

Statistics

What is Statistics?

→ It is a branch of Mathematics which helps to collect, organize, analyse the data.

Why use Statistics?

1. Summarization → which is mean, median, mode, count.

2. Test assumption →

3. Visualisation →

→ To find what really is going on.

→ To spot the pattern & problem.

→ To smart decision using data.

→ To simplify & understand large data.

→ For making future predict or

estimates,

When to use Statistics?

→ When we need to understand large data sets.

→ When we want to analyze the results of a survey or experiment.

→ When we need to make decisions based on trends or patterns.

Where to use statistics?

1. education \Rightarrow Analyzing student result, school performance.
2. Business \Rightarrow Market analysis, customer behaviour, profit/loss analysis.
3. government \Rightarrow Economic planning.
4. Medical field \Rightarrow Analysing drug effectiveness, patient recovery date.
5. sports \Rightarrow Player performance analysis, team statistics.

How to use statistics?

1. Data collection \Rightarrow using surveys, experiments.
2. Data organization \Rightarrow using table, charts, grouping.
3. Data Analysis \Rightarrow using averages, %, standard deviation, correlation.
4. Data interpretation \Rightarrow Understanding the meaning of the results.
5. Conclusion / Decision Making \Rightarrow Taking actions based on the results.

Types of statistics

① Descriptive

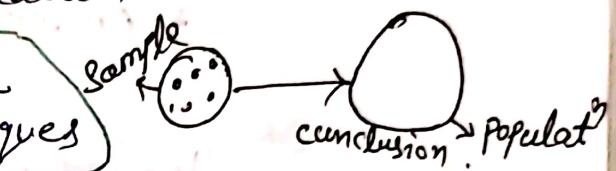
② Inferential

① Descriptive Statistics \Rightarrow Summary
They are some method for summarizing & understanding the key aspect of the data.
 \Rightarrow Summarizing data (Mean, Median, mode, graphs, tables etc.)
 \Rightarrow data ke bare mein kuch information mile.

② Inferential statistics \Rightarrow conclusions.
 \Rightarrow Drawing conclusions about the future of a larger population based on data.

Using data we can make some conclusions using some techniques

Types of ~~Variable~~ Variable



① Quantitative

Numerical

Discrete
int

Continuous
float

② Qualitative
Categorical

Nominal

Ordinal

① Quantitative Data \Rightarrow ~~(Numerical data)~~

\Rightarrow This type of data consists of numbers and can be measured.

\Rightarrow Use graph: line, histogram.

Discrete \Rightarrow countable data, usually whole numbers.
(whole no., int)

Ex No. of student, No of bank ac, no. of cars.

continuous \Rightarrow measurable data that can take on any value within a range.
(fractional, decimal)

Ex Height \rightarrow 79.1, weight \rightarrow 50.7, temp., time.

② Qualitative data \Rightarrow (categorical data) \rightarrow labels.

\Rightarrow This type of data describes non-numeric information.

\Rightarrow graph \rightarrow Bar, Pie.

Nominal \Rightarrow Data with categories that have no order.

\Rightarrow Order does not matter.

Ex gender, color, religion, blood type, subject

Ordinal \Rightarrow Data with categories that have a meaningful order, but the difference between values are not measurable.

\Rightarrow Order matters.

Ex feedback \rightarrow good, bad, Avg

size \rightarrow small, med., large

education \rightarrow High school, Bachelors, masters

Scales of Measurements

\Rightarrow It defines how the data is classified compared & used mathematically.

Properties

\rightarrow Identity (label)

\rightarrow Order (ranking)

\rightarrow Equal interval

\rightarrow True zero

Scales of Measurement

Nominal Scale

Ordinal Scale

Interval Scale

Ratio

① Nominal Scale

⇒ Naming data

⇒ Label on category data with no or

Ex: gender, Region, country

② ordinal scale

⇒ Label on category with order

Ex: customer feedback
educatⁿ level.

③ Interval Scale

- ⇒ Numeric scale with equal intervals, between values.
- ⇒ No true zero value.
- ⇒ No true zero matlab addition, subtraction karsakta hai; Temp, calendar. (*, / nehi karsakta hai)
- ⇒ Zero does not mean absence.
- ⇒ Mean, median, standard deviation find karsakta hai.
- ⇒ Histogram, Boxplot depends on data.
- ⇒ Correlation find karsakta hai. Numerical ke bich main se find karsakta hai.

④ Ratio

- ⇒ Equal interval & true zero value.
- ⇒ Yaha pe zero ka matlab zero hi hai, zero ka matlab nothing.
- ⇒ Height, weight, money
- ⇒ All mathematical operations can be performed.
- ⇒ Yaha pe line chart, box plot, histogram Sab banata hain.
- ⇒ Correlation bhi find karsakta hai.

Population & Sample

Population → whole data set.

→ A population is the entire group of individual items or data points that you study & draw conclusions from,

→ It is the complete set of data.

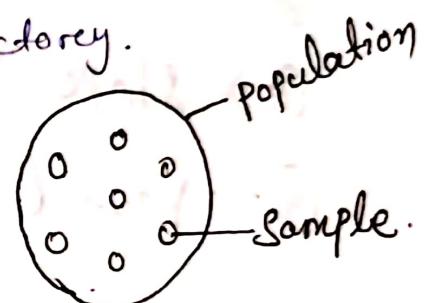
ex. • All students in a university.

• All citizens of a country.

• Every product from a factory.

→ Population represent (N)

Sample



→ A sample is a ~~size~~ subset of population, selected for the actual study or analysis.

→ Sample represent (n)

What is data

→ Facts or pieces of information that can be measured.

Sampling Techniques

① Simple Random Sampling.

② Stratified Sampling.

③ Systematic Sampling.

④ Convenience Sampling.

Descriptive Statistics

- Descriptive statistics are the methods for summarizing & understanding key aspects of data.
- key aspects means the most important parts, features, or characteristics of something.

Types of DS

- ① → Central Tendency → Mean, Median, mode.
- ② → Dispersion → Range, Variance, SD
- ③ → Shape → Skewness, Kurtosis
- ④ → Visual → Diff of charts.

* What is descriptive statistics? → Summarization

→ DS is the branch of statistics that focuses on summarizing & organizing data so it can be easily understood.

* Why used? → find pattern/info

- Summarize large datasets in a simple, meaningful way.
- Identify patterns or trends
- Make data easier to interpret.
- Present data visually for reports or presentations.

- * When is DS used? → before decision making.
- Before analysis to explore and understand the dataset
 - In reports and presentations to summarize result.
 - In research to describe characteristics of a population or sample.
 - In business to track performance (e.g., avg sales, customer age)

* Where is DS used? → every where.

- In education (e.g., average student scores)
- In business (e.g., monthly sales performance)
- In healthcare (e.g., patient age distribution)

① Measure of Central Tendency

1. # Measure of central Tendency

- ~~Mean~~ Measure of central tendency describes the "center" of the data.
- They provide summary of the data with a single value that represents the data as a whole.

① Mean → average

→ Refers to the center of distribution.

Used when:

Data is numeric and evenly distributed.

Population

Formula:

$$M = \frac{\sum_{i=1}^N x_i}{N}$$

Sample

Formula:

$$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$M = Mu$$

N = Total no. data points
in pop.

n = Total no. of
data points in
sample.

\sum = summation

x_i = particular data point

i = row no.

Ex

$$\text{Data } (1, 2, 2, 2, 3, 3, 4, 5, 5, 6)$$

$$\begin{aligned}\bar{x} &= (1+1+2+2+3+3+4+5+5+6) \div 10 \\ &= \frac{32}{10} = 3.2\end{aligned}$$

Outliers

→ a data point(s) who doesn't follow the pattern of data / or are extreme points.

→ Mean is sensitive of outliers.

Ex:

$$\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$$

outlier

outliers

$$\frac{132}{11} = 12$$

② Median

→ Middle value.

The median is a measure of central tendency and represents the middle value in a sorted list of numbers.

Even:

If the number of values is even, the median is the average of the two middle numbers.

Formula: $\frac{\left(\frac{n}{2}\right)^{\text{th}} + \left(\left(\frac{n}{2}\right) + 1\right)^{\text{th}}}{2}$

Ex: $\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$

$$\frac{10^{\text{th}}}{2} + \left(\frac{10}{2} + 1\right)^{\text{th}} = \frac{5^{\text{th}} + 6^{\text{th}}}{2}$$

$$= \frac{3 + 3}{2} = 3$$

Ex
List: 10, 2, 5, 8
Sorted: 2, 5, 8, 10

$$\text{Median} = \frac{5+8}{2} = 6.5$$

Odd:

⇒ If the number of values is odd, the median is the middle number!

Formula: $\left(\frac{n+1}{2}\right)^{\text{th}}$

ex: $\{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$

$\cdot \left(\frac{11+1}{2}\right)^{\text{th}} = 6^{\text{th}} = 3$

ex
List = 3, 1, 7
Shorted = 1, 3, 7
Median = 3

③ **Mode** ⇒ most frequent

⇒ The mode is the number that appears most frequently in a set of values.

ex Data set:

3, 7, 3, 2, 9, 3, 5

Here, 3 appears three times more than any other number.

So, the mode is 3.

Why use mode?

⇒ It's helpful in understanding the most common value.

* [1, 1, 2, 3, 5, 1, 1, 3, 1, 3, 5, 1]

1 = 6 3 = 3
2 = 1 5 = 2



2 # Measure of Dispersion

- ⇒ Measure of dispersion are statistical tools used to describe the spread or variability of a data set.
- ⇒ They show how much the data values differ from the average or mean.



① Range

- ⇒ Difference between the highest and lowest values.

Formula:-

$$\text{Range} = \text{Maximum} - \text{minimum}$$

- ⇒ Simple but sensitive to outliers.

Ex

Data: 5, 7, 8, 9, 100

$$\text{Range} = 100 - 5 = 95$$

yahan 100 ek outlier hai (baki sab value choti hai), aur ikki waajah se range bahut bda ho gaya isliye range ek simple but outlier-sensitive measure of dispersion hai.

- ⇒ isiliye hum data ko sahi se samajhne ke liye variance ka use karta hain.

② Variance

→ It measures how far each data point is from the mean.

High Variance → Spread ← (More data are far from mean)

Low Variance → Spread ← (More data around mean)

Population

$$\sigma^2 = \sum_{i=1}^N \frac{(x_i - \mu)^2}{N}$$

Sample

$$s^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n-1}$$

③ Standard Deviation

→ Standard deviation is measure of how spread out the values in a data set are from the mean (avg).

→ It tells you how much individual data points differ from the mean value.

Key Points:

→ A low standard deviation means data points are close to the mean.

→ A high standard deviation means data points are spread out over a wider range.

3. # IQR \rightarrow Inter Quartile Range

Percentile

\Rightarrow A Percentile is a value below which certain percentage of observation lie.

$$\text{Percentile} = \frac{I(\alpha) - 1}{n} \times 100$$

\curvearrowleft Index of α
 \curvearrowright Total observations

$$\frac{\text{Percentile}}{100} \times n = I(\alpha) - 1$$

$$\left(\frac{\text{Percentile}}{100} \times n \right) + 1 = I(\alpha)$$

Percentile rank?

10 \rightarrow value

$$\frac{\text{no. of values below } x}{n} \times 100$$

$$= \frac{16}{80} \times 100 = 80$$

Five Number Summary

1. Minimum $\rightarrow Q_0$
2. 25% \rightarrow First Quartile (Q_1)
3. Median \rightarrow 50 percentile Q_2
4. 75% \rightarrow Third Quartile (Q_3)
5. Maximum $\rightarrow Q_4$

Dataset = 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 9, 10,
 ↓
 Ascending
 n = 20

$$Q_1 \rightarrow 25\% \text{ tile} = \left(\frac{25}{100} \times 20 \right) + 1 = 6^{\text{th}} \text{ index}$$

→ 5

$$Q_3 \rightarrow 75\% \text{ tile} = \left(\frac{75}{100} \times 20 \right) + 1 = 15 + 1$$

= 16^{\text{th}} \text{ index}

→ 9

IQR = Inter quartile range
 $= Q_3 - Q_1 = 9 - 5 = 4$

$$\underline{\text{Min}} = \text{Lower fence} \Rightarrow Q - 1.5(\text{IQR})$$
 $= 5 - 1.5(4)$
 $= 5 - 6 = -1$

$$\underline{\text{Max}} = \text{Higher fence} = Q_3 + 1.5(\text{IQR})$$
 $= 9 + 1.5(4) = 9 + 6$
 $= 15$

$x < -1 \rightarrow \text{outlier}$

$x > 15 \rightarrow \text{Outlier}$

Box Plot

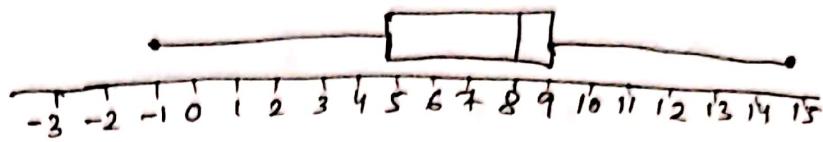
Min = -1

25% = 5

Median = 8

75% = 9

Max = 15



Shape

Shape

Skewness Kurtosis

Skewness

~~→ Tail is more towards right or left~~

→ It measures asymmetry of a distribution around mean.

* Symmetric means → value → around mean → evenly spread.

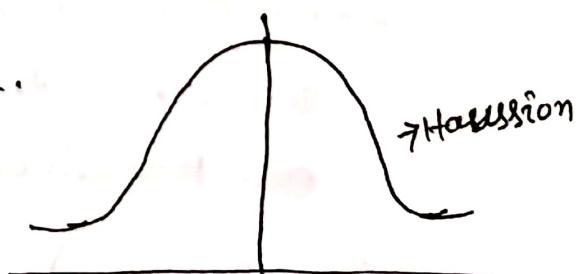
Why?

→ Understand the shape of data.

Ex

- Imagine people earn between 20K - 50K.

- A few people earn 2. lakhs +



Mean = median = mode
Symmetric

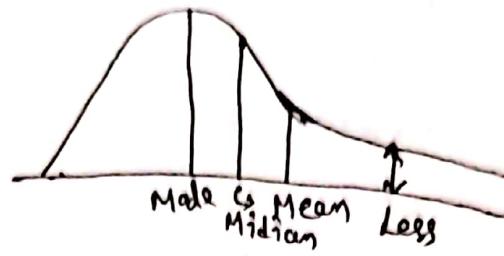
Types of skew

① +ve Skew

⇒ Tail is more towards right or +ve

⇒ Mean > median > Mode

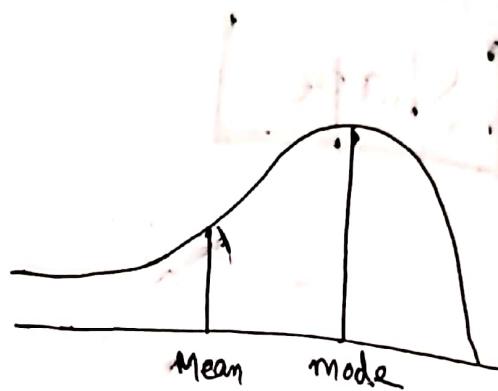
⇒ Skewness > 0



② -ve skew

⇒ Tail on the left side longer.

⇒ Mean < median < Mode



$$\text{Skewness} = \frac{n}{(n-1)(n-2)} \sqrt{\left[\frac{(x_i - \bar{x})^3}{s^3} \right]}$$

x_i - data pt

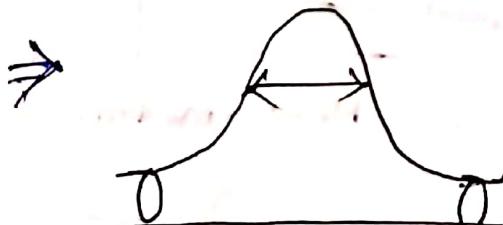
\bar{x} - mean

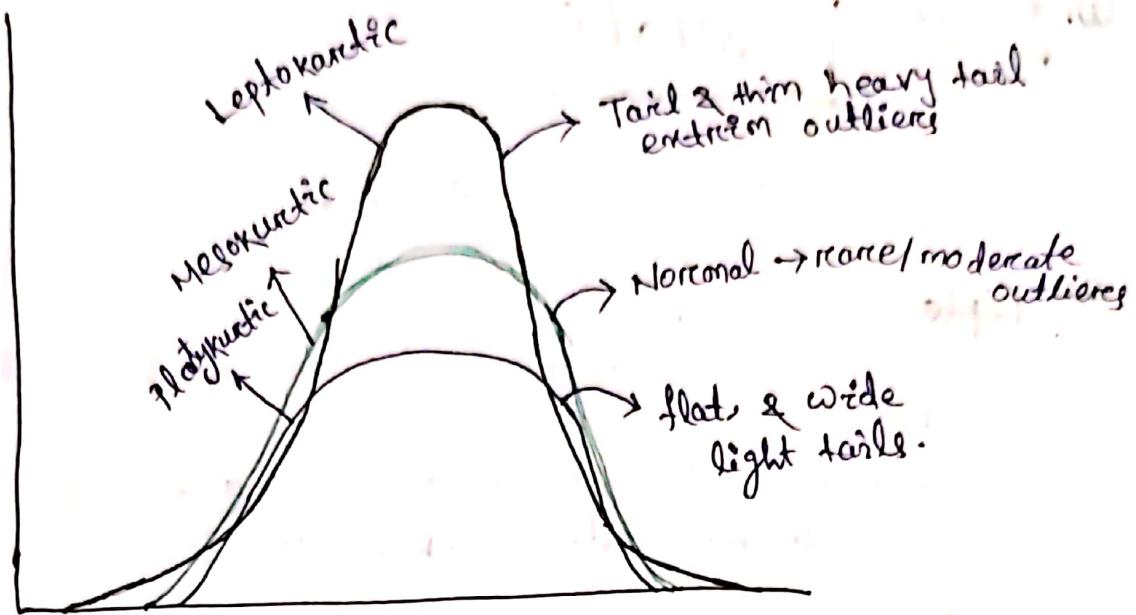
s - std

n - no of observations

Kurtosis

⇒ kurtosis is a statistical measure that describes the "tailedness" or "peakedness" of a probability distribution.





Platykurtic

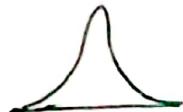
1. Mesokurtic

- ⇒ Distribution with normal kurtosis;
- ⇒ Moderate tails, similar to the normal distrib?



2. Leptokurtic

- ⇒ Distribution with high kurtosis ($Kurtosis > 0$)
- ⇒ Heavy tails and a sharp peak
- ⇒ most extreme outliers.
- ⇒ Higher probability of extreme values.



3. Platykurtic

- ⇒ Distribution with low kurtosis ($Kurtosis < 0$)
- ⇒ Light tails & a flatter peak - fewer outliers.



$$Kurtosis = \frac{n(n+1)}{(n-1)(n-2)(n-3)} \sum \frac{(x_i - \bar{x})^4}{S^4} - \frac{3(n-1)^2}{(n-2)(n-3)}$$

$K > 0 \rightarrow$ lepto

$K < 0 \rightarrow$ Platy

$K = 0 \rightarrow$ Meso

Why?

- To detect outliers.
- To understand ~~exist~~ in data.

Ex Lepto

- ① Stock returns → more chance of big profit or loss.
- ② Most students score 50 but a few get 0 or 100
→ extreme value → high kurtosis.

Platy

- ① Fewer extreme answers.
- ② Scores are more evenly spread between 40-60 low kurtosis.

4. #

Visual

- Statistics Visuals are graphical representations of data used to summarize, interpret and communicate information clearly and effectively.
- Jiske hume patterns, outliers, distribution Pata chalta hai.

① Histogram

② Bar chart

③ Pie chart

④ Box plot

⑤ Line chart

⑥ Scatter plot.

① Histogram \Rightarrow Univariate

\Rightarrow Histogram is the graphical representation of numeric variable / quantitative data.

\Rightarrow Use for continuous data & groups data into range.



why?

\Rightarrow To understand how data is spread out, patterns like skewness, central tendency or outliers.

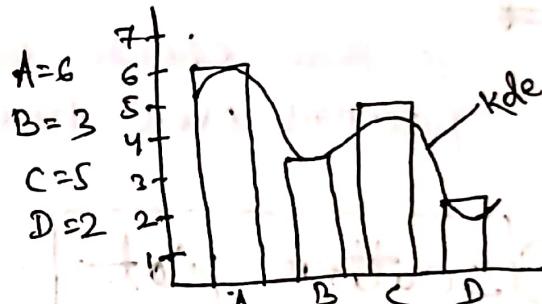


② Bar chart \Rightarrow Bivariate

\Rightarrow Bar chart represents a categorical data, and the length or height of the bar shows the value or frequency of that category.

why?

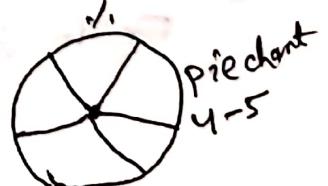
\Rightarrow Easy to understand, compare different categories or trends.



③ Pie chart

\Rightarrow A pie chart is a circular chart divided into slices, where each slice represents a percentage of a whole data.

why?



\Rightarrow Want to visualize percentage.

\Rightarrow easy to see which category is biggest or smallest.

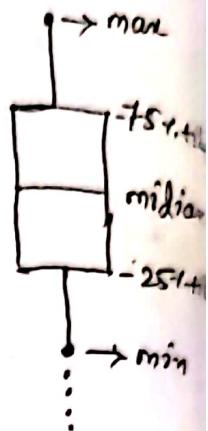
4. Box plot \Rightarrow univariate.

\Rightarrow Box plot shows the spread and skewness of numerical data.

\Rightarrow using five-number summary.

why?

\Rightarrow understand the spread and center of the data, to outlier detection.

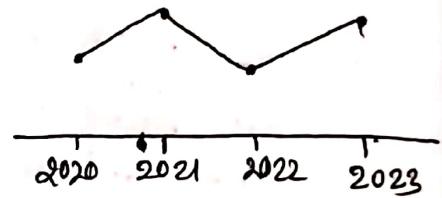


5. Line chart \Rightarrow Bivariate.

\Rightarrow Line chart shows trends over time by connecting data points with a line.

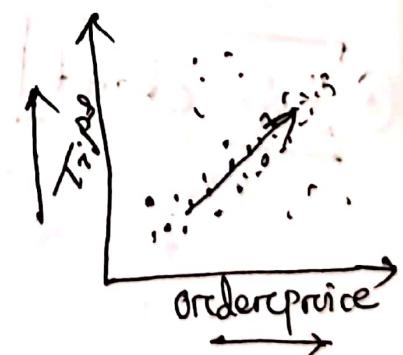
why?

To show change or trends patterns over time!



⑥ Scatter plot \Rightarrow Bivariate analysis tools.

\Rightarrow Scatter plot displays the CO- relationship between two numerical variables.



Why?

\Rightarrow To detect trends or outliers.

Distribution

(Probability Distribution)

Distribution

⇒ The possible values that a variable can take & how frequently they can occur.

* $y \rightarrow$ actual outcome

$y \rightarrow$ one of the possible outcome, (0-1)

Ex

$$\text{Coin} \begin{cases} H - \frac{1}{2} \\ T - \frac{1}{2} \end{cases}$$

$$\text{Dice} \begin{cases} 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{cases} n=6, \frac{1}{6}$$

$$* P(E) = P(y=3) = \frac{1}{6}$$

*

Quantitative

Discrete
(finite no.)

Continuous
(infinite no.)

follows:

$$X \sim N(\mu, \sigma)$$

variable

dist

characteristics

Why do we need distribution?

⇒ where it +

Discrete distribution \Rightarrow values that can be counted
(not measured)

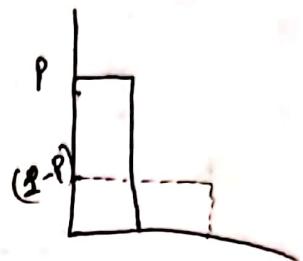
1. Bernoulli distribution

\Rightarrow It models a single experiment with two possible outcomes.

\Rightarrow 1 trial & outcomes.

\Rightarrow success ($x=1$), failure ($x=0$)
 $x \sim \text{Bern}(P)$

* P = Probability of success



$1-P$ = probability of failure [Binary outcome 0 or 1]

2. Binomial distribution

\Rightarrow It is a sequence of identical Bernoulli events.

$$x \sim B(n, p)$$

\Rightarrow fixed no of trials each with Bernoulli result.

Probability concept:

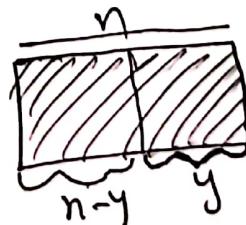
$$\cdot P(\text{desired outcome}) = p$$

$$\cdot P(\text{alternate outcome}) = 1-p$$

$$y \rightarrow P(y) = p^y$$

$$P(n-y) = 1-p^{n-y}$$

$$P(y) = nC_y \cdot p^y \cdot (1-p)^{n-y}$$



$$\text{PDD} \quad nC_y = \frac{n!}{(n-y)! y!}$$

3. Poisson distribution

⇒ Poisson distribution is a probability distribution that describes the number of events that occur within a fixed interval of time.

⇒ Count of events in fixed time, event happen independently, scarcely.

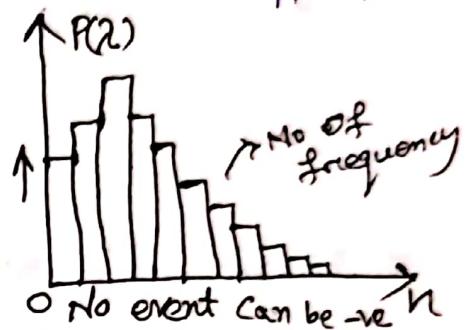
$$X \sim Po(\lambda)$$

interval

e.g. each day 4 question ans.. asked, ~~each day 7 question asked.~~

$$P(Y = 7) = \frac{\lambda^7 e^{-\lambda}}{7!} \quad e = \text{Euler's no.}$$

$$= \frac{4^7 \times e^{-4}}{7!} = \frac{4^7 \times 2.72^{-4}}{7!}$$



Continuous distribution

⇒ It shows how different values occur in data set & describes the overall pattern of the data.

I. Normal Distribution (Gaussian Distribution)

⇒ A normal distribution is a continuous probability distribution, that is symmetric around its mean, forming a bell-shaped curve.

It is one of the most important distributions in statistics because many natural and social phenomena follow it.

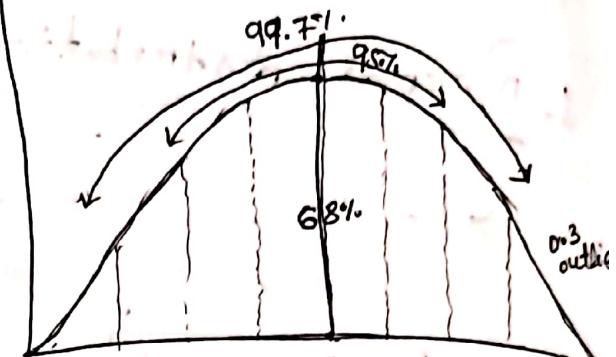
Empirical rule

$$\mu \pm 1\sigma \rightarrow 68\%$$

$$\mu \pm 2\sigma \rightarrow 95\%$$

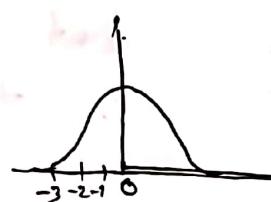
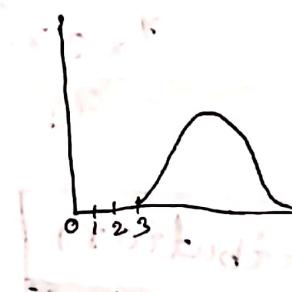
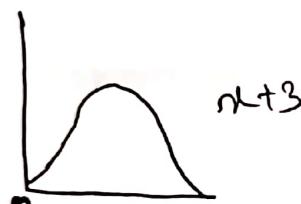
$$\mu \pm 3\sigma \rightarrow 99.7\%$$

0.03 → outliers

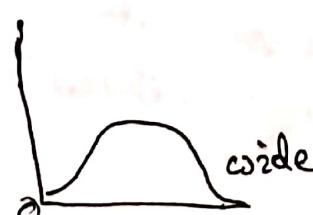
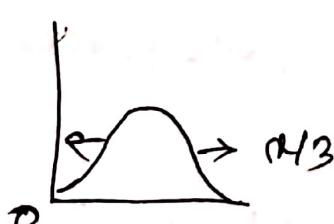


Mean = median = mode

* Distributⁿ ke uparne transformation bhi karega tain
hai : value add ya subtract karne. to distribution shift hoga.



shrink



Special kind of transformation jisko hum
Z-score kehatain hai !



Standardized \Rightarrow compare data to others / data point + others.

$$Z \text{ score} = \frac{x - \mu}{\sigma} \xrightarrow{\text{data point}} \begin{aligned} \text{Mean}(\mu) &= 0 \\ \text{SD}(\sigma) &= 1 \end{aligned}$$

- * Used to standardise a data point so you can compare it with others.

Mean = 0
+ve value = above avg
-ve value = below avg.

Ex

$$\mu = 70 \quad Z = \frac{85 - 70}{10} = \frac{15}{10} = 1.5$$

$$SND = 85 \quad M = 0 \quad SD = 1$$



2. Student t-distribution

what \Rightarrow Bell shape distribution similar to ND, but wider & flatter.



\Rightarrow use small sample

Sample SD \rightarrow known

Pop SD \rightarrow unknown

why

\Rightarrow Sample jab small hote hai, use hum inference
like ana hai!

$\Rightarrow n < 30$ ho tab hum isko use Karoain hai,

Inference means drawing a conclusion based on evidence.

3. Chi-square distribution (χ^2)

- The Chi-square distribution is the distribution of a sum of ~~two~~ or the squares.
- It is used in hypothesis testing, confidence interval estimation for categorical data.

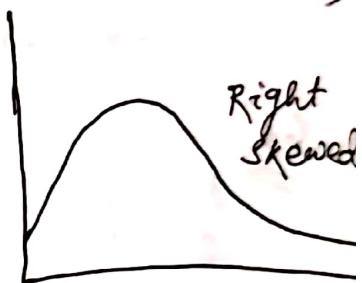
Use • Chi-square Test of independence (e.g. gender vs voting preference)

- The sketch shows a right skewed curve.

$$df = (r-1) \times (c-1)$$

r = no. of rows

c = no. of columns



Hypothesis testing

- Hypothesis Testing is a statistical method used to make decisions or inferences about a population based on sample data.
- It is a process where you try to gather the evidence to reject the null hypothesis and to prove that alternate hypothesis is true (accept).

- Hypothesis → kind of a statement
- Two type
 - null hypothesis
 - alternate hypothesis.

Null hypothesis (H_0)

- It is a kind of statement which says nothing new is happening.
- There is no significant relationship between two variables.

Alternate hypothesis (H_1 / H_a)

- There is significant relationship between two variables.
- 50-point estimate \rightarrow avg of all matches (sample)
 - 50 ± 5 [$\pm 5 \rightarrow$ margin of error]
 - $[45, 55] \rightarrow$ confidence interval.

confidence level $\Rightarrow (1 - \alpha)$

$\alpha =$ significance level

$$= 0.05 = 5\%$$

$$\alpha = 0.05$$

$$\text{confidence level } 1 - \alpha \\ = 1 - 0.05 = 0.95 = 95\%$$

- confidence level $\Rightarrow 90\% = 0.9 = [0.9, 0.95, 0.99]$
 $[40-60] = 0.99 \Rightarrow$ precision less, large range
[type-2]

$[45-55] = 0.95 \Rightarrow$ Moderate = mostly used

$[48-52] = 0.90 \Rightarrow$ precision more, orange
less \Rightarrow confidence less [type 1]

- Type I \Rightarrow It is a type of error where you are rejecting H_0 when it is actually true.

- Type 2 \Rightarrow It is a type of error where you are accepting H_0 when it is actually false.

- Loss of significance :-

It is the probability that we are taking the risk to reject H_0 when it is actually true, (type I), [1 - confidence level], [0.05]

Sampling Techniques

Types of Sampling

① Simple Random Sampling :-

- Every member has an equal chance of selection.
- Example :- Drawing names from a hat.

② Stratified Sampling :-

- Population is divided into subgroups
- and random samples are taken from each subgroups.

Example :- Age group

10 - 20 }
21 - 30 } male/female,
31 - 40 }

③ Systematic Sampling :-

- every k^{th} (like 3rd, 5th, 10th) member is selected from a list.

Example :-

choose 7th, 14th, 21th, 28th person in part time queue,

- Selecting "mth" number follows same

Z - TEST

A Z-test is a statistical test used to determine whether there is a significant difference between sample and population means (or between two sample means).

When the population variance is known and the sample size is large (typically $n \geq 30$)

Formula :-

pop

$$Z = \frac{\bar{x} - \mu}{\sigma}$$

Sample

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

When?

- The sample size is large ($n \geq 30$)
- The population standard deviation (σ) is known.
- Random Sampling [normalized]
- The data is approximately normally distributed (especially important for small samples)

* μ = population mean

σ = population std dev

\bar{x} = sample mean

n = sample size

P-Value	Kya karne ha?	Matlab
$P \leq 0.05$	H ₀ ko reject karo	H ₀ ko support karo
$P > 0.05$	H ₀ ko reject nahi kar sakta	H ₀ ko support nahi milta

T-TEST

T-test is a statistical test used to determine whether there is a significant difference between the means of two groups.

It is commonly used when:

- The sample size is small ($n < 30$)
- The population is unknown, σ ($\sigma_d = \text{unknown}$)

① One Sample t-test

Compares the sample mean to a known value (often a population mean).

Ex: Checking if a class's average score differs from the national average.

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

df = degree of freedom

$$df = n - 1$$

② Two Independent Sample t-test.

Compares the means of two independent groups.

Ex: Comparing test scores of two different classes.

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$t = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

$$df = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\frac{(s_1^2/n_1)^2}{n-1} + \frac{(s_2^2/n_2)^2}{n-2}}$$

③ Paired t-test (dependent)

Compares the means of two related groups.

Ex: Pre-test and post-test scores from the same students.

Formula:

$$t = \frac{\bar{d}}{S_d / \sqrt{n}}$$

\bar{d} → mean of the differences
 S_d → standard deviation of the differences

Chi-Square Test

The Chi-Square test (χ^2 test) is a statistical method used to test "relatedness" between categorical variables.

It helps answer questions like:

- Is this distribution what we expected?
- Are two categorical variables related or independent?

Why?

To compare observed data (what is collected) with expected data (what is theoretically expected).

$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

O = Observed frequency

E = Expected frequency

$$E = \frac{\text{Row total} \times \text{Column total}}{\text{Grand total}}$$

$$df = (r-1) \times (c-1)$$

r = row total

c = column total

Example of "Education" background & job satisfaction

Education	Satisfied	Not Satisfied	Total
High school	50	70	120
Graduation	90	60	150
Post graduation	20	10	30
	160	140	300

Of = observed frequency

Education	Satisfied	not satisfied.
High school	$\frac{120 \times 160}{300} = 64$	$\frac{120 \times 140}{300} = 56$
Graduation	$\frac{150 \times 160}{300} = 80$	$\frac{150 \times 140}{300} = 70$
Post graduation	$\frac{30 \times 160}{300} = 16$	$\frac{30 \times 140}{300} = 14$

ER = Expected frequency or Result

$$\chi^2 / HS/S = \frac{(O-E)^2}{E} = \frac{(50-64)^2}{64} = \frac{196}{64} = 3.06$$

$$HS/NS = \frac{(70-56)^2}{56} = \frac{196}{56} = 3.5$$

$$G/S = \frac{100}{80} = 1.25$$

$$G/NS = \frac{100}{70} = 1.43$$

$$PG/S = \frac{16}{16} = 1$$

$$PG/NS = \frac{16}{14} = 1.14$$

$$\chi^2 = 3.06 + 3.5 + 1 + 1.14 + 1.25 + 1.43$$

$$\boxed{\chi^2 = 11.38}$$

$$\begin{aligned}df &= (\alpha - 1) \times (c - 1) \\&= (3 - 1) \times (2 - 1) \\&= 2 \times 1 = 2\end{aligned}$$

$$\boxed{df = 2}$$

