction and strength of relationship between variables

$$\text{Corr}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

Q.1) 2) How to check the Normality of data.

→ ~~Look at~~ Value of the Shapiro-wilk Test is greater than 0.05, the data is normal. If it is below 0.05, the data significantly deviate from a normal distribution. If you need to use skewness and kurtosis values to determine normality, rather the Shapiro-wilk test you will find these in our enhanced testing for normality guide.

Q.1) 5) Point Estimate vs ~~an~~ Interval Estimate

→ A point estimate is a single number that is used to estimate an unknown population parameter.

Ex. While watching a football match, you say "I bet their line must weigh average 250 pounds"

- A point estimation is often insufficient, because it is either right or wrong.
- You do not know how right or how wrong is your answer (i.e. you do not know the margin of error)
- It is much more useful if it is accompanied by an estimate of error.
- An interval Estimate is a range of value used to estimate a population parameter.
- It indicates the error in two ways:-
 - By the extent of its range
 - By the probability of true population parameter lying within that range.

Ex. The dept head would say something like "I estimate that true enrollment in this course will fall between 330 and 380."

Q.1) \rightarrow t) F Test vs ANOVA

\rightarrow The Analysis of Variance (ANOVA) is the process of assigning the total variation into its components, in an experimental setup and each component is tested separately for its significance.

\rightarrow The test is a test for variance, a function of the corresponding sum of squares (SS). Each component SS is tested with respect to the Error sum of squares (ESS).

\rightarrow The ratio of these two sums of squares is the same as the ratio of two variances and therefore is a F-Test.

\rightarrow A simple example is the ANOVA for a two-way classification experimental set up.

\rightarrow The resultant Total sum of squares (TSS) is split up into two components as the Between sum squares (BSS) and within sum of squares (WSS).

\rightarrow The remainder is the ESS which is used for testing each of the other two.

\rightarrow The F-Test used for testing

any null hypothesis based on
the equality of any two vari-
ances, that may not be the
results of an ANOVA.

- Q. 2) Write the steps involved in testing of hypothesis. Also explain the Critical value approach and P-value approach.
- Test Procedure:- (Classical Method)
- 1) Determine and state H_0 (i.e. Null Hypothesis) and H_1 (i.e. Alternate Hypothesis)
 - 2) Decide the Significance level α and the critical region

- 3) Based on the parameter, choose the test statistic
- 4) Using available data Compute the test statistic
- 5) Make the statistical Accept or Reject Decision based on
 - a) Computed value of the test statistic
 - b) The critical region identified in step 2

• P-value method

- 1) Determine and state H_0 (i.e. Null Hypothesis) and H_1 (i.e. Alternate Hypothesis)
- 2) Decide the significance level α
- 3) Based on the parameter, choose the test statistic
- 4) Using available data Compute the test statistic and p value
- ~~p-value~~
5) Make the statistical Accept or Reject decision based on
 - $\alpha \geq p\text{-value}$
 - a) p-value less than α should reject H_0 .
 - b) p-value greater than α should not reject H_0 .