

Date \_\_\_\_\_

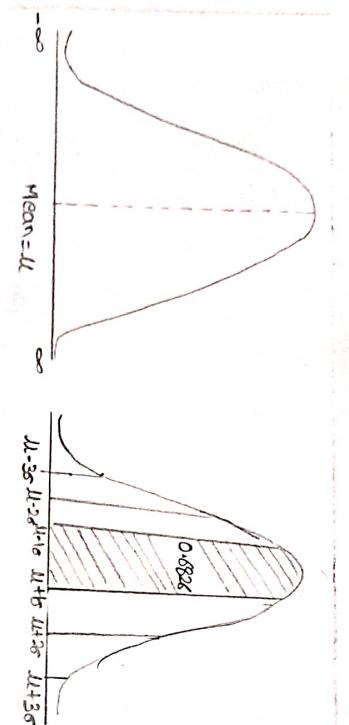
Topic Assignment No:-1

Name - Adarsh Kumar (L.E.)  
Roll No - 00676807223  
Branch - CSE3

Q-1 :-

Properties of Normal Distribution.

POLYtechnic



- Normal distribution is a type of continuous probability distribution that describes how random variable are distributed around their mean. It has the following properties
  - The normal distribution curve is symmetric and bell-shaped, meaning that the left and right halves of the curve are mirror images of each other.
  - The mean, median, and mode of the normal distribution are equal and located at the center of the curve.
  - The standard deviation of the normal distribution measure how spread out the data are from the mean. A smaller standard deviation means that data are more clustered around the mean, while a larger standard deviation means the data are more dispersed.
  - The total area under the normal distribution curve is equal to 1 or 100%. About 68% of the data are within one standard deviation of the mean, about 95% of the data are within two standard deviations of the mean, and about 99.7% of the data are within three standard deviations of the mean, and about 99.7% of the data are within three standard deviation of the mean. This is known as the empirical rule or the 68-95-99.7 rule.
  - The normal distribution is defined by two parameters: The mean ( $\mu$ ) and the standard deviation ( $\sigma$ ). The formula for the probability density function of the normal distribution is:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

where  $x$  is the random variable,  $\mu$  is the mean, and  $\sigma$  is the standard deviation.

Date \_\_\_\_\_

Topic \_\_\_\_\_

Q-2 In a test on 2000 Electric bulbs, it was found that the life of a particular make was normally distributed with an average life of 2040 hours and S.D. of 60 hours. Estimate the number of bulbs likely to burn for.

- ii) More than 2150 hours.
- iii) Less than 1950 hours.
- (iii) More than 1920 hours and but less than 2160 hours.

S-2 Apply the formula for the standard score (Z-score):

$$Z = \frac{X - \mu}{\sigma}$$

where  $X$  is the raw score,  $\mu$  is the mean, and  $\sigma$  is the standard deviation.

Using the given data, we can calculate the Z-scores of each of the cutoff points and then use a normal distribution table or a calculator to find the corresponding probabilities.

- (i) For more than 2,150 hours, the Z-score is:

$$Z = \frac{2150 - 2040}{60} = 1.83$$

The probability of getting a Z-score less than or equal to 1.83 is 0.9664. Therefore the prob. of getting a Z-score greater than 1.83 is  $1 - 0.9664 = 0.0336$ . This means that about 3.36% of the bulbs are likely to burn for more than 2,150 hours. Since there are 2,000 bulbs in total, we can estimate the number of bulbs by multiplying the prob. by the total number:  $0.0336 \times 2000 = 67.2$

Therefore, about 67 bulbs are likely to burn for more than 2,150 hours.

- (iii) For less than 1,950 hours, the Z-score is:

$$Z = \frac{1950 - 2040}{60} = -1.5$$

The probability of getting a Z-score less than or equal to -1.5 is 0.0668

Date \_\_\_\_\_

Topic \_\_\_\_\_

This means that about 6.68% of the bulbs are likely to burn for less than 1,950 hours. Again, we can estimate the number of bulbs by multiplying the prob. by the total number:

$$0.0668 \times 2000 = 133.6$$

about 134 bulbs are likely to burn for less than 1,950 hours

- (iii) For more than 1,920 hours but less than 2,160 hours, the Z-score are:-

$$Z_1 = \frac{1920 - 2040}{60} = -2$$

$$Z_2 = \frac{2160 - 2040}{60} = 2$$

The prob. of getting a Z-score between -2 and 2 is 0.9545. This means that about 95.45% of the bulbs are likely to burn for more than 1,920 hours but less than 2,160 hours.

$$0.9545 \times 2000 = 1909$$

Therefore, the about 1909 bulbs are likely to burn for more than 1,920 hours but less than 2,160 hours.

Q.3

Let the Joint P.d.f of X and Y be

$$f(x,y) = \begin{cases} (x+y) & 0 < x \leq 1, 0 \leq y \leq 1 \\ 0 & \text{otherwise} \end{cases}$$

Sol-3

Find  $P[0 < x < \frac{1}{2}, 0 < y < \frac{1}{4}]$ ,  $E[X]$ ,  $E[Y]$ ,  $E[X^4]$ ,  $E[X+Y]$ ,  $P[x, y]$

Date \_\_\_\_\_

Topic \_\_\_\_\_

$$P[0 < X < \frac{1}{2}, 0 < Y < \frac{1}{4}] = \int_0^{\frac{1}{2}} \left[ \int_0^{\frac{1}{4}} (x+y) dy \right] dx = \int_0^{\frac{1}{2}} \left[ xy + \frac{1}{2}y^2 \right]_0^{\frac{1}{4}} dx \\ = \int_0^{\frac{1}{2}} \left( \frac{1}{4}x + \frac{1}{32} \right) dx = \frac{1}{4}x^2 + \frac{1}{64}x^3 \Big|_0^{\frac{1}{2}} = \frac{3}{64}$$

$$\text{Now } E[X] = \int_0^1 \int_0^1 xf(x,y) dx dy = \int_0^1 \int_0^1 x(x+y) dx dy \\ = \int_0^1 \left[ x^2 y + x \left( \frac{1}{2} y^2 \right) \right]_0^1 dx = \int_0^1 (x^2 + \frac{1}{2}x) dx \\ = \frac{1}{3} + \frac{1}{4} = \frac{7}{12}$$

$$\text{Similarly, } E[Y] = \int_0^1 \int_0^1 yf(x,y) dx dy = \frac{7}{12}$$

$$\text{Now } E[XY] = \int_0^1 \int_0^1 xy f(x,y) dx dy = \int_0^1 \int_0^1 xy(x+y) dx dy \\ = \int_0^1 \left[ x^2 \left( \frac{1}{2} y^2 \right) + x \left( \frac{1}{3} y^3 \right) \right]_0^1 dx = \int_0^1 \left( \frac{1}{2}x^2 + \frac{1}{3}x \right) dx \\ = \frac{1}{6} + \frac{1}{6} = \frac{1}{3}$$

$$E[X+Y] = \int_0^1 \int_0^1 (x+y) f(x,y) dx dy = \int_0^1 \int_0^1 (x+y)^2 dx dy \\ = \int_0^1 \int_0^1 (x^2 + 2xy + y^2) dx dy \\ = \int_0^1 \left[ x^2 y + xy^2 + \left( \frac{1}{3} y^3 \right) \right]_0^1 dx \\ = \int_0^1 (x^2 + x + \frac{1}{3}) dx = \frac{1}{3} + \frac{1}{2} + \frac{1}{3} = \frac{7}{6}$$

$$\text{we have cov}[X,Y] = E[XY] - E[X]E[Y] = \frac{1}{3} - \left( \frac{7}{12} \right)^2 = -\frac{1}{144}$$

$$\text{Now } E[X^2] = \int_0^1 \int_0^1 x^2 f(x,y) dx dy = \int_0^1 \int_0^1 x^2(x+y) dx dy \\ = \int_0^1 \left[ x^3 y + x^2 \left( \frac{1}{2} y^2 \right) \right]_0^1 dx \\ = \int_0^1 (x^3 + \frac{1}{2}x^2) dx = \frac{1}{4} + \frac{1}{8} = \frac{5}{12}$$

Date \_\_\_\_\_

Topic \_\_\_\_\_

$$\text{similarly, } E[y^2] = \frac{s}{12}$$

$$\text{Now } \text{Var}[x] = E[x^2] - (E[x])^2 = \frac{s}{12} - \left(\frac{7}{12}\right)^2 = \frac{11}{144}$$

$$\text{similarly, } \text{Var}[y] = \frac{11}{144}$$

$$\text{Hence, } P[x,y] = \frac{\text{cov}[xy]}{\sqrt{\text{Var}[x]\text{Var}[y]}} = \frac{\frac{-1}{144}}{\sqrt{\frac{11}{144} \times \frac{11}{144}}} = \frac{-1}{11}$$

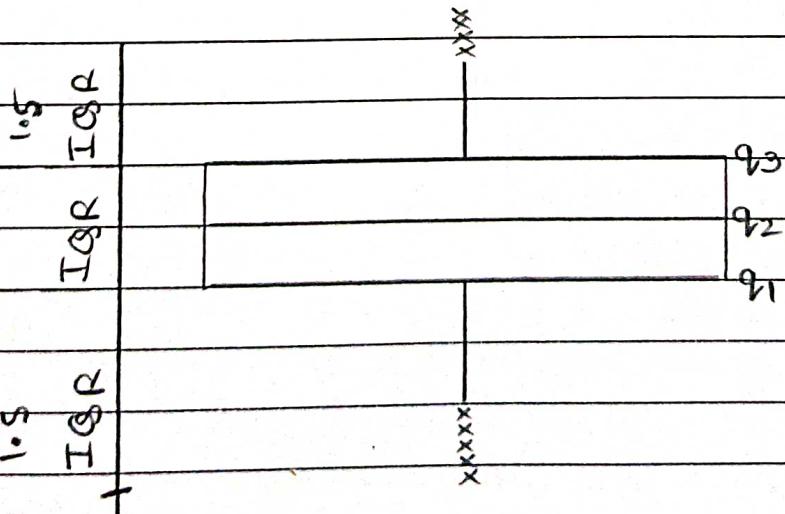
#### -4 BOX Plots

for a given data set, we determine all the quartiles  $q_1, q_2, q_3$  the difference  $q_3 - q_1$ . (upper quartile, lower quartile) is determined and is called IQR (Inner Quartile Range).

we draw a diagram with the data on the y-axis (in the increasing orders) we draw a box whose lower side is  $q_1$  & upper side is  $q_3$  and a middle line  $q_2$ .

ignoring the outliers we draw two straight lines above & below the box for the data within the 1.5 IQR range

The data beyond this is denoted as signal point and the called outliers.

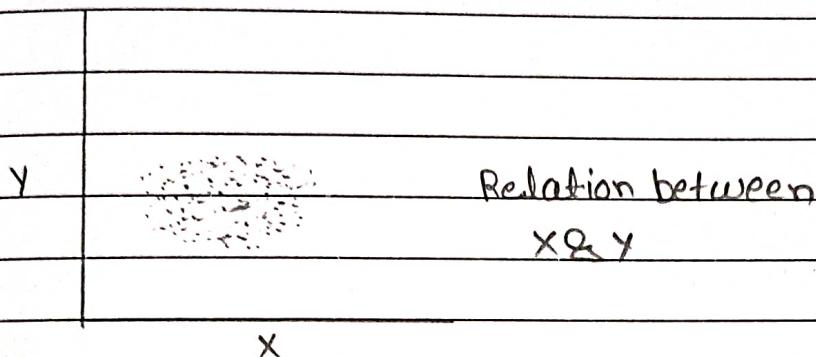


Date \_\_\_\_\_

Topic \_\_\_\_\_

## Q-5 Scatter Diagrams

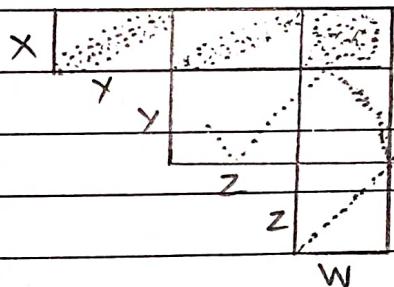
SI-5 A scatter diagram is conducted by plotting each pairs of observations with one measurement in the pair on the vertical axis of the graph and the other measurement in the pair on the horizontal axis.



Relation between  
X & Y

when two or more variables exists, the matrix of scatter diagram may be useful.

\* Sample correlation coefficient.



$$R_{xy} = \frac{\sum \Delta x_i \Delta y_i}{\sqrt{\sum \Delta x_i^2 \sum \Delta y_i^2}} \quad \Delta x_i = x_i - \bar{x}$$

$$\Delta y_i = y_i - \bar{y}$$

Q6 The mean of a certain normal population is Equal to the standard error of the mean of the standard sample of 100 form that distribution. Find the prob. that the mean of the sample of 25 from the distribution will be negative.

SI-6 we need to use the mean of the formula for the standard error of the mean:

$$SE = \frac{\sigma}{\sqrt{n}}$$

Date \_\_\_\_\_

Topic \_\_\_\_\_

where,  $\sigma$  is the standard deviation and  $n$  is the sample size

$\therefore$  The standard deviation error of the mean is  $\frac{\sigma}{\sqrt{n}} = \frac{\sigma}{\sqrt{10}}$

Given mean is equal to the standard error of the mean, therefore

$$\text{For ex. } \bar{x} = \frac{\sigma}{\sqrt{10}} = 1$$

The z-test statistical formula is given by  $z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$

using equation(1) in z-test statistical formula for the sample size 25,

$$z = \frac{\bar{x} - 10}{\frac{\sigma}{\sqrt{25}}} = \frac{10 - \sigma}{10} \times \frac{5}{\sigma} = \frac{(10 - \sigma)5}{10\sigma} = \frac{50\bar{x} - 5\sigma}{10\sigma} = \frac{5\bar{x} - \sigma}{2\sigma}$$

$$z = \frac{5\bar{x} - \sigma}{2\sigma}$$

Since, the population is normally distributed, and if  $z < 0.5$ ,

$\bar{x}$  is negative. therefore the probability is given by

$$P(z < 0.5) = P(-\infty < z < 0.5)$$

Using the table of z-values,

$$P(z < 0.5) = 0.5 + P(0 < z < 0.5)$$

### Ex-7 Statistical Interval for a single sample

Sol-7 Statistical intervals provide a range of values within which a population parameter is estimated to lie with a certain level of confidence.

- for a single sample, intervals are constructed to estimate population parameters such as the mean, proportion, or variance.

#### ① Type of Intervals.

Confidence Intervals (CI)  $\hat{=}$  Estimate the range within which the population parameter is likely to fall

$$\text{mean: } CI = \bar{x} \pm z^* (\sigma/\sqrt{n}), \text{ proportion } CI = p \pm z^* / (p(1-p))$$

Date \_\_\_\_\_

Topic \_\_\_\_\_

Variance : ( $I = [(n-1)s^2 / \chi^2_{\text{Upper}}, (n-1)s^2 / \chi^2_{\text{Lower}}]$ )

2) prediction Intervals (PI) :- provide a range within which the population parameter value of a future observation

$$\text{mean}(\mu) := PI = \bar{x} + (*s^*)\left(1 + 1/n\right)$$

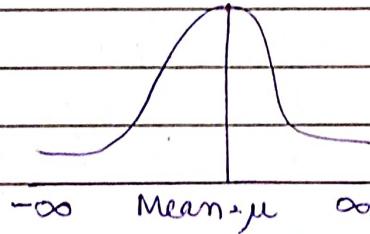
$$\text{Individual } (\xi) := \bar{x} + t^* (s^*)$$

## Ques 1 Properties of Normal Distribution

- (1) The normal probability curve with mean  $\mu$  and standard deviation  $\sigma$  is given by

$$f(x) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} \quad -\infty < x < \infty$$

- (2) The graph of the curve is



- (3) The curve is bell-shaped and symmetrical about the line  $x = \mu$ .

- (4) Mean, median and mode of the normal distribution coincide and the normal distribution is unimodal.

- (5)  $f(x)$  decreases rapidly as  $x$  increases.

- (6) X-axis is an asymptote to the curve (the tangent to the curve at  $\infty$  is X-axis)

- (7) The maximum probability occurs at the point  $x = \mu$  and is

$$\frac{1}{\sigma \sqrt{2\pi}}$$

(8) Mean deviation about mean =  $\frac{4}{5} \sigma$

(9) Since  $f(x)$ , being the probability, can never be negative, so that no portion of the curve lies below the  $x$ -axis.

(10) A linear function of independent normal variates is also normal variate.

Ques 2 In normal distribution, 31.1% of the items are under 45 and 8.1% are over 64. Find the mean and standard deviation of the distribution.

$$z_1 = 0.31$$

$$z_2 = 0.92 \quad (1 - 0.08 = 0.92, \text{ since } 8.1\% \text{ are over } 64, 92\% \text{ are below } 64)$$

$$n_1 = 45^{\text{th}} \text{ percentile}$$

$$n_2 = \frac{92}{64}^{\text{th}} \text{ percentile}$$

From the standard normal distribution table

$$z_1 \approx -0.48$$

$$z_2 \approx 1.41$$

$$\rightarrow z = \frac{x - \mu}{\sigma}$$

$$n_2 = \mu + z \times \sigma$$

$$\rightarrow n_1 = \mu + (-0.48) \times \sigma$$

$$\rightarrow n_2 = \mu + (1.41) \times \sigma$$

$$n_1 = 45, n_2 = 64$$

Topic..... Date.....

$$45 = \mu - 0.48 \times \sigma \rightarrow ①$$

$$64 = \mu + 1.41 \times \sigma \rightarrow ②$$

$$64 = 45 + 0.48 \times \sigma + 1.41 \times \sigma$$

$$64 = 45 + 1.89 \times \sigma$$

$$19 = 1.89 \times \sigma$$

$$\boxed{\sigma \approx 10.05} \quad \underline{\text{Ans}}$$

$$\mu = 45 + 0.48 \times 10.05$$

$$\boxed{\mu \approx 49.82} \quad \underline{\text{Ans}}$$

Ques 3

X and Y are two random variable having joint density function  $\frac{1}{27}(2n+y)$ , where n and y can assume only

integer values 0, 1, 2. Find conditional distribution of Y for  $X=n$ .

To find the conditional distribution of Y given  $X=n$ , we need to calculate PMF of Y when  $X=n$ , denoted as  $P(Y|X=n)$ .

$$P(Y|X=n) = \frac{P(X=n, Y)}{P(X=n)}$$

$$f(n, y) = \frac{1}{27}(2n+y)$$

$$P(X=0) = \sum_{y=0}^2 f(0, y) = \sum_{y=0}^2 \frac{1}{27} (2(0)+y) = \frac{1}{27} (0+1+2) = \frac{3}{27}$$

$$P(X=1) = \sum_{y=0}^2 f(1, y) = \sum_{y=0}^2 \frac{1}{27} (2(1)+y) = \frac{1}{27} (2+3+4) = \frac{9}{27}$$

$$P(X=2) = \sum_{y=0}^2 f(2, y) = \sum_{y=0}^2 \frac{1}{27} (2(2)+y) = \frac{1}{27} (4+5+6) = \frac{15}{27}$$

Teacher's Sign.....

Now let's calculate  $P(X=n, Y)$  for  $X=0, 1, 2$ :

- For  $X=0$ ,

$$P(X=0, Y) = \frac{1}{27} (2(0)+Y) = \frac{1}{27} Y$$

- For  $X=1$ :

$$P(X=1, Y) = \frac{1}{27} (2(1)+Y) = \frac{1}{27} (2+Y)$$

- For  $X=2$ :

$$P(X=2, Y) = \frac{1}{27} (2(2)+Y) = \frac{1}{27} (4+Y)$$

$$\Rightarrow P(Y|X=0) = \frac{P(X=0, Y)}{P(X=0)} = \frac{1/27 Y}{3/27} = \frac{Y}{3}$$

$$P(Y|X=1) = \frac{P(X=1, Y)}{P(X=1)} = \frac{1/27 (2+Y)}{9/27} = \frac{2+Y}{9}$$

$$P(Y|X=2) = \frac{P(X=2, Y)}{P(X=2)} = \frac{1/27 (4+Y)}{15/27} = \frac{4+Y}{15}$$

So, the conditional distribution of  $Y$  given  $X=n$  is

- $P(Y|X=0) = \frac{Y}{3}$  for  $Y=0, 1, 2$

- $P(Y|X=1) = \frac{2+Y}{9}$  for  $Y=0, 1, 2$

- $P(Y|X=2) = \frac{4+Y}{15}$  for  $Y=0, 1, 2$

#### Ques4 Stem and Leaf Diagram

A stem and leaf diagram is a simple and effective method for displaying a set of quantitative data. It allows for a quick visual representation of the distribution of the data while retaining the individual datapoints.

Topic..... Date.....

- (1) Stem: The stem consists of the leading digits of the data values. It represents the larger part of the number. For example, if the data values are 23, 27, 34, 36 the stems are 2, 2, 3, 3.
- (2) Leaf : The leaf consists of the last digit of each data value. It represents the finer part of the number. Using the same example data, the leaves would be 3, 7, 4, 6.
- (3) Organization : - The stems are written vertically in ascending order on the left side of a vertical line, and the corresponding leaves are placed to the right of each stem. Leaves are typically arranged in ascending order with each stem.
- (4) Example : - using the data above, a stem and leaf plot might look like this :

2 | 3 7

3 | 1 4 6

This indicates that there are two data points in the 20s range (23 and 27) and two data points in the 30s range (34 and 36).

### Ques5 Time Sequence Plots

A time sequence plot, known as a time series plot, is a graphical representation of data points collected or recorded over time. It's a fundamental tool in analyzing time based data in various fields such as economics, finance, weather forecasting etc.

- (1) Representation of Time:- The horizontal axis of a time sequence plot represents time, typically with equally spaced

Teacher's Sign.....

intervals. Time can be measured in various units depending on the nature of the data, such as minutes, hours, days, months or years.

(2) Vertical Axis :- The vertical axis represents the values of the data being measured or observed over time. This could include variables like temperature, stock prices, sales figures.

(3) Plotting Data Points :- Each data point is plotted on the graph at its corresponding time and value. Connecting these points with lines helps visualize trends and patterns in the data over time.

Ques 6 What is point estimation?

Point estimation is a method used in statistics to estimate an unknown population parameter based on sample data. It involves using a single value, called a point estimator, to approximate the true value of the parameter of interest. The goal of point estimation is to provide a best guess or approximation of the parameter, given the available sample information.

Ques 7 Central Limit Theorem

The central limit theorem (CLT) is one of the most important and fundamental concepts in statistics. It states that regardless of the distribution of the population from which a sample is drawn, the sampling distribution of the sample mean approaches a normal distribution as the sample size increases, under certain conditions.

↳ The Central Limit Theorem holds under the following

Topic..... Date.....

conditions:-

- The sample is drawn from a population with a finite mean ( $\mu$ ) and a finite standard deviation ( $\sigma$ )
- The sample size is sufficiently large. While there is no strict rule for what constitutes "sufficiently large", a commonly used guideline is that the sample size should be at least 30. However, the CLT often starts to apply even for smaller sample sizes if the population distribution is not heavily skewed.

Teacher's Sign.....

## Assignment

- (Q1) X is a normal variate with mean 30 and standard deviation 50, Find the probability that
- $26 \leq x \leq 40$
  - $x \geq 45$
  - $|x - 30| > 5$

Ans- Now, to find the probability the basic formula used is -

$$z = \frac{x - \mu}{\sigma}$$

where,  $\mu$  is the mean and  
 $\sigma$  is the standard deviation

(i) For  $x = 26$ :

$$z_1 = \frac{26 - 30}{50} = -0.08$$

For  $x = 40$ :

$$z_2 = \frac{40 - 30}{50} = 0.20$$

$$P(-0.08 \leq z \leq 0.20) \cong P(z \leq 0.20) - P(z \leq -0.08)$$

$P(-0.08 \leq z \leq 0.20)$  is approximately 0.4515.

(ii) Probability that  $x \geq 45$ :

here,  $x = 45$ :

$$z = \frac{45 - 30}{50} = 0.30$$

$P(z \geq 0.30)$  is approximately 0.3821.

(iii) Probability that  $|x - 30| > 5$ :

Here, we standardize  $x = 25 \& x = 35$ :

$$z_1 = \frac{25 - 30}{5} = -0.10$$

$$z_2 = \frac{35 - 30}{5} = 1.0$$

$$\begin{aligned} P(z < -0.10) + P(z > 1.0) \\ = 0.4602 + 0.4602 = 0.9204 \end{aligned}$$

So, the probabilities obtained are -

$$P(20 \leq x \leq 40) \approx 0.4515 \quad \text{and} \quad P(x \geq 45) \approx 0.3821$$

$$P(|x - 30| > 5) \approx 0.9204$$

Q2 In a referendum 60% of voters voted in favour. A random sample of 200 voters was selected what is the probability that in the sample -

(i) More than 130 voted in favour?

(ii) Between 105 & 130 inclusive voted in favour?

(iii) 120 voted in favour?

Ans-

$$P(x = k) = \binom{n}{k} \times p^k \times (1-p)^{n-k}$$

•  $P(x=k)$  is the probability of getting exactly  $k$  successes in  $n$  trials.

•  $\binom{n}{k}$  is the number of combinations of  $n$  items taken  $k$  at a time.

•  $p$  is the probability of success in one trial.

•  $n$  is the no. of trials.

Also,  $p = 0.60$  (probability of voting in favour)

$$n = 200$$

(i) Probability that more than 130 voted in favour:

$$P(X > 130) = 1 - P(X \leq 130)$$

Using the binomial probability formula:

$$P(X \leq 130) = \sum_{x=0}^{130} \binom{200}{x} (0.60)^x (0.40)^{200-x}$$

Calculating the sum:

$$P(X \leq 130) \approx 0.0048$$

$$P(X > 130) = 1 - 0.0048 = 0.9952$$

(ii) Probability that between 105 and 130 inclusive voted in favour -

$$P(105 \leq X \leq 130) = \sum_{x=105}^{130} \binom{200}{x} (0.60)^x (0.40)^{200-x}$$

Calculating the sum -

$$P(105 \leq X \leq 130) \approx 0.9609$$

(iii) Probability that exactly 120 voted in favour:

$$P(X = 120) = \binom{200}{120} (0.60)^{120} (0.40)^{80}$$

Calculating the probability -

$$P(X = 120) \approx 0.0650$$

### Q3 Frequency distribution and Histograms

Ans- Frequency Distribution → The frequency of a value is the number of times it occurs in a dataset. A frequency distribution is the pattern of frequencies of a variable. It's the no. of times each possible value of a variable occurs in a dataset.

There are different types of frequency distributions some of them are listed below -

(i) Grouped Frequency Distribution ⇒

The number of observations of each class interval of a variable class intervals are ordered groupings of a variables values.

We can also use this type of frequency distribution for quantitative variables.

Relative frequency distribution → The number of observations of each class or the proportion of observation of each value in class interval of a variable.

We can use this type of frequency distribution for any type of variable when we are more interested in comparing frequencies than the actual no. of observations.

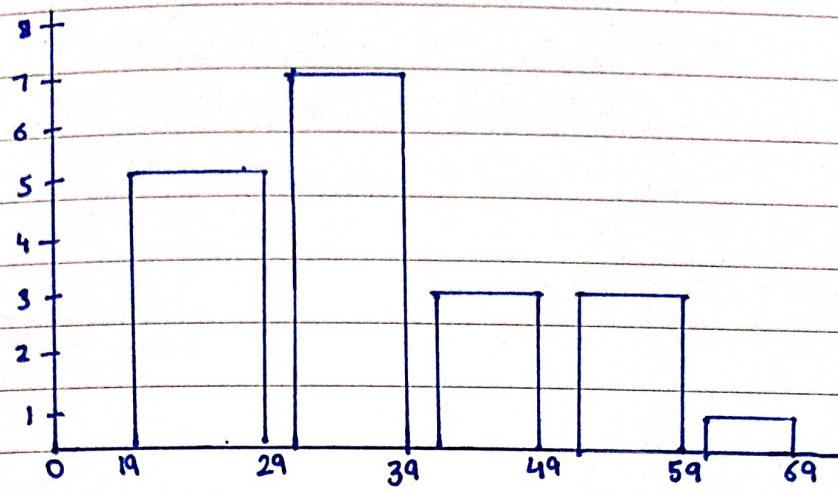
Ungrouped frequency distribution → The number of observations of each value of a variable. We can use this type of frequency distribution for categorical variables.

Cumulative frequency distribution → The sum of the frequencies less than or equal to each value or class interval of a variable.

We can use this type of variable or frequency distribution for ordinal or quantitative variables when you want to understand how often observations fall below certain values.

Histogram → A histogram is a graph that shows the frequency or relative frequency distribution of a quantitative variable. It looks similar to a bar chart. The continuous variable is grouped into interval classes, just like a grouped frequency table. The y-axis of the bars shows the frequencies or relative frequencies, and the x-axis shows the interval classes. Each interval class is represented by a bar and the height of the bar shows the frequency or relative frequency of the interval class.

A histogram is an effective visual summary of several important characteristics of a variable.



#### (Q4) Probability Plots

Ans - The probability plot is a way to visually comparing the data coming from different distributions. These data can be of empirical dataset or theoretical dataset. The probability plots can be of two types-

- P-P Plot → The p-p plot is the way to visualize the comparing of cumulative distribution function (CDF's) of the two distributions against each other.
- Q-Q Plot → The q-q plot is used to compare the quantities of two distributions. The quantities can be defined as continuous intervals with equal probabilities or dividing the samples between a similar way. The distributions may be theoretical or sample distributions from a process etc. The normal probability plot is a case of the q-q plot.

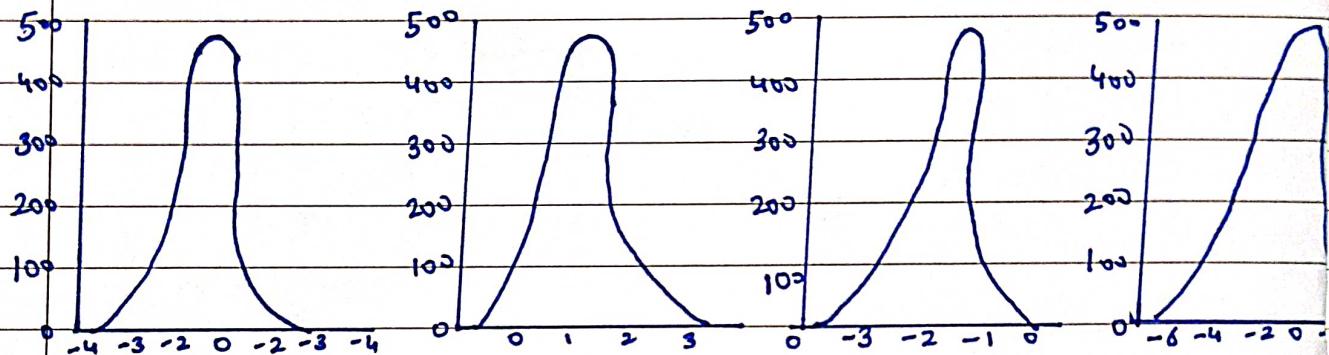
Normal Probability plot → The normal probability plot is a way of knowing whether the dataset is normally distributed or not. In this plot, data is plotted against the theoretical normal distribution plot in a way such that if a given dataset is

normally distributed it should form an approximate straight line.

The normal probability case is a case of the probability plot.

The normal probability plot has the following axis -

- Horizontal axis → Normal-order statistic medians.
- Vertical axis → order response values.



Standard

Right skewed

left skewed

Heavy-tail

### Q5 General concept of point estimation.

Ans- A point estimate is defined as a calculation where a sample statistic is used to estimate or approximate an unknown population parameter.

Point estimators are defined as functions that can be used to find the approximate value of a particular point from a given population parameters. The sample data of a population is used to find the point estimate or a statistic that can act as the best estimate of an unknown parameter that is given for a population.

Properties of point estimators →

- consistent
- unbiased
- most efficient

The maximum likelihood method is a popularly used way to calculate point estimators. This method uses differential calculus to understand the probability function from a given no. of sample parameters.

Formulae that can be used to measure point estimators -

- Maximum likelihood estimation
- Jeffrey estimation
- Wilson estimation
- Laplace estimation

Point estimate for  $p$  statistics →

$\hat{p} = x/n$ , the proportion of success in the sample, to be the point estimate of  $p$ .

Point estimation of the mean →

It is determined by calculating the mean of a sample drawn from the population.

## Q6 Sampling Distribution

Ans- A sampling distribution refers to the distribution of a statistic such as the mean, standard deviation, or proportion, calculated from multiple samples drawn from the same population. It provides insights into the behaviour & variability of the statistic under repeated sampling.

Example → Consider a population of heights of all adult males in a city. We want to estimate the average height of adult males in the city.

- (i) Sampling → We take multiple random samples of, say, 50 to adults males each from the population.
- (ii) Population parameter → Let,  $\mu$  be the population mean height of adult males.
- (iii) Sample statistic → For each sample, we calculate the sample mean height, denoted as  $\bar{X}$ .
- (iv) Sampling Distribution → The distribution of these sample means ( $\bar{X}$ ) across all samples forms the sampling distribution of the sample mean.

## Q7 Central limit Theorem

Ans- The central limit theorem relies on the concept of a sampling distribution, which is the probability distribution of a statistic for a large no. of samples taken from a population.

The central limit theorem says that the sampling distribution of the mean will always be normally distributed, as long as the sample size is large enough. Regardless of whether the population has a normal, poisson, binomial, or any other distribution, the sampling distribution of the mean will be normal.

Formula → mean -  $M\bar{X} = \mu$

Standard deviation -  $\sigma\bar{X} = \frac{\sigma}{\sqrt{n}}$

$$\bar{X} \sim N(\mu, \frac{\sigma}{\sqrt{n}})$$

where,  $\bar{X}$  is the sampling distribution of the sample means.

$\sim$  means "follow the distribution"

$\mu_{\bar{x}}$  is the mean of that population

$\sigma$  is the standard deviation of the population

$n$  is the sample size

Conditions of the central limit theorem  $\rightarrow$

- (i) The sample size is sufficiently large. This condition is usually met if the sample size is  $n \geq 30$ .
- (ii) The samples are independent & identically distributed random variables. This condition is usually met if the sampling is random.
- (iii) The population's distribution has finite variance. Central limit theorem doesn't apply to distributions with infinite variance, such as the cauchy distribution. Most distribution have finite variance.