

## What is Statistics

Statistics is a branch of Mathematics that involves collecting, analysis, interpreting and presenting data. It provides tools and methods to understand and make sense of large amount of data and to draw conclusions and make decisions based on the data.

In practice, statistics is used in a wide range of fields, such as business, economics, social science, medicine and engineering. It is used to conduct research studies, analyse market trends, evaluate the effectiveness of treatments and interventions and make forecasts and predictions.

### Example :

- ① Business → Data analysis (Identifying customer behaviour) and Demand forecasting.
- ② Medical → Identify efficacy of new medicines (clinical trials), Identifying risk factor for diseases.
- ③ Government and politics → Conducting surveys, Polling.
- ④ Environmental science → Climate research

## TYPES OF STATISTICS

### ① Descriptive statistics

↳ Descriptive statistics deals with the collection, organization, analysis, interpretation and presentation of data. It focuses on summarizing and describing the main features of set of data, without making inferences or predictions about the larger population.

### ② Inferential statistics

↳ Inferential statistics deals with making conclusions and predictions about a population based on a sample. It involves the use of probability theory to estimate the likelihood of certain events occurring, hypothesis testing to determine if a certain claim about a population is supported by the data, and regression analysis to examine the relationship between variables.

Population → Population is the entire group of individual

ex

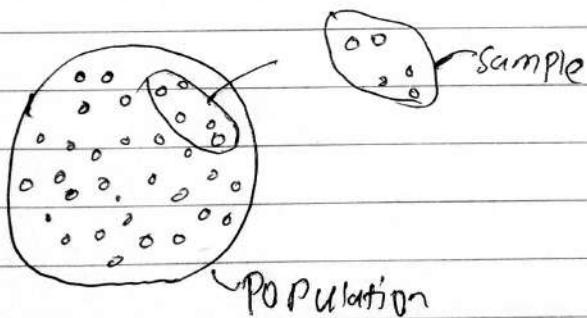
- ① All People of India
- ② All the User of Netflix.

Sample → A small subset of Population that represents the whole.

→ Used because studying the whole Population is often too costly or time consuming.

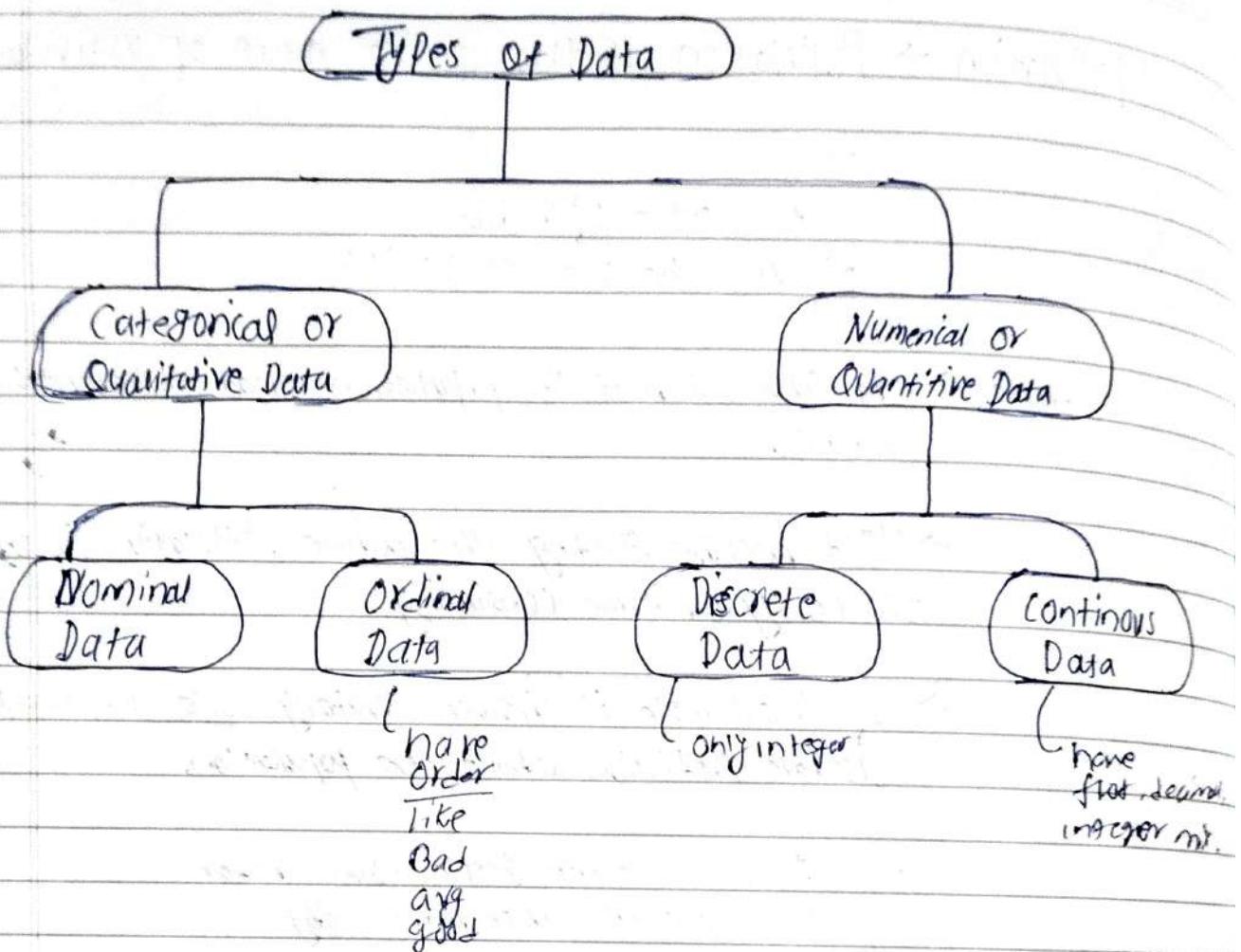
→ If the sample is chosen properly, we can make accurate prediction about the population.

- ex
- ① 100 random people from bitor
  - ② 50 Netflix user from fbg



Population → Whole group

Sample → small part



## Measure Of Central Tendency

→ A measure of central tendency is a statistical measure that represents a typical or central value for a dataset. It provides a summary of the data by identifying a single value that is most representative of the dataset as a whole.

### ① Mean

→ The mean is the sum of all values in the dataset divided by the number of values.

Population mean	Sample mean
$\mu = \frac{\sum_{i=1}^N x_i}{N}$ $N = \text{No of items in the population}$	$\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ $n = \text{no of items in sample}$

Population mean is represented by  $\mu$  (mu)

Sample mean is represented by  $\bar{x}$  (x bar)

### ② Median

→ The median is the middle value in the dataset when the data is arranged in order.

### ③ mode

→ The mode is the value that appears most frequently in the dataset.

### ④ Weighted mean

→ The weighted mean is the sum of products of each value and its weight, divide by the sum of the weights. It is used to calculate a mean when the data values in the dataset have different importance or frequency.

### ⑤ Trimmed mean

→ A trimmed mean is calculated by removing a certain percentage of the smallest and largest values from the dataset and then taking the mean of the remaining values. The percentage of values removed is called trimming percentage.

$$\{ 20000, 22000, 23000 \} \{ 21000, 28000, 30000, 32000, 35000, 35000, 80000 \}$$

$$\{ 36500 \}$$

$$\{ 25000, 28000, 30000, 32000, 38000 \}$$

$$\{ 30000 \}$$

full mean  
you had

trimmed  
mean

## Measure of Dispersion

→ A measure of dispersion is a statistical measure that describes the spread or variability of a dataset. It provides information about how the dataset is distributed around the central tendency (mean, median, mode) of the dataset.

### (1) Range

→ The range is the difference between the maximum and minimum values in the dataset. It is a simple measure of dispersion that is easy to calculate but can be affected by outliers.

Small range → Means all your values are close to each other - no much difference between min and max.

#### Interpretation

- └ Low variation
- └ Data is stable and predictable
- └ Indicate consistency.

Large range → Means values are far apart - big gap between min and max

#### Interpretation

- └ High variation
- └ Data is inconsistent and unpredictable
- └ Indicates diversity or instability.

② Variance → The variance is the squared differences between each data point and the mean. It measures the average distance of each point from the mean. It is useful in comparing the dispersion of datasets with different means.

	$n - \text{mean}$	$(n - \text{mean})^2$
3	3-3	0
2	2-3	1
1	1-3	4
5	5-3	4
4	4-3	1

$$\sigma^2 = \frac{\sum (n - \mu)^2}{N} \quad \begin{array}{l} \text{Population} \\ \text{variance} \end{array}$$

$$s^2 = \frac{\sum (n - \bar{n})^2}{n-1} \quad \begin{array}{l} \text{Sample} \\ \text{variance} \end{array}$$

### Mean absolute Deviation

$$MAD = \frac{\sum |n_i - \bar{n}|}{n}$$

③ Standard Deviation → The standard deviation is the square root of the variance. It is a widely used measure of dispersion that is useful in describing the shape distribution.

$$\sigma = \sqrt{\frac{\sum (x - \mu)^2}{N}} \quad \begin{array}{l} \text{Population} \\ \cancel{\text{SD}} \end{array}$$

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n-1}} \quad \begin{array}{l} \text{Sample} \\ SD \end{array}$$

④ Coefficient of Variation

↳ Coefficient of variation (CV) is the ratio of the SD to the mean expressed as a ~~percentage~~ percentage. It is used to compare the variability of dataset with different means and is commonly used in field such as biology, chemistry and engineering.

→ CV is a statistical measure that expresses the amount of variability in a dataset relative to the mean. It is a dimensionless quantity that is expressed as percentage.

formula of CV

$$CV = \text{Standard deviation / mean} \times 100\%$$

$$CV = \left[ \frac{\sigma}{M} \right] \times 100$$

Interpretation:

- It is expressed in Percentage
- Higher CV → more spread → More inconsistent data

↓  
CV yeh batata hai ke  
mean k respect me data

Kitne spread hai;

- Lower CV → less spread → More consistent data

## Graphs for univariate analysis

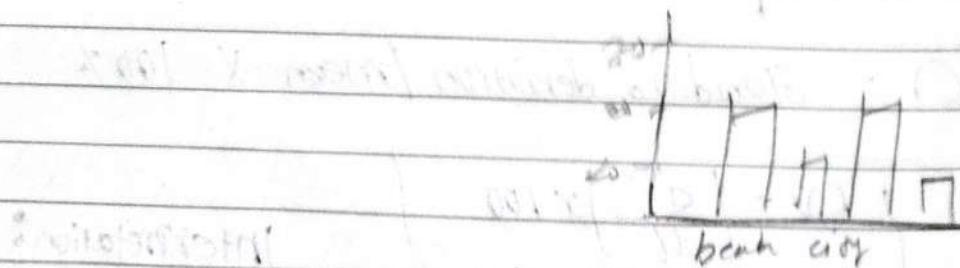
- ① Categorical - frequency Distribution table and Cumulative frequency.

A frequency distribution table is a table that summarizes the number of times (or frequency) that each values occur in dataset.

lets say we have a survey of 200 people & we ask them about their favourite types of vacation, which could be one of six categories: Beach, city, Adventure, Nature, cruises or other.

Type Of Vacation	
Beach	40
City	60
Adventure	30
Nature	35
Cruise	20
Other	15

Bar



## Quantiles

→ Quantiles are statistical measure used to divide a set of numerical data into equal sized group , with each group containing an equal Number of observations.

Quantiles are important measure of Variability and can be used to, Understand distribution of data , summarize and compare different datasets. They can also be used to identify outliers.

There are several types of quantiles used in statistical analysis including :

① Quartiles → Divide the Data into four equal parts , Q<sub>1</sub> (25<sup>th</sup> percentile), Q<sub>2</sub> (50<sup>th</sup> percentile) and Q<sub>3</sub> (75<sup>th</sup> percentile)

② Deciles → Divide the Data into 10 equal parts , D<sub>1</sub> (10<sup>th</sup> percentile), D<sub>2</sub> (20<sup>th</sup> percentile)..... D<sub>9</sub> (90<sup>th</sup> percentile).

③ Percentile → Divide the data into 100 equal parts , P<sub>1</sub> (1<sup>st</sup> percentile) P<sub>2</sub> (2<sup>nd</sup> percentile) .... P<sub>99</sub> (99<sup>th</sup> percentile).

④ quantiles → Divide the Data into ~~equal~~ s equal parts

Things to remember while calculating these measures:

- ① Data Should be Sorted from low to high
- ② You are basically finding the location of an observation
- ③ They are not actual Values in the Data
- ④ All other tiles can be easily ~~be~~ derived from Percentiles

Percentile → A Percentile is a statistical measure that represent the percentage of observation in a dataset that fall below a particular value.

e.g. 75th percentile is the value below which 75% of the observation in the dataset fall.

formula to calculate percentile value:

$$PL = \frac{P}{100} \cdot (N+1)$$

Where

P = the desired value/location

N = the total no of observation in the dataset

P = the percentile rank (expressed as a percentage)

Example

find the 75th percentile score for the below data.

78, 82, 84, 88, 91, 93, 94, 96, 98, 99

Step1 : sort the data (ASC)

78, 82, 84, 88, 91, 93, 94, 96, 98, 99

$$PL = \frac{75}{100} (10+1) = \frac{3}{4} \times 11 = \frac{33}{4} = 8.25$$

96 - 98

$$96 + 0.25 (2) = 96.5$$

75th percentile = 96.5

## Percentile of a value

$$\text{Percentile} = \frac{n + 0.5Y}{N}$$

$n$  = No of value below the given value

$Y$  = No of value equal to the given value

$N$  = total no of values in the dataset.

Ex

78, 82, 84, 88, 81, 93, 94, 96, 88, 89

find percentile

$$\text{Percentile} = \frac{3 + 0.5 \times 1}{10} = \frac{3.5}{10} = 0.35 \quad [35]$$

## 5 Number Summary

The five number summary is a descriptive statistic that provide a summary of a dataset. It consist of five values that divide the dataset into four equal parts, also known as quartiles. The five-number summary includes the following values:

- ① Minimum value → The smallest value in the dataset
- ② first quartile ( $Q_1$ ) → The value that separate the lowest 25% of the data from the rest of the data.
- ③ median ( $Q_2$ ) → The value that separates the lowest 50% from the highest 50% of the data.
- ④ Third quartile ( $Q_3$ ) → The value that separate the lowest 75% of the data from the highest 25% of the data.
- ⑤ Maximum value → The largest value in the dataset.

The five number summary is often represented visually using a box plot which display the range of the dataset, the median and the quartiles. The five number summary is a useful way to quickly summarize the center tendency, variability and distribution.

Minimum	Lower quartile	Median	Upper quartile	Maximum
25%	25%	25%	25%	

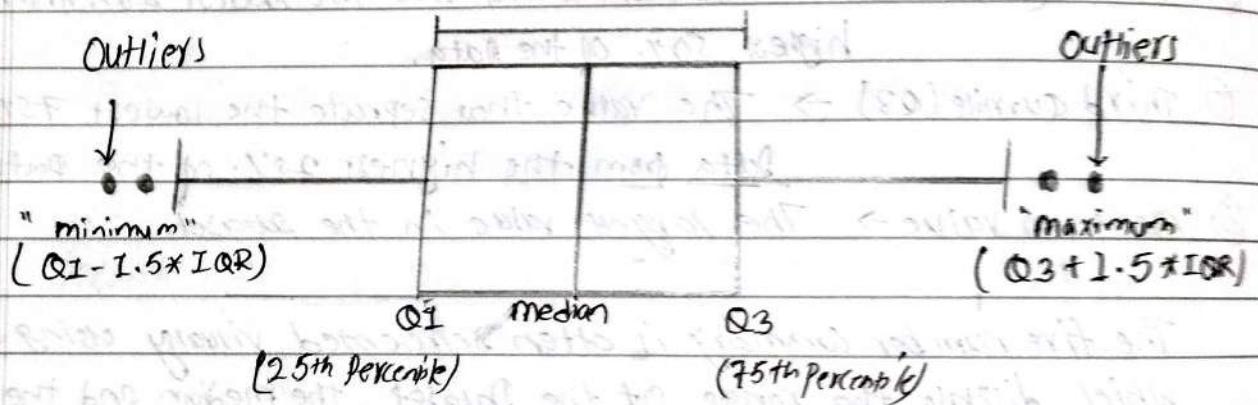
### Interquartile Range (IQR)

- ↳ IQR is a measure of variability that is based on five number summary of a dataset. Specifically, the IQR is defined as the difference between third quartile ( $Q_3$ ) and the first quartile ( $Q_1$ ) of a dataset.

## What is a box plot?

→ A box plot, also known as box-and-whisker plot, is a graphical representation of a dataset that shows the distribution of the data. The box plot displays a summary of the data, including the minimum and maximum value, the first quartile (Q1), the median (Q2), and the third quartile (Q3).

Interquartile Range  
(IQR)



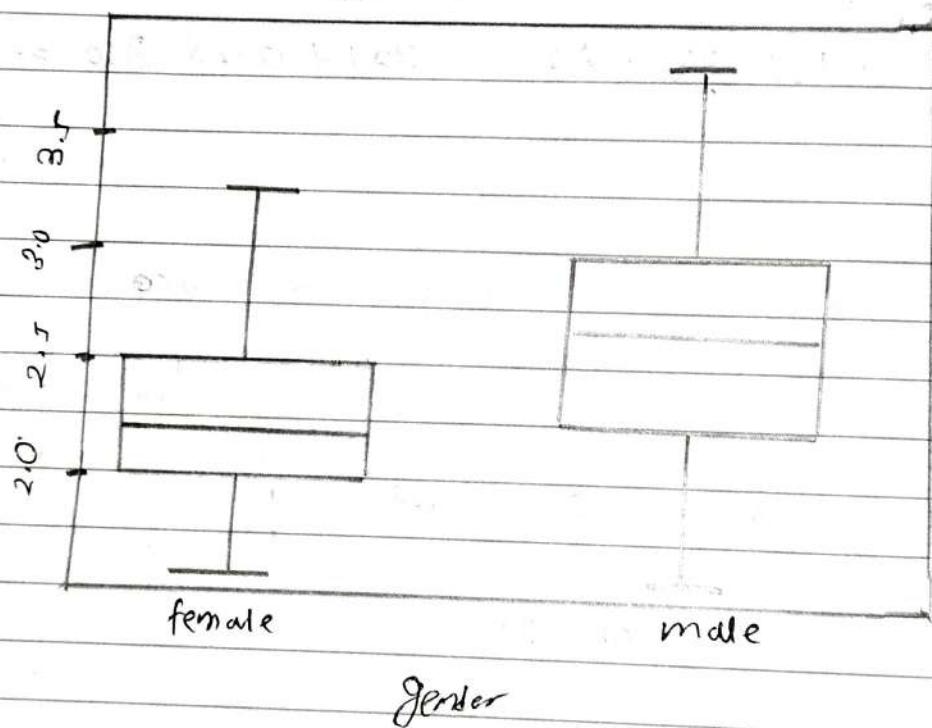
Minimum - Maximum both values  
Value outliers have

## Benefits of a box plot

- ↳ easy way to see distribution of data
- ↳ Tells about skewness of data
- ↳ can identify outliers
- ↳ compare 2 category of data

## ② Side by side box plot.

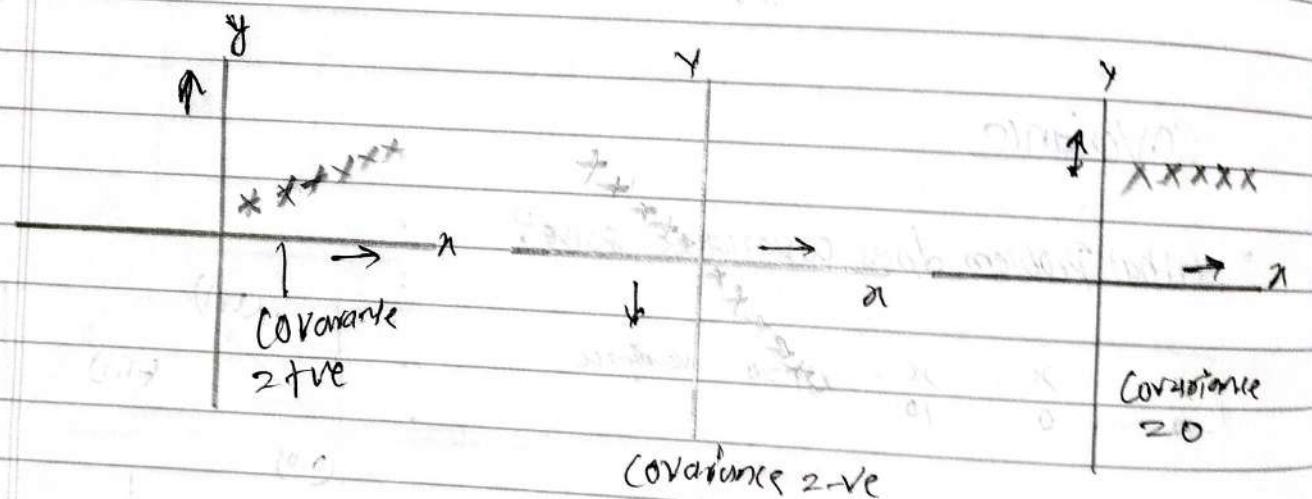
Cost weight by gender



- ~~Ques~~ What is Covariance and How it is interpreted?

→ Covariance is a statistical measure that describes the degree to which two variables are linearly related. It measures how much two variables change together, such that when ~~DATA~~ one variable increase, does the other variable also increase, or does it decrease?

If the covariance between two variables is positive, it means the variables tend to move together in the same direction. If covariance is negative, it means the variables tend to move in the opposite directions. A covariance of zero indicates that the variables are not linearly related.



covariance will not give an answer) & others  
↳ linear relation) interpretation,  
↳ need to understand ins and outs w.r.t.  
↳ statistics for stats with us

How is it calculated?

### Covariance formula

Population

Sample

$$\sigma_{xy} = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{N}$$

$$s_{xy} = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

$x_i, y_i \rightarrow$  The value of  $x$  &  $y$  in the Population  $\bar{x}, \bar{y} \rightarrow$  The value of  $x$  &  $y$  in the sample

$\mu_x, \mu_y \rightarrow$  The Population mean of  $x$  and  $y$ .  $\bar{x}, \bar{y} \rightarrow$  The sample mean of  $x$  and  $y$ .

$N \rightarrow$  Total no of Observation

$n \rightarrow$  Total number of observation.

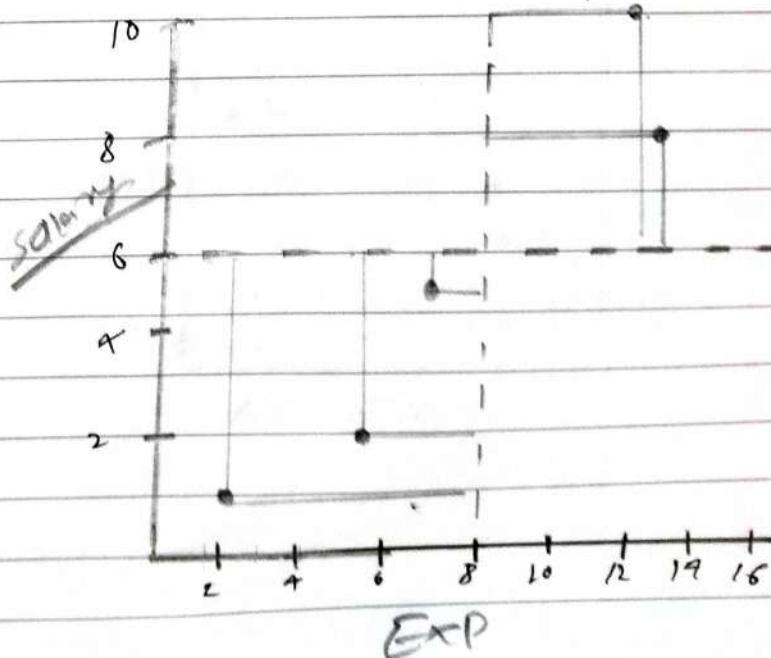
Emp (x)	Salary (y)	$x - x\text{mean}$	$y - y\text{mean}$	$(x - x\text{mean})(y - y\text{mean})$
2	1	-6	-5	30
5	2	-3	-4	12
8	5	0	-1	-1
12	12	4	6	24
13	10	5	-4	20

$$\bar{x} = 8 \quad \bar{y} = 6$$

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

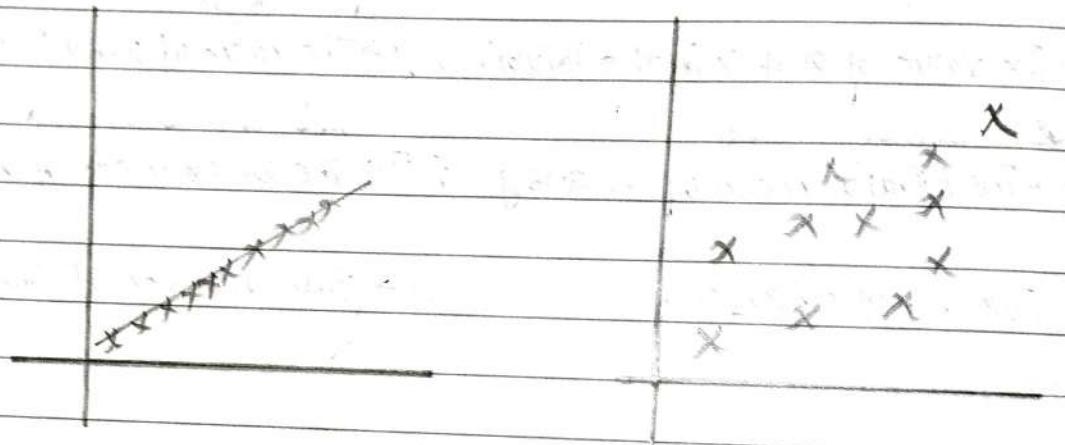
$$\frac{85}{9}$$

$\text{cov} = 21.5$



## Disadvantage of Using Covariance

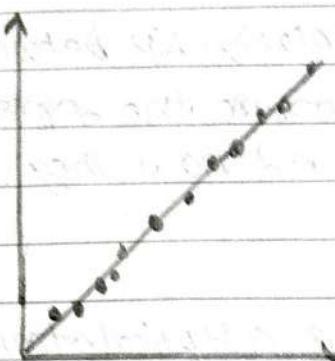
One limitation of Covariance is that it does not tell us about the strength of relationship between two variable, since the magnitude of Covariance is affected by the scale of the Variable.



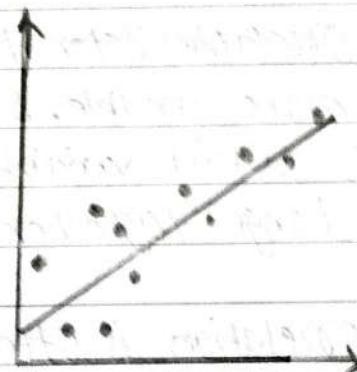
Note → Dono the hai katt ek ka relationship  
bohot accha hai aur ek sirf the hai.

## Correlation

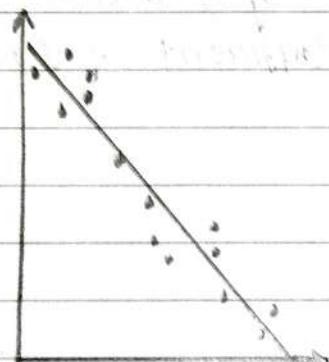
① What problem does the Correlation solve?



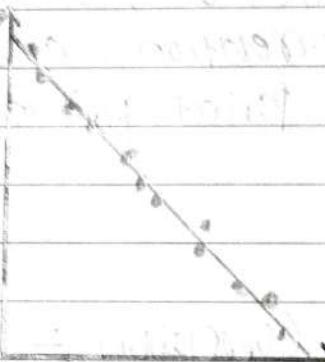
Strong Positive  
Correlation



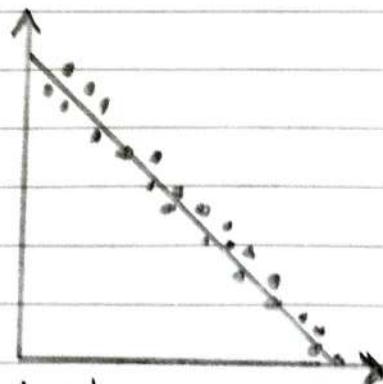
Weak Positive  
Correlation



Weak negative  
Correlation



Strong negative  
Correlation



Moderate Negative  
Correlation



No correlation

## ② What is Correlation?

→ Correlation refers to a statistical relationship between two or more variables. Specifically, it measures the degree to which variables are related and how they tend to change together.

Correlation is often measured using a statistical tool called the correlation coefficient, which ranges from -1 to 1.

A correlation coefficient of -1 indicates a perfect negative correlation, a correlation coefficient of 0 indicates no correlation, and a correlation coefficient of 1 indicates a perfect positive correlation.

$$\text{Correlation} = \frac{\text{Cov}(x, y)}{\sigma_x \times \sigma_y}$$

+1 = Perfect Positive Correlation

0 = No Correlation

-1 = Perfect Negative Correlation

## Correlation and Causation

The phrase "correlation does not imply causation" means that just because two variables are associated with each other, it does not necessarily mean that one causes the other. In other words, a correlation between two variables does not necessarily imply that one variable is the reason for the other variable's behaviour.

Suppose there is a positive correlation between the number of firefighters present at a fire and the amount of damage caused by the fire. One might be tempted to conclude that the presence of firefighters causes more damage. However, this correlation could be explained by a third variable - the severity of the fire. More severe fires might require more firefighters to be present, and also cause more damage.

Thus, while correlation can provide valuable insights into how different variables are related, they cannot be used to establish causality. Establishing causality often requires additional evidence such as experiments, randomized controlled trials, or well-designed observational studies.

## Random Variable

- What is Algebraic variable?

In algebra a variable is like  $x$ , is an unknown value.

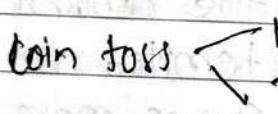
$$x + 5 = 10$$

$$\underline{x = 5}$$

- What is Random variables in statistic and Probability?

↳ A random Variable is set of Possible values from a random experiment.

coin toss



Dice

$$\begin{cases} 1 \\ 2 \\ 3 \end{cases} - \begin{cases} 4 \\ 5 \\ 6 \end{cases}$$

$$X = \{1, 0, 3\}$$

$$Y = \{1, 2, 3, 4, 5, 6\}$$

$$H = \{1, 2, 3\}$$

Sample Space

① What are Probability Distribution.

→ A Probability Distribution is a list of all possible outcome of random variable along with their corresponding Probability values.

comes	H	T
Probability	$\frac{1}{2}$	$\frac{1}{2}$

1 dice

1	2	3	4	5	6
$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$

2 dice

	1	2	3	4	5	6
1	2	3	4	5	6	7
2	3	4	5	6	7	8
3	4	5	6	7	8	9
4	5	6	7	8	9	10
5	6	7	8	9	10	11
6	7	8	9	10	11	12

$$\{ 2 \rightarrow 1/36$$

$$3 \rightarrow 2/36$$

$$4 \rightarrow 3/36$$

$$5 \rightarrow 4/36$$

$$6 \rightarrow 5/36$$

$$7 \rightarrow 6/36$$

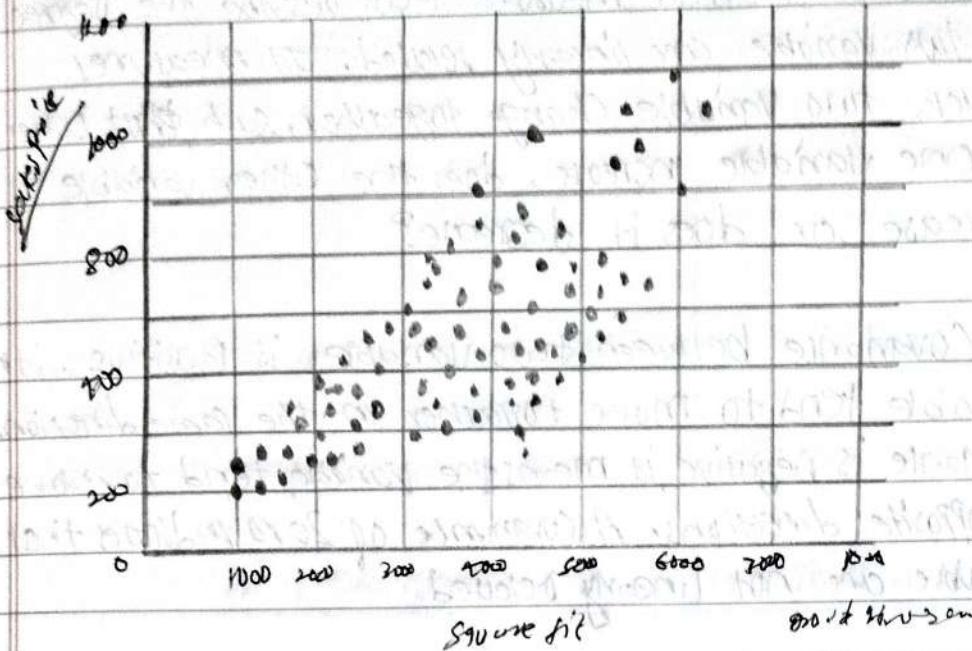
$$8 \rightarrow 5/36$$

$$9 \rightarrow 4/36$$

$$10 \rightarrow 3/36$$

$$11 \rightarrow 2/36$$

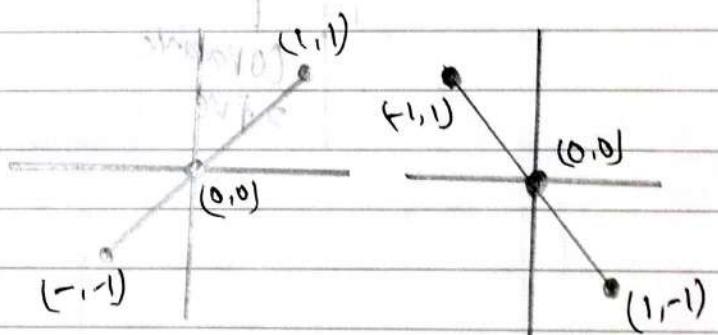
$$12 \rightarrow 1/36$$

Scatter PlotCovariance

- What problem does covariance solve?

$$\begin{array}{ccc} x & x & x \\ 10 & 0 & 10 \end{array} \quad m=0 \quad \text{variance}$$

$$\begin{array}{ccc} x & x & x \\ -20 & 0 & 20 \end{array} \quad m=0$$



$$\frac{1^2 + 0^2 + 1^2}{2} = \frac{2}{2} = 1$$

$$\frac{1^2 + 0^2 + 1^2}{2} = \frac{2}{2} = 1$$

Note → Covariance only tells the direction of relationship (Positive or Negative)

But not the strength, because it depends on the scale of variable.

## Problem with Distribution

- In many scenarios, the number of outcomes can be much larger and hence a table would be tedious to write down. Worse still, the number of possible outcomes could be infinite, in which case, good luck writing a table for that.

Example → Height of people, rolling 10 dice together.

Solution → function?

- What if we use a mathematical function to model the relationship between outcome and probability?

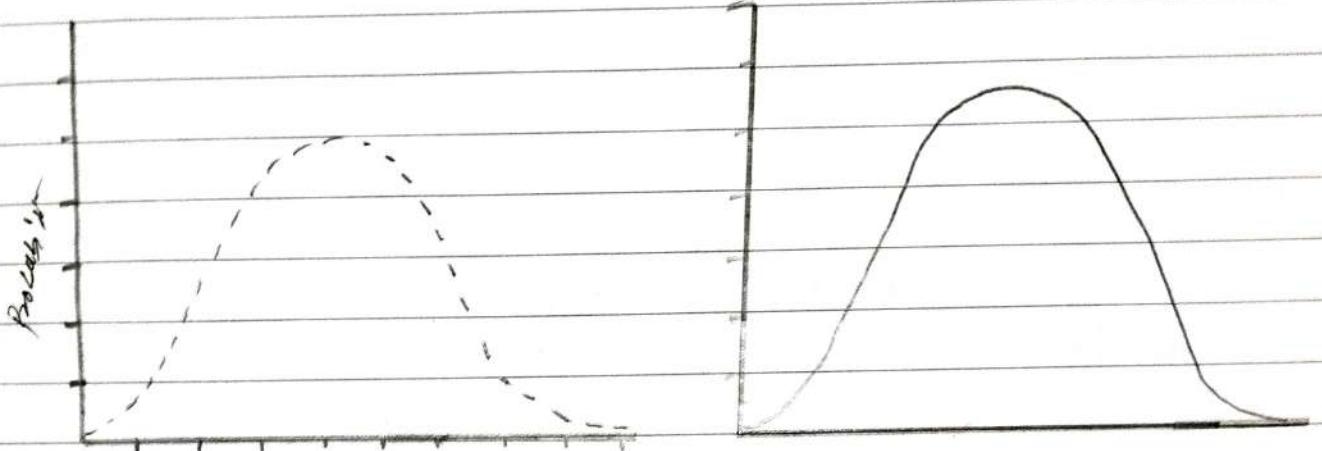
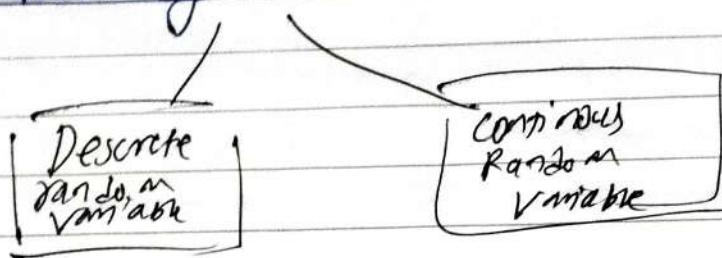
$n \rightarrow$  outcome

$y \rightarrow$  Probability

$$y = f(n)$$

**Note** → A lot of time probability distribution and probability distribution functions are used interchangeably.

## Types of Probability Distribution (PDF)



## Why are Probability Distribution important ?

- ↳ Gives an idea about the shape / distribution of the data.
- ↳ And if our data follows a famous distribution we automatically know a lot about the data.

## Note on a parameter

- Parameter in probability distribution are numerical values that determine the shape, location, and scale of the distribution.
- Different probability distributions have different sets of parameters that determine their shape and characteristics, and understanding these parameters is essential in statistical analysis and inference.

## Probability Distribution function

A Probability distribution function is a mathematical function that describes the probability of obtaining different values of random variable in a particular Probability distribution.

$y = f(x)$

Probability  
Mass function (PMF)

Probability density  
function (PDF)

Cumulative  
Distribution function (CDF)

PMF  $\rightarrow$  CDF  
PDF  $\rightarrow$  LDF

## Probability Mass function (Pmf)

Pmf stands for Probability mass function. It is a mathematical that describe the probability distribution of a discrete random variable.

The Pmf of discrete random variable assigns a probability to each possible value of the random variable. The Probabilities assigned by the Pmf must satisfy two conditions:

- ① The probability assigned to each value must be non-negative (ie., greater than or equal to zero).
- ② The sum of the probabilities assigned to all possible value must equal 1.

$$Y = f(n) \quad Y = \begin{cases} \frac{1}{8} & \text{if } n \in \{1, 2, 3, 4, 5, 6\} \\ 0 & \text{otherwise} \end{cases}$$

$$\text{Pmf} \rightarrow Y = \begin{cases} \frac{1}{36} & \text{at } \{2, 12\} \\ \frac{1}{36} & \text{at } \{3, 11\} \\ 0 & \text{otherwise} \end{cases}$$

Note  $\Rightarrow$  Pmf se ek discrete random variable ke particular value ki probability nikal sakte hain.

## Cumulative Distribution function (CDF) of Pmf

→ The Cumulative distribution function (CDF)  $f(x)$  describes the Probability that a random variable  $X$  with a given probability distribution will be found at a value less than equal to  $x$ .

$$F(x) = P(X \leq x)$$

Discrete/continuous

Note → CDF batata hai ke random variable ke kisi value se chotta ya uske barabar value aane ki total probability tya hai

### Example

Agar ek fair dice hai  $\rightarrow x = \text{number on dice}$ .

Possible values : 1, 2, 3, 4, 5, 6

$$\text{Pmf} : P(X=x) = \frac{1}{6}$$

Toh CDF hoga :

$$\bullet f(1) = P(X \leq 1) = \frac{1}{6}$$

$$\bullet f(2) = P(X \leq 2) = \frac{2}{6}$$

$$\bullet f(3) = P(X \leq 3) = \frac{3}{6}$$

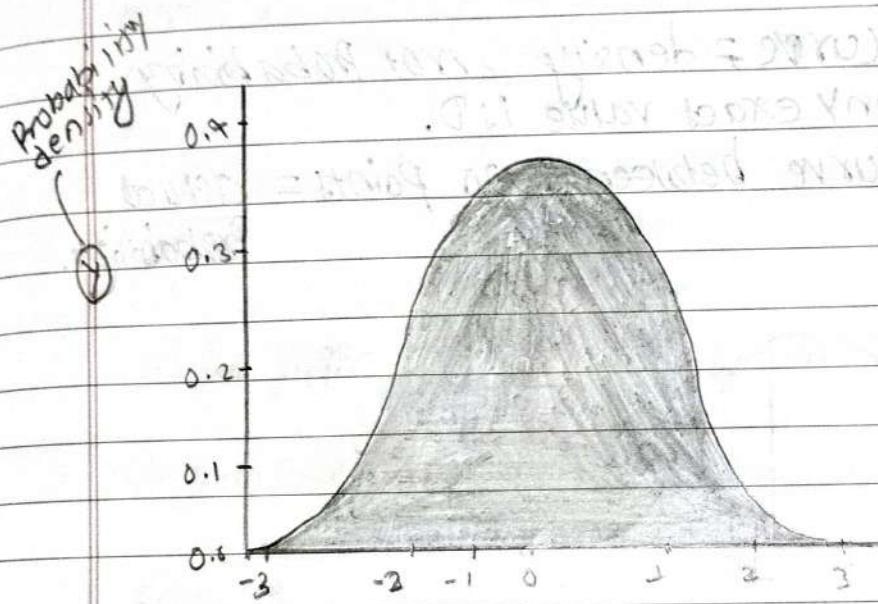
$$\bullet f(6) = P(X \leq 6) = 1$$

### Marks

→ CDF batata hai ke dice ka number 3 ya usse chotta value aane ka chance =  $3/6 = 0.5 (50\%)$

## Probability Density function

→ PDF Stands for probability Density function. It is a mathematical function that describe the probability distribution of a continuous random variable.



Why "Probability Density" not just "Probability"?

↳ Probability ke liye Probability Density use karte hai kyuki continuous ma ek exact value ka probability 0 hota hai. Jaise agar humne ek particular <sup>like 2.5</sup> value ka probability chahiye, aur 1 aur 2 ke beech infinite values hain to.

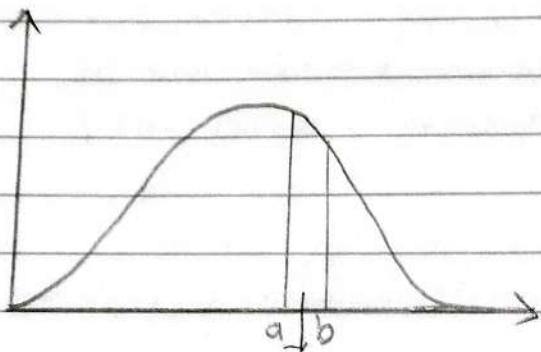
$$\frac{1}{\infty} = 0$$

Isleya ham Probability nahi, Probability Density use karte hain.

- What is Probability Density?

Probability Density  $\Rightarrow$  ek function (curve) jo batata hai kisi value ke ass pass Probability kitne "packed" hua

- Height of the curve = density, not probability
- Probability of any exact value is 0.
- Area under curve between two points = actual Probability.



$$P(a \leq X \leq b)$$

In another  
words

$\Rightarrow$  Probability Density. Describe how dense the Probability is at a certain value of a random variable, but not the actual Probability at a single point. For continuous variable, the Probability at any exact value is zero; instead you calculate the Probability for a range by finding the area under the curve between two points.

Note  $\rightarrow$  Probability Density batata hai ke 2 value ke bech na koi else value come ka probability kya hoga.

- How to calculate Probability then?

Ham koi ek exact point ka Probability nikalenge to n o hoga

isega ham PDF nikalte hai jo batata hai ke 2 points ke  
bacha mazaana ka ek no ka probability kya

$$\text{PDF formula: } P(a < X \leq b) = \int_a^b f(x) dx$$

- $f(x)$  ke jagha uska formula laga do (like distribution hai, Uniform, normal/exponential etc.).
- Integral solve karo  $\rightarrow$  Jo number aayega, wahi probability hai

### Example

Uniform [1, 5], range 2-3:

$$f(x) = 0.25 \Rightarrow P(2 \leq X \leq 3) = \int_2^3 0.25 dx = 0.25 \times (3-2) = 0.25$$

yahi hain Probability  $\Rightarrow 0.25$

- How graph is Calculated?

↳ Jaw Data hai.

↳ Parha nahi kavasa distribution follow karta hai

↳ Density estimation use karte hai taaki Probability Distribution ka shape sumajhi aaye

~~Design~~

## Density estimation

Density estimation ka matlab hai data se probability distribution ka estimate bannana. Matlab, thumore pass sirf data points hain, aur tumhe ~~is~~ nahi pata ke data konsa theoretical distribution follow karta hai (like Normal, Uniform, exponential etc). Toh ham ek curve banate hoga jo batya ki data probability kaha zada hai aur kaha kam.

Density ~~estimation~~ → raw data se distribution ka shape nikalne ka tanka.

- Density estimation can be used for a variety of purposes such as hypothesis testing, data analysis, and data visualization. It is particularly useful in areas such as machine learning, where it is often used to estimate the probability distribution of input data or model the likelihood of certain events or outcomes.

There are various methods for density estimation, including parametric and non-parametric approaches. ~~parametric~~

Commonly used techniques for density estimation include kernel density estimation (KDE), histogram estimation, Gaussian mixture model (GMMs). The choice of method depends on the specific characteristics of the data and the intended use of the density estimate.

## Parametric Density Estimation

→ Parametric density estimation is a method of estimating the probability density function (pdf) of a random variable by assuming that the underlying distribution belongs to a specific parametric family of probability distributions, such as normal, exponential or Poisson distribution.

### Explanation:

→ Parametric Density Estimation Steps :

- Data ko observe Karen aur assume Karen ki "Yeh kisi parametric distribution (e.g. normal distribution) ko follow karta hai"  $\rightarrow$  Plot histogram to assume distribution.
- Sample data se parameters calculate Karen, for example mean ( $\mu$ ) aur standard deviation ( $\sigma$ ) agar normal distribution assume kar raho ho
- Uske baad app us parametric function ko use karte data ke liye probability density values nikal sakte hai.

$$\text{Jaise Normal Distribution ke liya } f(x; \mu, \sigma) = \frac{1}{\sigma \sqrt{2\pi}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Value  $x$   
liye

- In values KO Plot Karen to visualize the pdf.

## Non-Parametric Density Estimation

But sometimes the distribution is not clear, it's not one of the famous distributions.

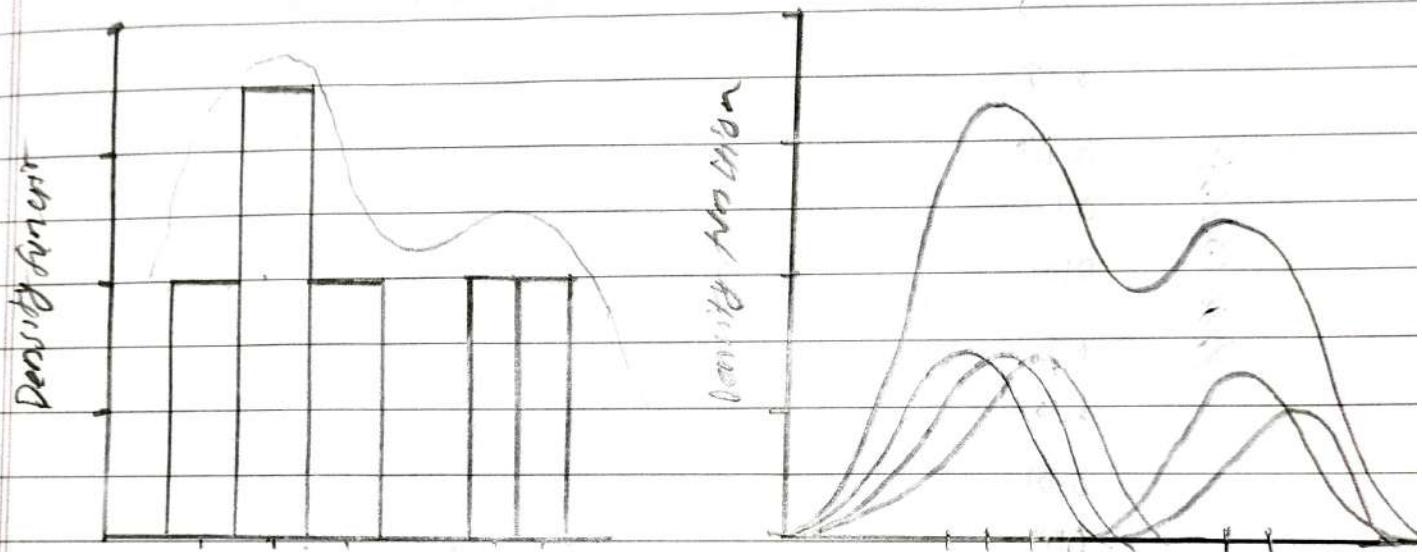
Non-Parametric density estimation is a statistical technique used to estimate the PDF of a random variable without making any assumption about the underlying distribution. It is also referred to as non-parametric density estimation because it does not require the use of predefined probability distribution function, as opposed to parametric methods such as the gaussian distribution.

The Non-Parametric density estimation techniques involves constructing an estimate of the PDF using the available data. This is typically done by creating a kernel density estimate.

Non Parametric density estimation has several advantages over Parametric density estimation. One of the main advantages, which allows for more flexible and accurate estimation in situations where the underlying distribution is unknown or complex. However, non-parametric density estimates compared to parametric methods,

## Kernel Density Estimate (KDE)

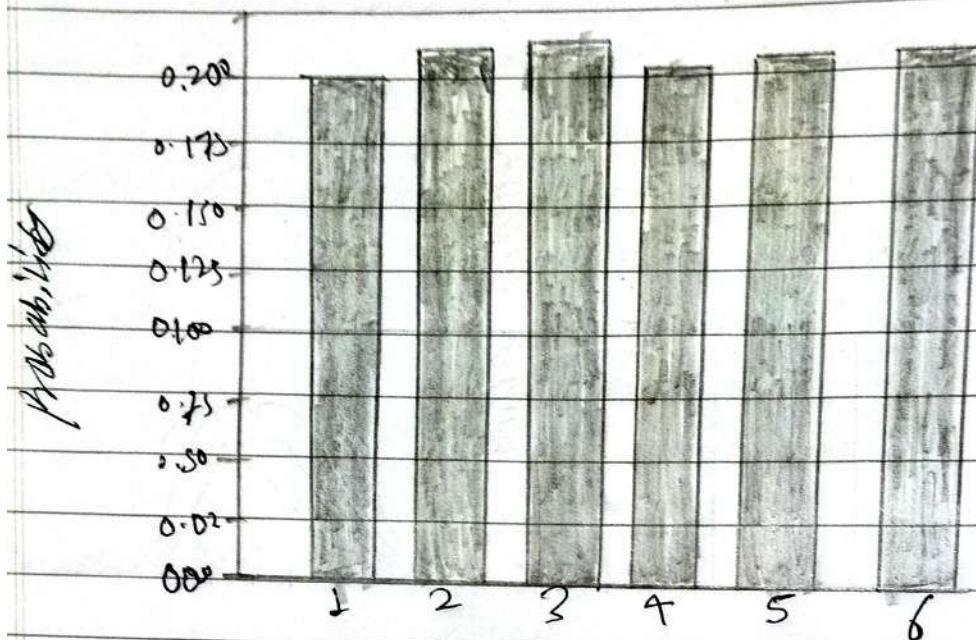
→ The KDE technique involves using a kernel function to smooth out the data and create a continuous estimate of the underlying density function.



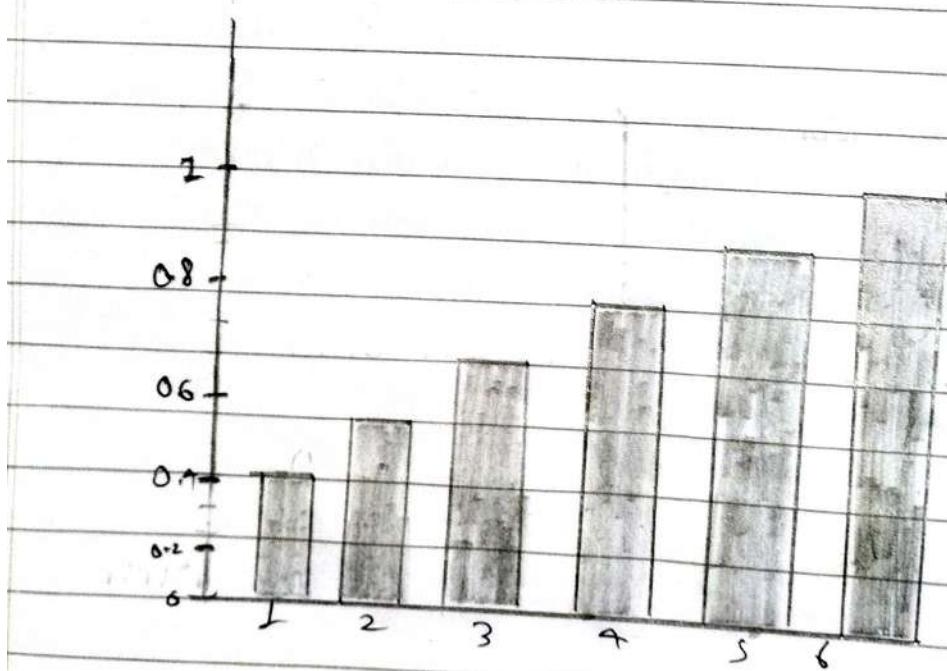
Bandwidth Jithna Jack hoger vtha smooth hoger

## Cumulative Distribution function function (CDF of PDF)

Rolling a Dice and calculate PDF and CDF.

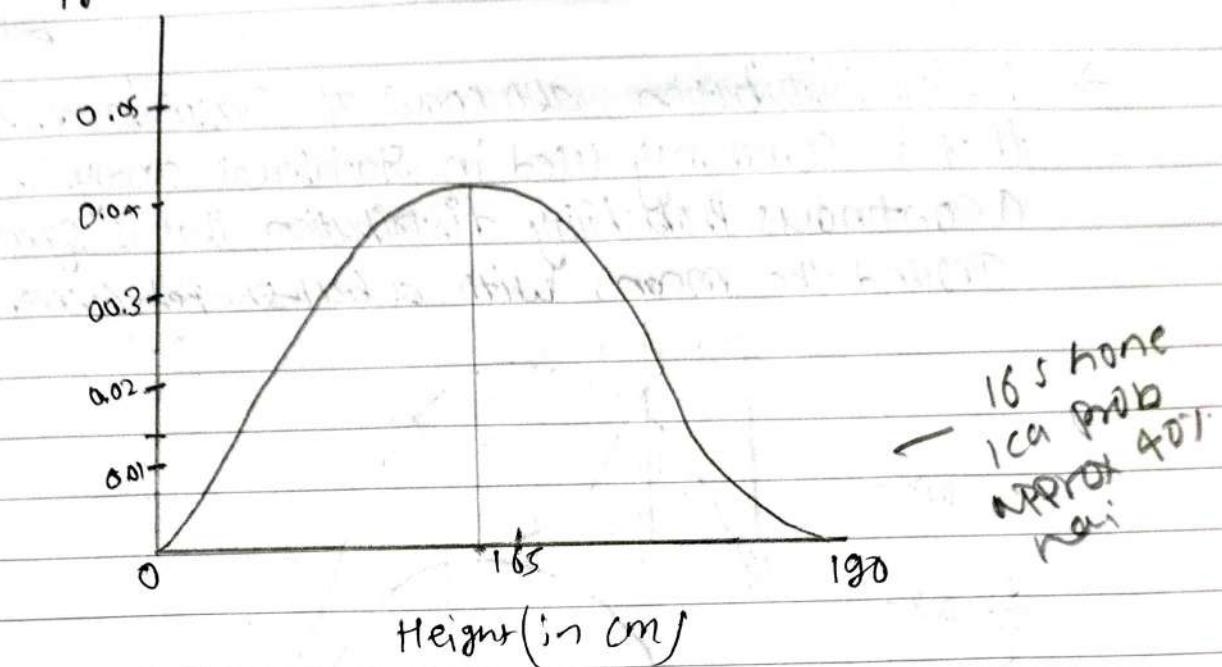


{ PMf of dice }  
Discrete Variable

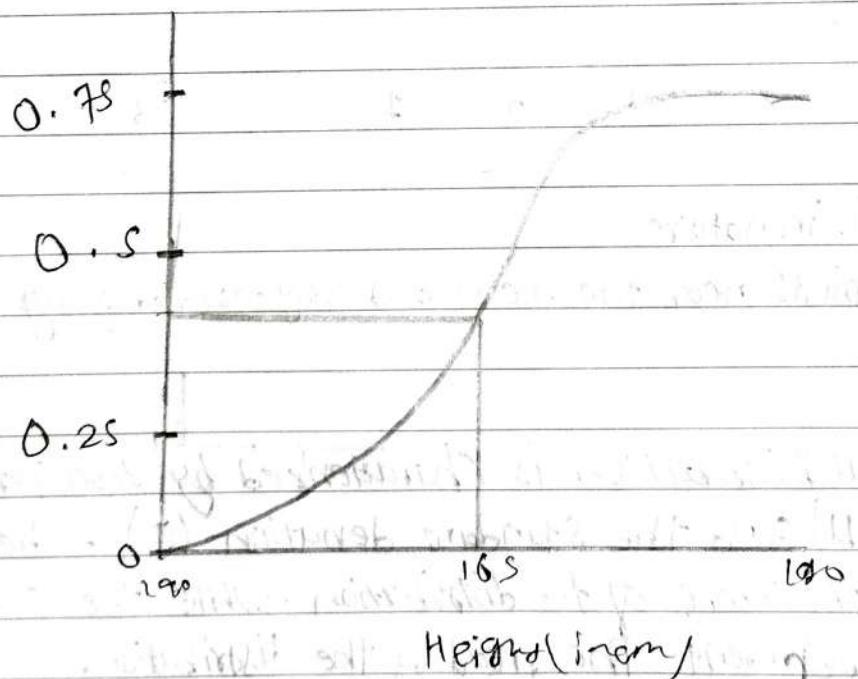


CDF of Dice and discrete random variable.

## Probability Density



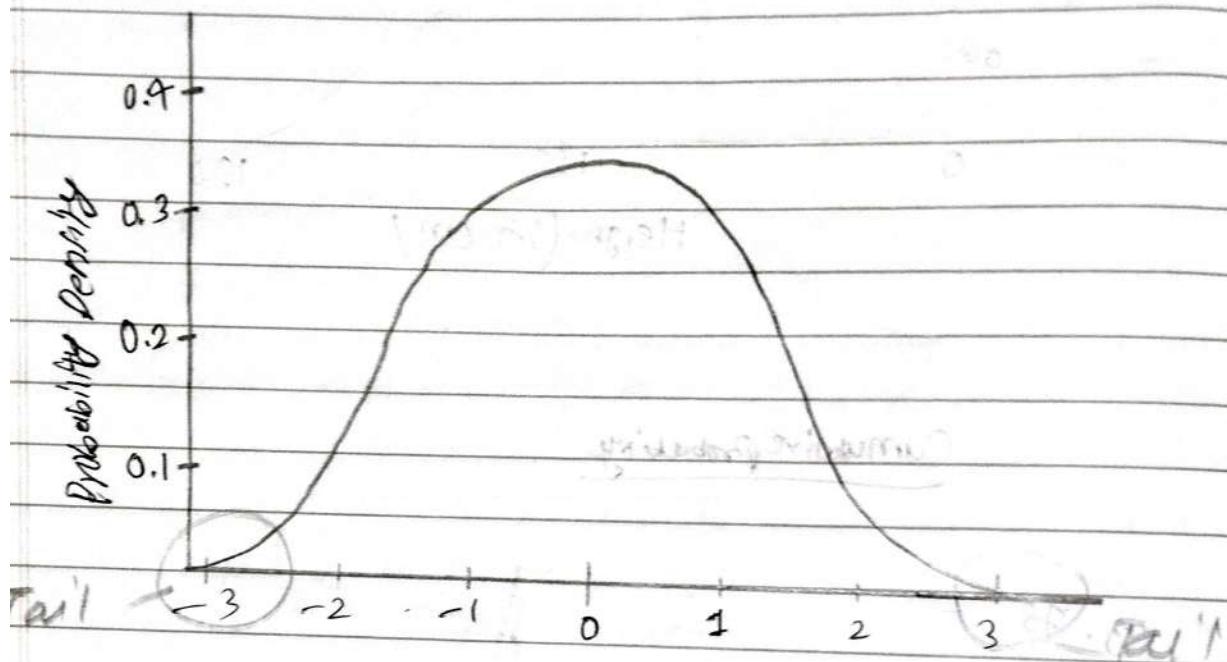
## Cumulative Probability



## Normal Distribution.

bell curve

Normal Distribution, also known as Gaussian distribution that is commonly used in statistical analysis. It is a continuous probability distribution that is symmetrical around the mean, with a bell-shaped curve.



Asymptotic in nature

Lots of points near the mean and very few far away

The Normal Distribution is characterized by two parameters: the mean ( $\mu$ ) and the standard deviation ( $\sigma$ ). The mean represents the centre of the distribution, while the standard deviation represents the spread of the distribution.

Denoted as:

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

Normal  
Dist

## Why it is so important?

→ Commonly in Nature : Many Natural Phenomena follow a normal Distribution, such as the height of people, the weights of objects, the IQ scores of population and many more. Thus, the normal distribution provides a convenient way to model and analyse such data.

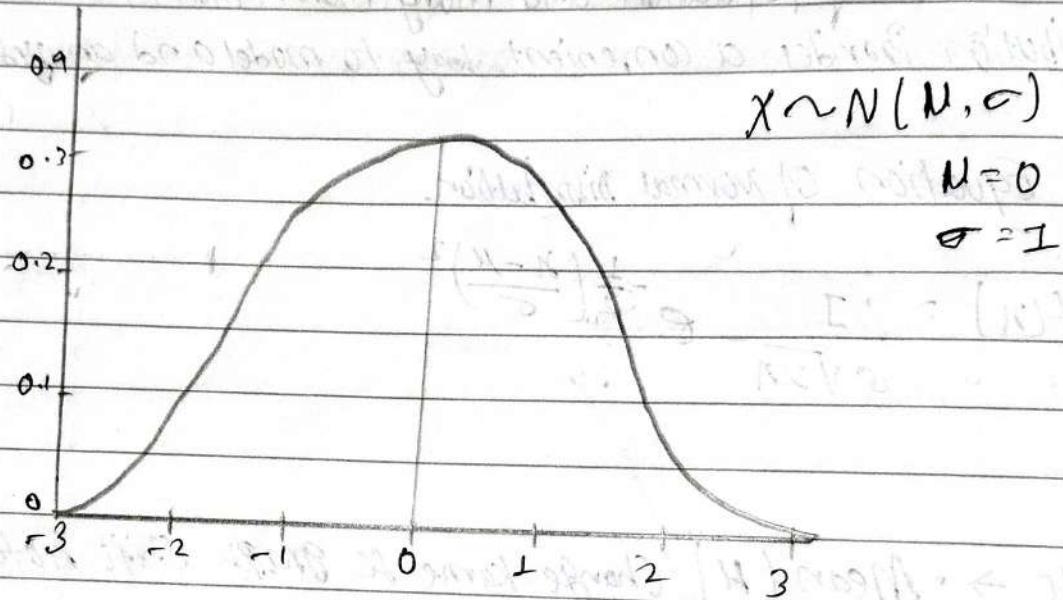
PDF Equation of Normal Distribution.

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

- Note → • Mean ( $\mu$ ) Change karne se graph shift hote hain  
 • Std ( $\sigma$ ) Change karne se graph ka spread change hota hain

## Standard Normal Variate (z)

→ A Standard Normal Variate (z) is a standardize form of the normal distribution with mean = 0 and std = 1.



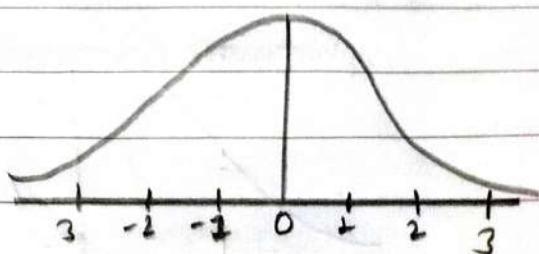
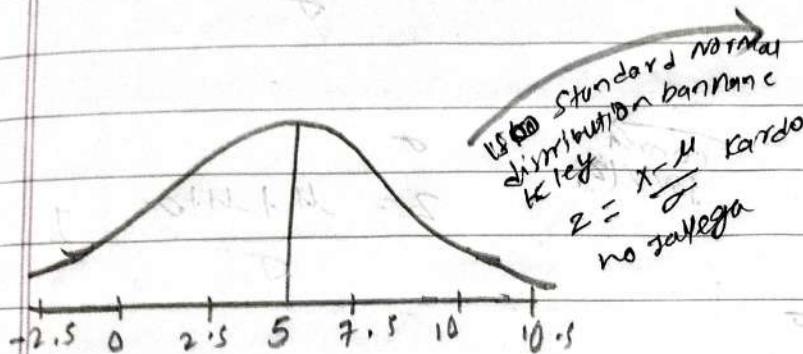
Standardizing a normal distribution allows you to compare different distribution with each other, and to calculate probability using standardized table or software.

Equation:

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

- How to transform a normal distribution to standardize Normal variable.

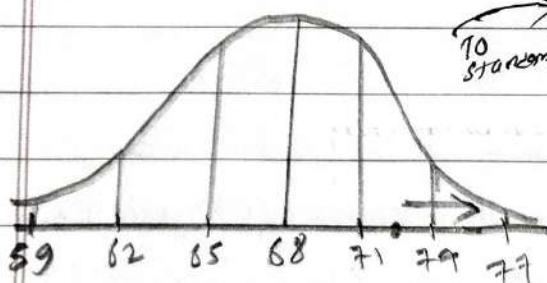
$$Z = \frac{X - \mu}{\sigma}$$



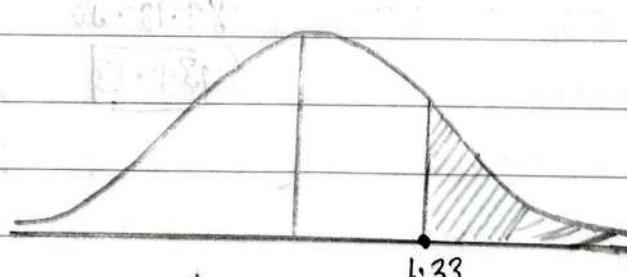
Kia fida hai karne ka?

Suppose the height of adult males in a certain population follow normal distribution with mean of 68 and a std of 3 inches. What is the Probability that a random selected male from this population is taller than 72 inches.

$$X \sim N(68, 3)$$



$$\rightarrow (Z) Z = \frac{72 - 68}{3} = \frac{4}{3} = 1.33$$



Z-table

→ A Z-table is used in Statistic to find the Probability that a standard Normal variable (Z) falls below a given value.

→ It is Standard Normal variable  
Ma  $1 - 0.90824$  se kam gane ka probability  $0.90824$  hai

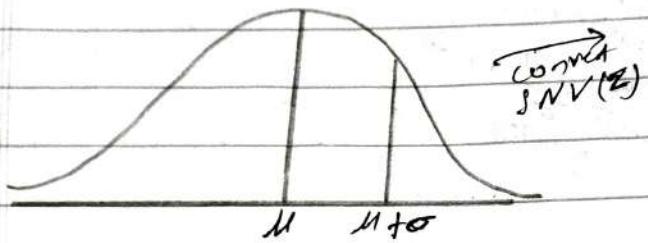
90.82 home chahiye 1.33 se jada

$$\begin{aligned} & 100 - 90.82 \\ & = 9.17\% \end{aligned}$$

for a Normal Distribution  $X \sim N(\mu, \sigma)$  What Percentage of population lie between mean and 1 standard deviation, 2 std and 3std?

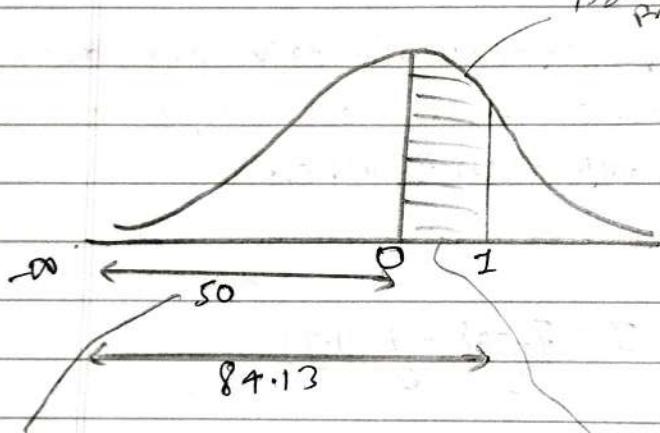
$$X \sim N(\mu, \sigma)$$

$$Z = \frac{M - \mu}{\sigma} = 0$$



$\text{Correct } SNV(Z)$

$$Z = \frac{\mu + \mu + \sigma}{\sigma} = 1$$



1 Standard deviation  
Probability

0.16  
0.04  
0.02  
0.01  
0.001  
0.0001

$$\begin{aligned} & 84.13 - 80 \\ & \boxed{34.13} \end{aligned}$$

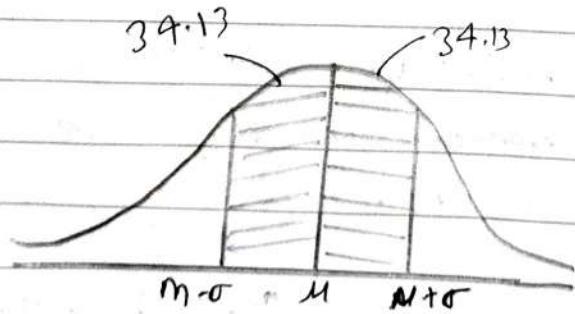
0.16  
0.04  
0.02  
0.01  
0.001  
0.0001

34.13

34.13

-1σ      +1σ

64.2

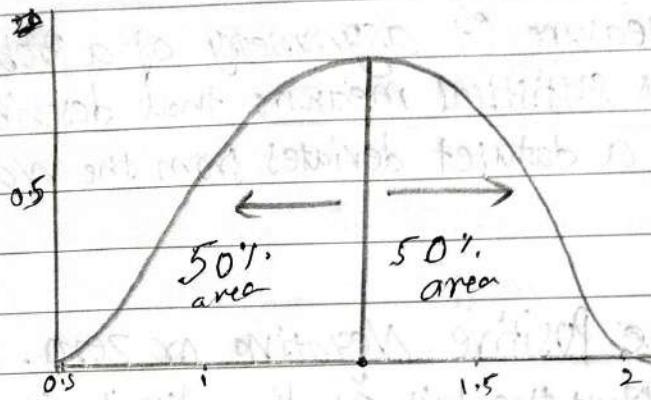


-1σ set +1σ tak 64.2%  
Data dalam  
empat

## Properties of Normal Distribution

### ① Symmetry

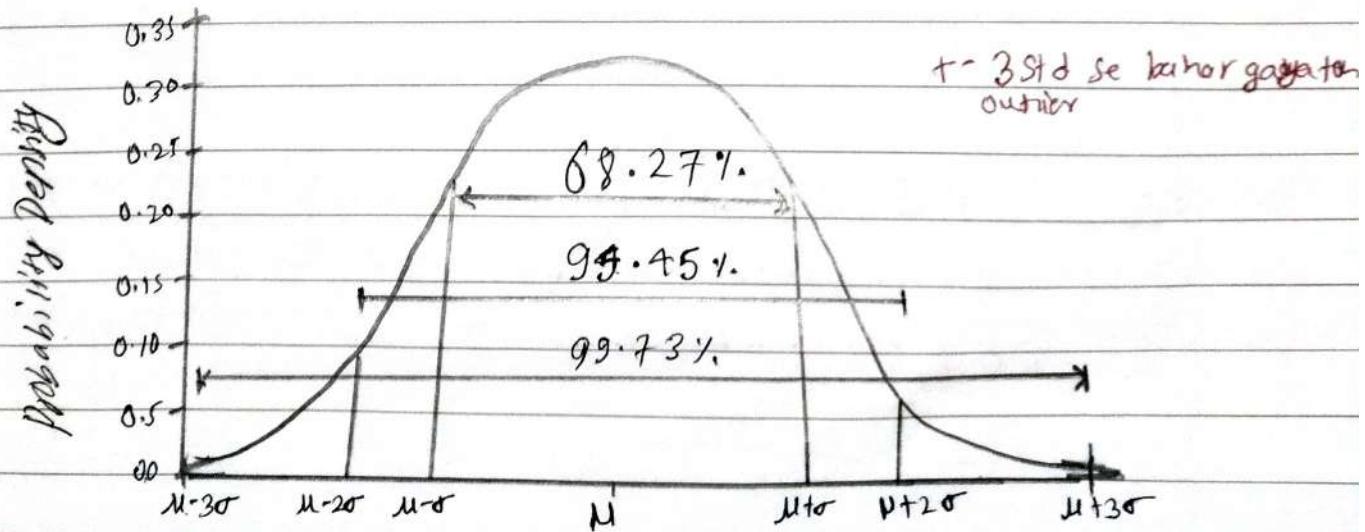
→ The normal distribution is symmetric about its mean which means that the probability of observing a value above the mean is the same as the probability of observing a value below the mean. The bell-shaped curve of the normal distribution reflects this symmetry.



- ② Measures of central tendencies are equal → mean = mode = median.  
 ④ The area under graph is 1.

### ③ Empirical rule

→ The normal distribution has a well-known empirical rule, also called the 68-95-99.7 rule, which states that approximately 68% of the data falls within one std of the mean, 95% of the data falls within 2 std of the mean, and about 99.7% of the data falls under the 3 std of the mean.

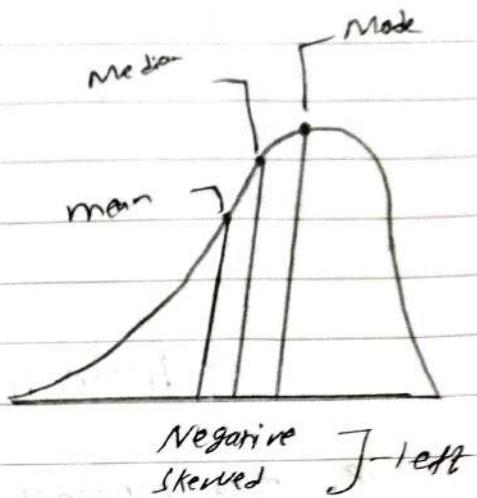
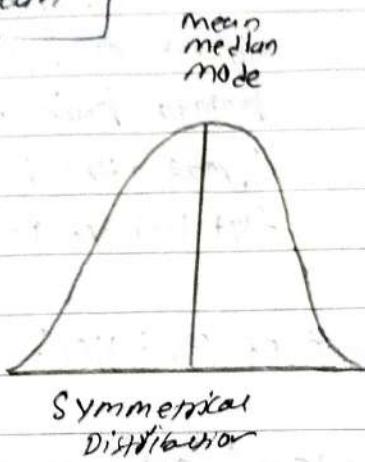
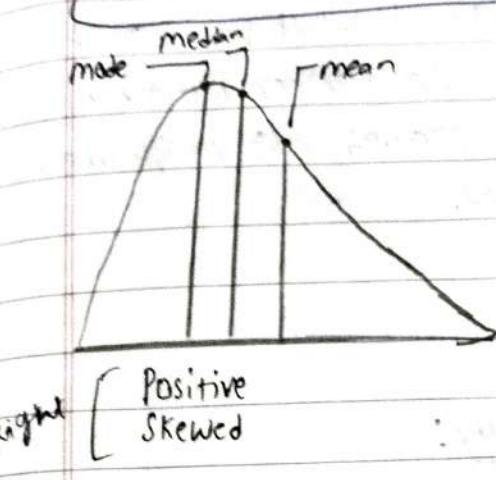


- What is Skewness?

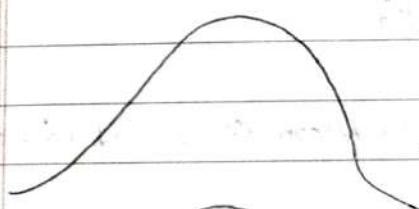
symmetric

- A Normal Distribution is a bell-shaped distribution with specific mathematical formula that describe how data is spread out. Skewness indicates that the data is NOT symmetrical, which means it is not normally distributed.
- Skewness is a measure of asymmetry of a probability distribution. It is statistical measure that describe the degree to which a dataset deviates from the Normal Distribution.
- Skewness can be positive, negative or zero. A positive skewness means that the tail of the distribution is longer on the right side, while negative skewness means that the tail is longer on left side. A zero skewness indicates a perfect symmetrical distribution.
- In a symmetric distribution mean, median, mode are equal. In contrast, in a skewed distribution, the mean, median, mode are not equal, and the distribution tends to have a longer tail on one side than the other.

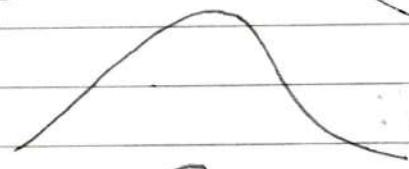
Mode < Median < mean



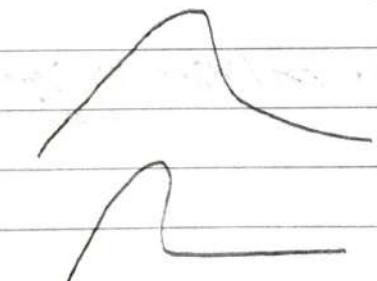
- The greater the skew the greater distance between mean, median and mode.



skew = 0



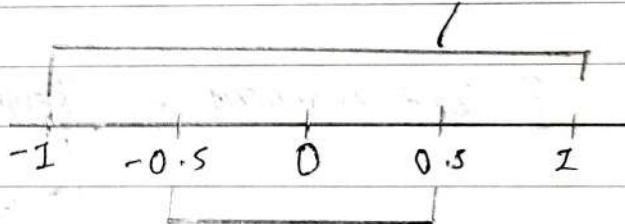
skew = 0.5



skew = 1



skew = 1.5



the Jada hua  
from skew

Use of Normal distribution

↳ Outliers detection

↳ Hypothesis testing

↳ Central limit theorem

↳ Assumption on data for ml algorithms

Note → agar mean -  $3 \times \text{std}$   
or  
mean +  $3 \times \text{std}$

Ke ya toh outlier find kar  
Sakre wo iss se

Statistical moments  $\rightarrow$  Statistical moments basically wo quantities hain jo distribution ka shape batata hain - Jaise center kaha hai, spread kitna hai, skewed ya symmetric, flat hai ya peaked.

#### ↑ Main moments of statistics:

① 1st moment - Mean (central tendency):

Yeh batata hai 'data ka centre kaha hai'

② 2nd moment - Variance (Dispersion):

Yeh batata hai 'data mean ke around kitna spread hai'

③ 3rd moment - Skewness (symmetry):

Yeh batata hai 'data left/right ya symmetric hai'

④ 4th moment - Kurtosis (Peakedness):

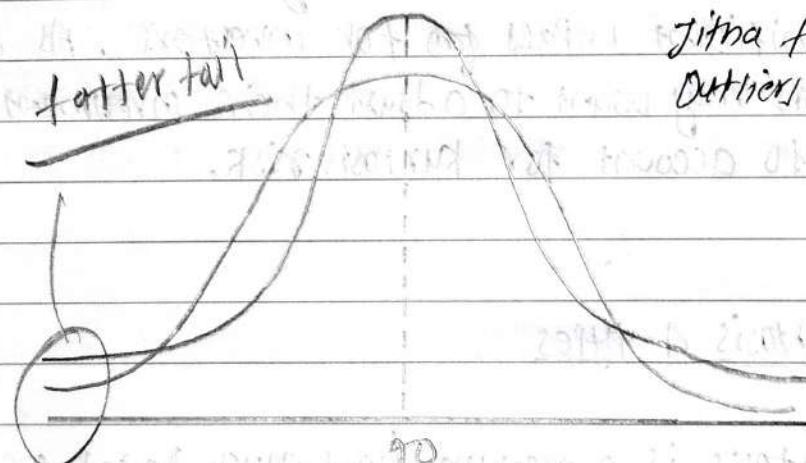
## What is kurtosis?

↳ Kurtosis is the 4th statistical moment. In probability and statistics, kurtosis (meaning "curved, arching") is measure of the "tailedness" of the probability distribution of a real-valued random variable. Like skewness, kurtosis describe a particular aspect of a probability distribution.

→ Data ke tails mai kitne heavy ya light (ie, extreme values kitne zada ya tam aate hain) yeh描述ता है।

fatness = presence of outliers

Jitna fatter tail utna satara  
Outliers hoga, ka chheses



formula simple kurtosis

$$\frac{n \times (n+1)}{(n-1) \times (n-2) \times (n-3)} \times \sum_{i=1}^n \left( \frac{x_i - \bar{x}}{s} \right)^4 - \frac{3 \times (n-1)^2}{(n-2) \times (n-3)}$$

- Practical Use - Case

In finance, kurtosis risk refers to the risk associated with the possibility of extreme outcome or "fat tails" in the distribution of returns of a particular asset or portfolio.

If a distribution has high kurtosis, it means that there is a higher likelihood of extreme events occurring, either positive or negative, compared to a normal distribution.

In finance, kurtosis risk is important to consider because it indicates that there is greater probability of large losses or gains occurring, which can have significant impact ~~for~~ for investors. As a result, the investors may want to adjust their investment strategies to account for kurtosis risk.

- Excess kurtosis & types

↳ Excess kurtosis is a measure how much peaked or flat distribution is compared to normal distribution, which is considered to have a kurtosis of 0. It is calculated by subtracting 3 from the sample kurtosis coefficient.

In my word

Kurtosis batata hou ke tails kitna fat hal (kitne extreme values hain) dikhana  
Lekin "excess kurtosis" sirf ek adjust kya version hai kurtosis ka -  
$$\text{Excess Kurtosis} = \text{Kurtosis} - 3$$

- 3 minus karne hain aur ke normal distribution ka kurtosis = 3 hota hai
- Toh agar usme se 3 hata do, toh normal distribution ka excess kurtosis 20 ban jata hai  
- easy comparison k lage,

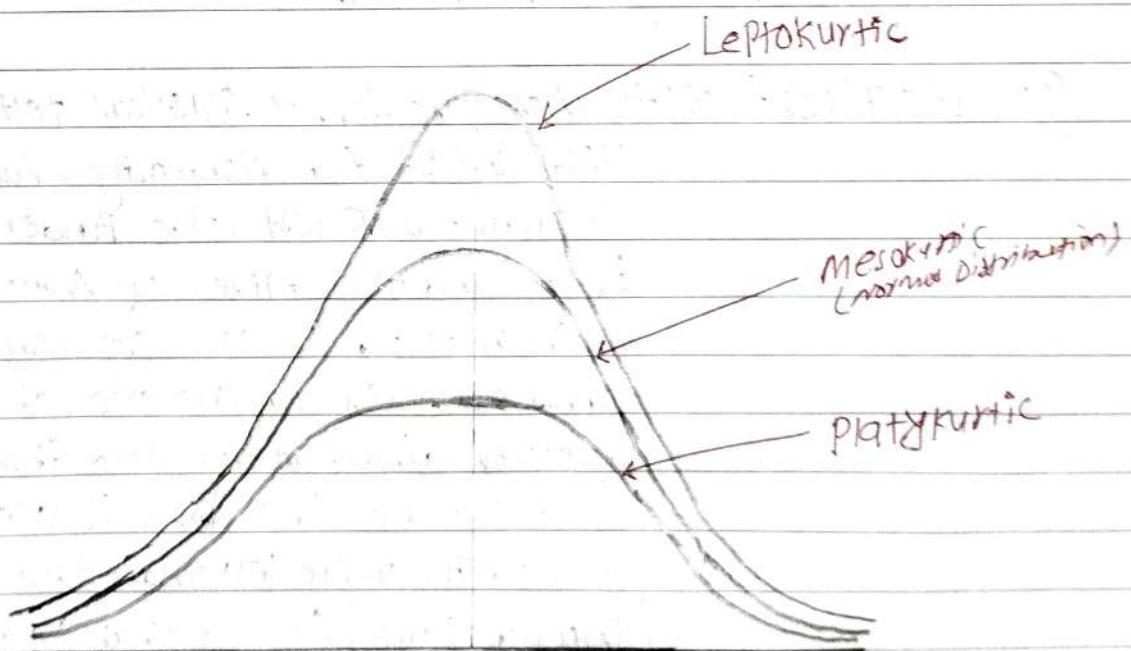
## Types of Kurtosis

Types	Excess Kurtosis	Tail Shape	Meaning
Mesokurtic	0	Normal tails	Normal distribution
Leptokurtic	>0	fat tails, sharp peak	Outlier zada
Platykurtic	<0	flat tails, broad shape	outliers kam

→ Excess Kurtosis tell how heavier or lighter the tails are compared to normal distribution.

Agar +ve hai → tails heavy → zyada extreme values.

Agar -ve hai → tails light → kam extreme values.



- How to find if a given distribution is normal or not?

(1) Visual Inspection → One of the easiest ways to check for Normality is to visually inspect a histogram or a density plot of the data. A Normal distribution has bell-shaped curve, which means that the majority of the data falls in the middle, and the tails taper off symmetrically. If the distribution looks approximately bell-shaped, it is likely to be normal.

(2) P-P Plot → Another way to check for Normality is to create a Normal Probability Plot (also known as Q-Q Plot) of the data. A Normal Probability Plot plots the observed data against the expected value of a Normal distribution. If the data points fall along a straight line, the distribution is likely to be normal.

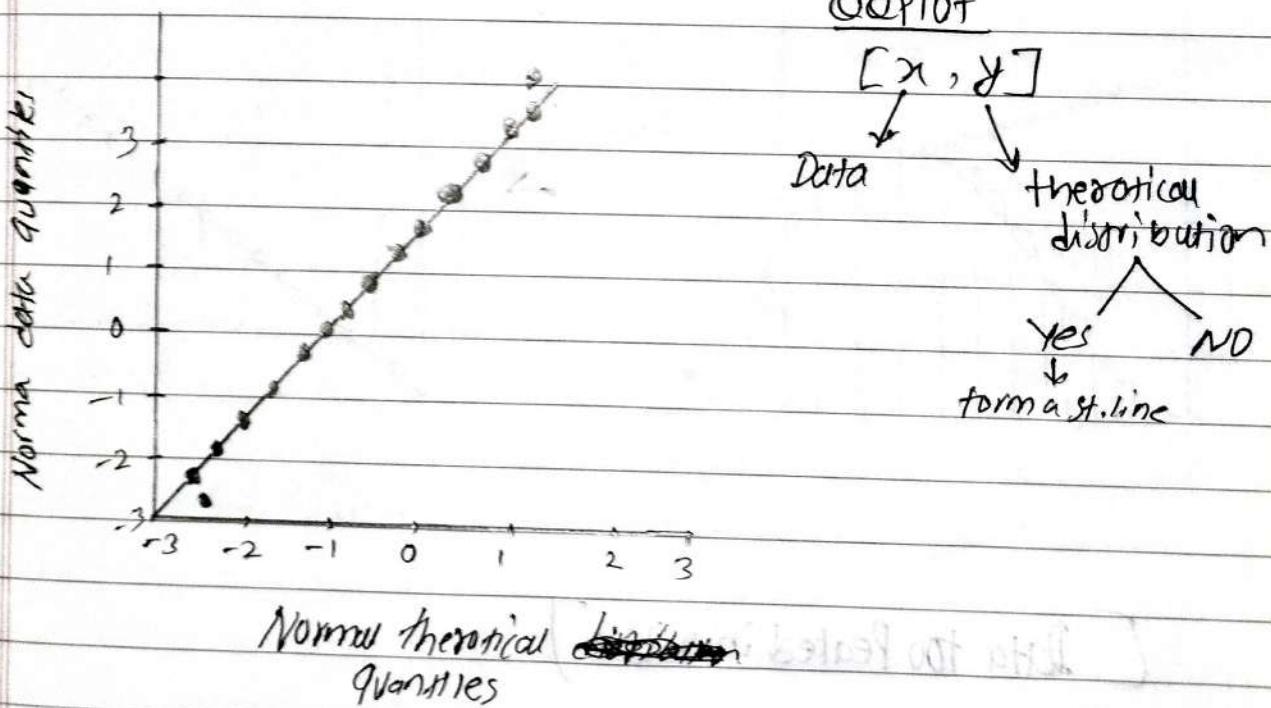
(3) Statistical tests → There are several statistical tests that can be used to test for Normality, such as the Shapiro-Wilk test, the Anderson-Darling test, and the Kolmogorov-Smirnov test. These tests compare the observed data to the expected values of a Normal distribution and provide a p-value that indicates the data is likely to be normal or not. A p-value less than the significance level (usually 0.05) suggests that the normal data is normal.

- What is a QQ plot and how it is plotted?

↳ A Q-Q Plot (Quantile - Quantileplot) is a graphical tool to check whether your data follows a specific theoretical distribution - usually the normal.

### Concept

- You plot your data quantiles (sorted values) against the quantiles of theoretical distribution (e.g., normal).
- If your data actually follows that distribution, the points will fall roughly on straight line (nearly).

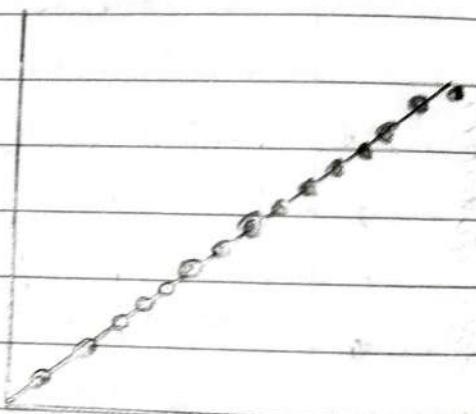
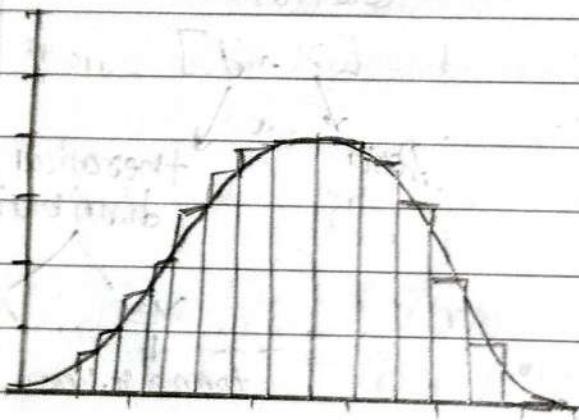


- QQ plot batata hai ke tumhare data kis distribution k jaise dikhta hai - agar seedhi line hai, sab set hai; Agar tedha hai, kuch gadbad hai.

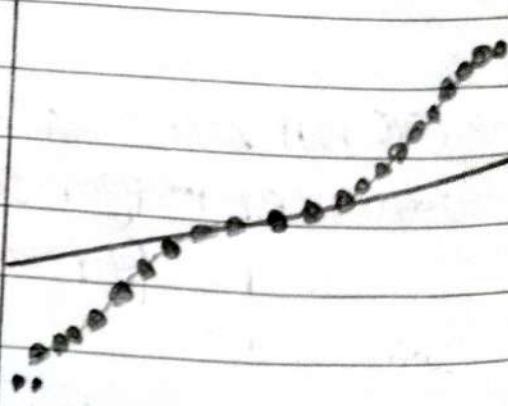
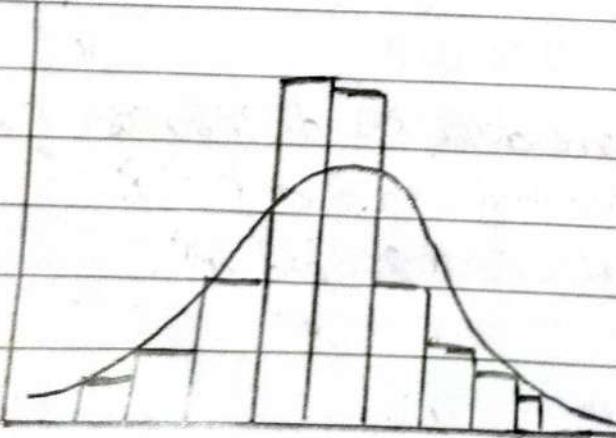
→ In QQ Plot, the quantiles of two set of data are plotted against each other. The quantiles of one set of data are plotted on x-axis, while the quantiles of the other set of data plotted on the y-axis. If two set have the same distribution, the QQ Plot will fall on a straight or straight line. If two sets of data do not have the same distribution, the points will deviate from the straight line.

- How to interpret QQ Plot.

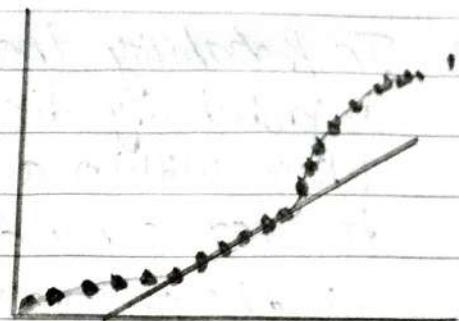
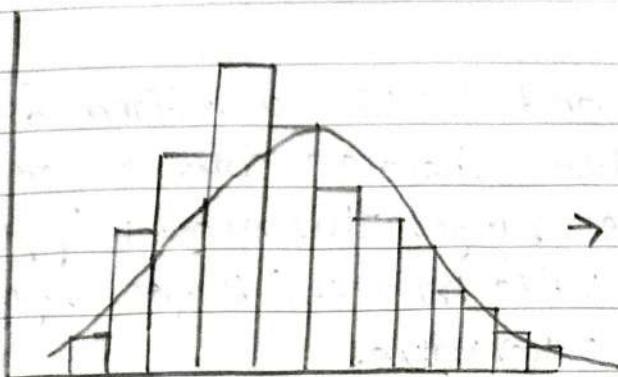
[Normal Distribution Data]



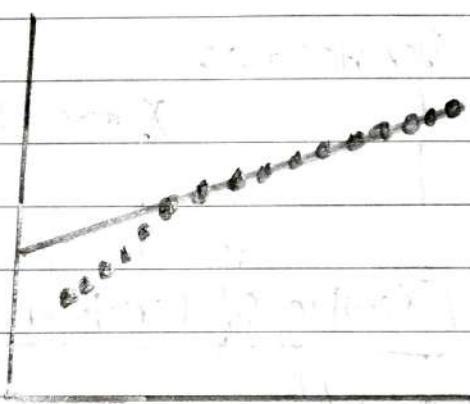
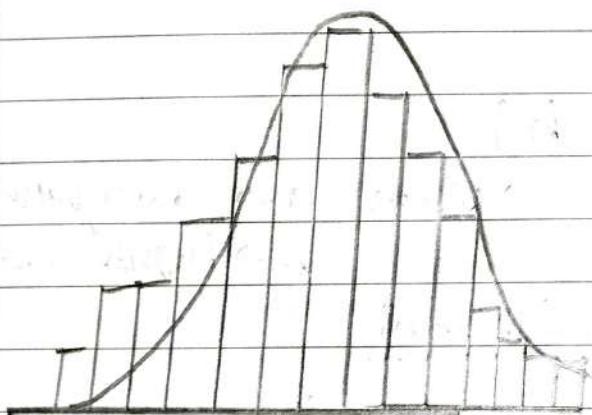
[Data too peaked in middle]



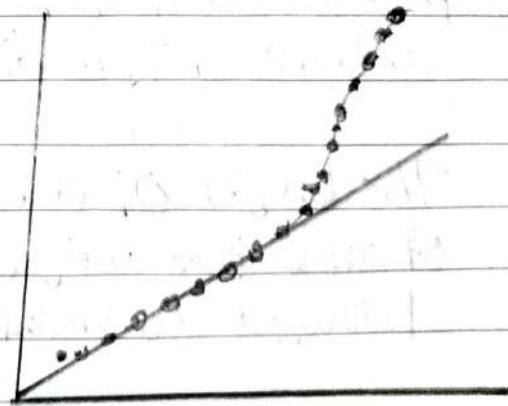
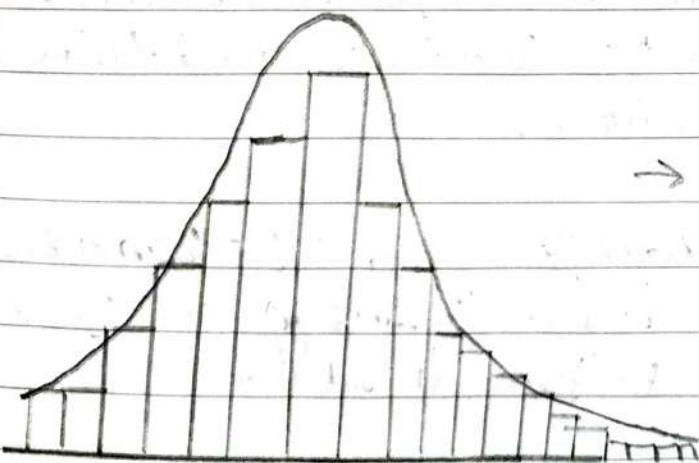
[Skewed Data]



[Skewed left]

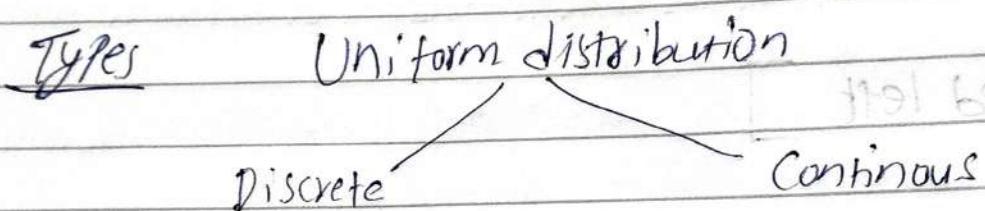


[Skewed right]



## What is Uniform Distribution

↳ In Probability theory and statistics, a uniform distribution is probability distribution where all outcomes are equally likely within a given range. This means if you were to select a random value from this range, any value would be as likely as other value.



Denotes as:

$$X \sim U(a, b)$$

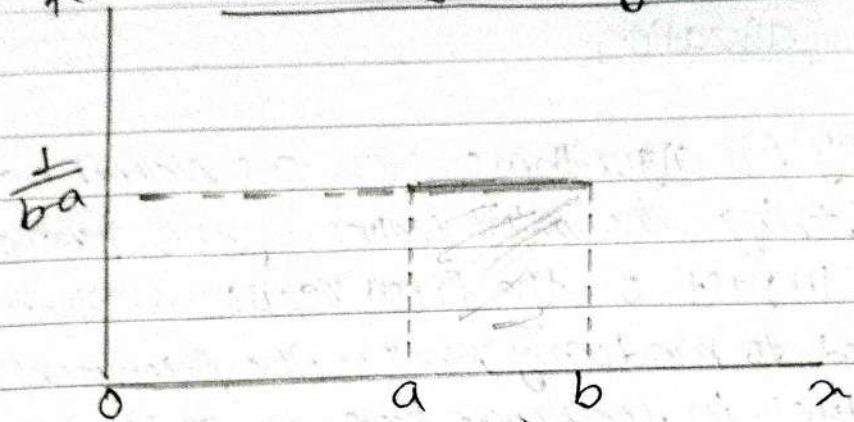
where  $a \rightarrow$  lower value  
 $b \rightarrow$  higher value

Example of Continuous Uniform distribution:

- (1) The height of a person randomly selected from a group of individuals whose heights range from 5'8" to 6'0" would follow a continuous uniform distribution.
- (2) The time it takes for a machine to produce a product, where the production time range from 5 to 10 minutes, would follow a continuous uniform distribution.
- (3) The weight of a randomly selected apple from a basket of apple that weights between 100 and 200 gm, would follow a continuous uniform distribution.

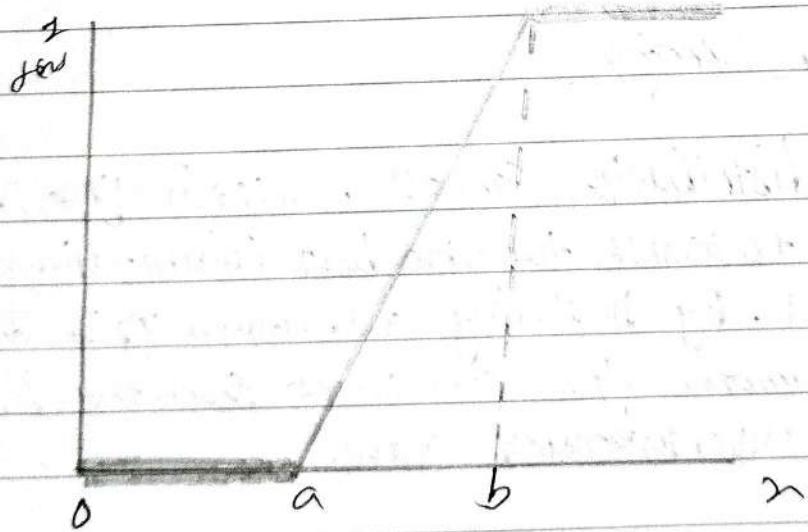
## PDF CDF and graphs.

for Probability Density function



Uniform equation  $f(x) = \begin{cases} \frac{1}{b-a} & \text{for } a \leq x \leq b \\ 0 & \text{for otherwise} \end{cases}$

Cumulative distribution function



Skewness  $\rightarrow 0 \rightarrow$  symmetric  $\rightarrow$  normal.

- Uniform distribution
- Application in machine learning and Data science.

### ① Random initialization

↳ In many ML algorithms, such as neural networks and k-mean clustering, the initial values of the parameters can have a significant impact on the final result. Uniform distribution is often used to randomly initialize the parameters, as it ensures that all values in the range have an equal probability of being selected.

### ② Sampling

↳ Uniform distribution can be used for sampling. For example, if you have a dataset with an equal no of samples from each class, you can use uniform distribution to randomly select a subset of the data that is representative of all the classes.

### ③ Hyperparameter tuning

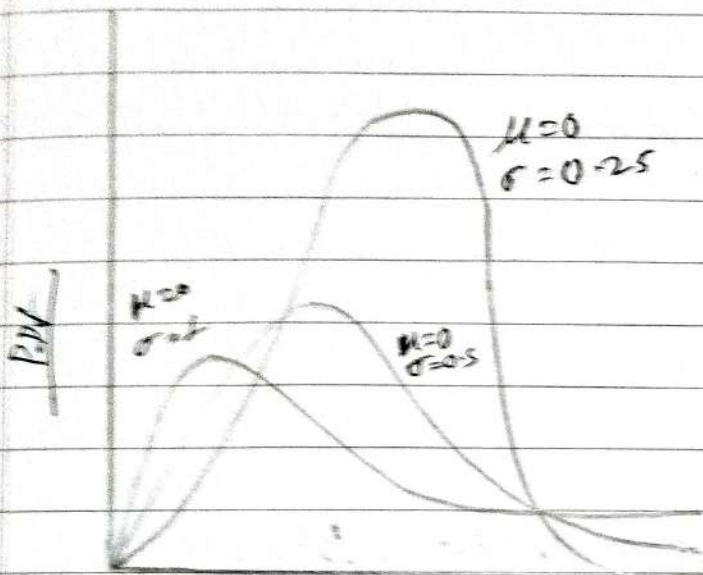
↳ Uniform distribution can also be used in Hyperparameter tuning, where you need to search for the best combination of hyperparameters for a ML model. By defining a uniform prior distribution for each hyperparameter, you can sample from the distribution to explore the hyperparameter space.

### ④ Data augmentation

↳ In some cases, you may want to artificially increase the size of your dataset by adding new examples that are similar to original data. Uniform distribution can be used to generate new data points that are within a specified range of the original data.

## Log Normal Distribution

- In probability ~~Probability~~ theory and statistics, a lognormal distribution is heavily tailed continuous probability distribution of a random variable whose logarithm is normally distributed.



### Examples

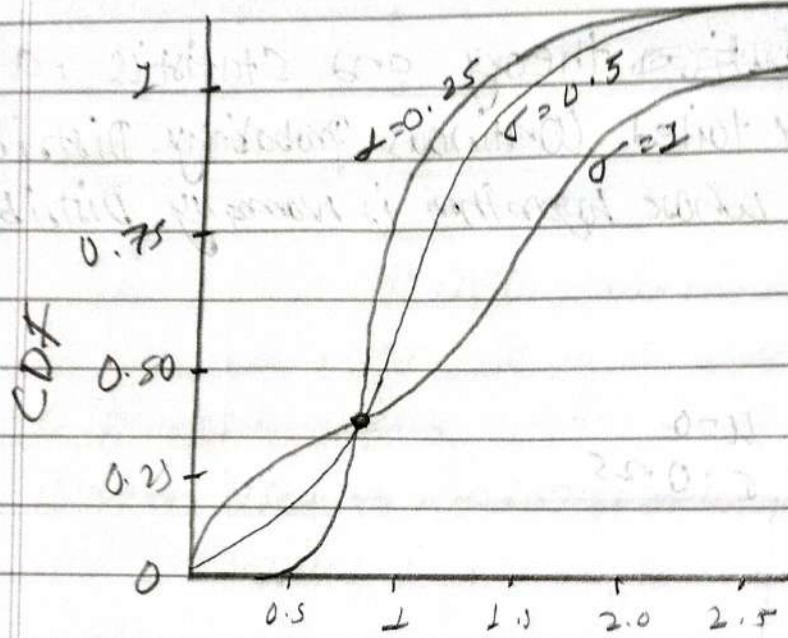
- The length of comment posted in internet discussion forums follow a log-normal distribution.
- User dwell time on online article follow a log-normal distribution.
- The length of chess game tends to follow a log-normal distribution.
- In economic, there is evidence that the income of 97% - 98% of the population is distributed log-normally.

Denoted by:  $X \sim \text{lognormal}(\mu, \sigma)$

### PDF equation

$$f(x) = \frac{1}{x\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{\ln x - \mu}{\sigma}\right)^2}, \text{ for } 0 < x < \infty$$

CDF

Skewness  $\rightarrow$  Skewed

- How to check a random variable is log normally distributed?

If  $X$  is log normally distributed, then  $Y = \log(X)$  should be normally distributed. To check normal distribution of  $Y$  plot a Q-Q plot if straight bands then conform log normally distributed law.

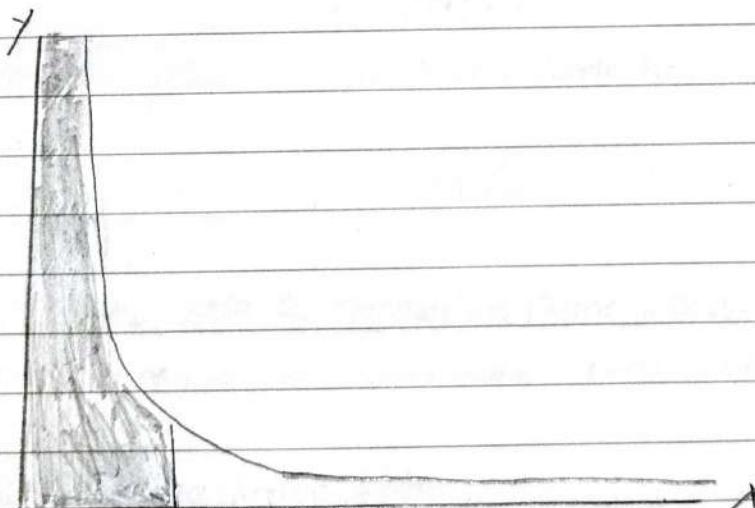
## Pareto Distribution

- The Pareto Distribution is a type of Probability Distribution that is commonly used to model the distribution of wealth, income and other quantities that exhibit a similar Power-law behaviour

### What is Powerlaw

- In Mathematics, a power law is a functional relationship between two variables, where one variable is proportional to power of the other. Specifically, if  $y$  and  $x$  are two variables related by a power law, then the relationship can be written as:

$$y = k^x x^\alpha$$



Vilfredo Pareto originally used this distribution to describe the allocation of wealth among individuals since it seemed to show well the way that larger portion of the wealth of society is owned by smaller % of the people in the society. He also used it to describe distribution of income. The idea is something expressed more simply as the Pareto principle or the "80-20" rule which says that 20% of the population controls 80% of the wealth.

## Mathematical transformation

- Mathematical transformation like log, sqrt, reciprocal, power, Box-Cox, Yeo-Johnson can make a distribution more normal, if it is skewed.

But yeh guarantee nahi hai ke har transformation se distribution Perfect Normal ho jaye. Kuch distribution extreme hain (bimodal, heavy tail), unko normal banana impossible hai with just a simple transform.

### Log transform

- ↳ Log transform ek mathematical transformation hai jisme ham kisi variable  $x$  ka logarithm leta hai:
- $$y = \log(x)$$
- Sirf positive values  $x$  liya kaam karta hai

left skewed  
reya  
square  
vir  
karo

Use kaha hota hai / kya karte hai

- ② Right-skewed data ko normalize karne kya:

↳ ex income, population, sales data - bahut bade values ko compress karta hai

- ③ Variance Stabilize karne kya

↳ jab high values spread bohat tha hota hai, log transform se variance relatively stable ho jate hai

- ④ Outliers ko reduce karne kya

↳ Extreme value ko compress karke analysis aur modeling main influence karte hai.

## Discrete Distribution

### Bernoulli Distribution

↳ Bernoulli Distribution ek discrete Probability distribution hai jisme sirf do possible outcome hote hain:  
Success (1)  
failure (0)

Example: Cointoss  $\rightarrow$  Head (1), tail (0)

Product test  $\rightarrow$  Pass (1), fail (0)

SPAM detection  $\rightarrow$  SPAM (1), Not SPAM (0)

P  $\rightarrow$  Probability of success (1)

P-1  $\rightarrow$  Probability of failure (0)

Only 1 Parameter P

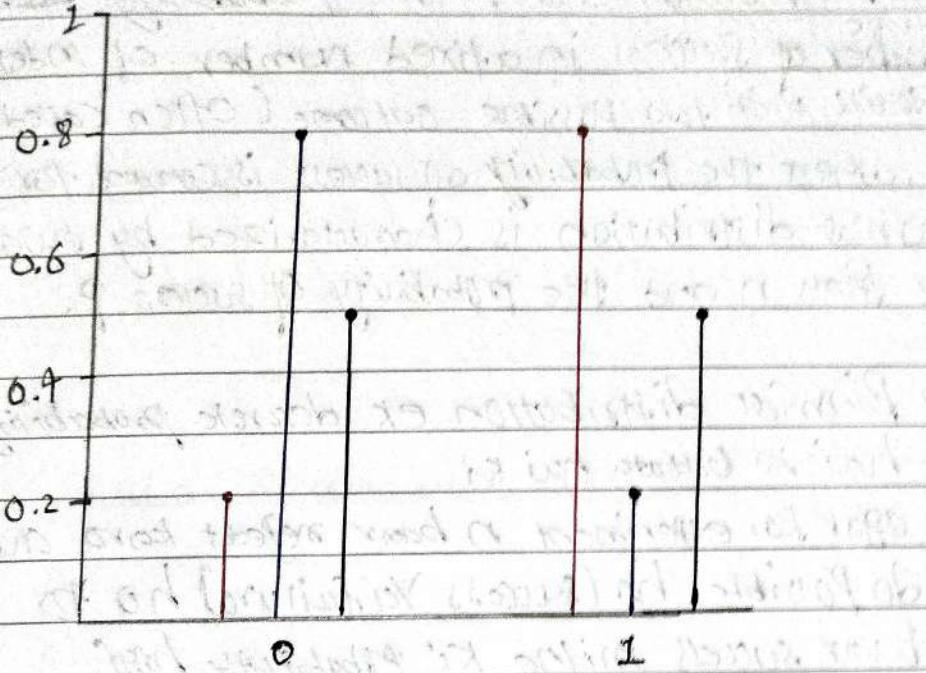
PMF :

$$P(X=x) = p^x(1-p)^{1-x}, \quad x \in \{0, 1\}$$

$P(X=x)$  batata hai ke success (1) ya failure (0) hone ki kitne chance hain.

The Bernoulli distribution is commonly used in ml for modeling binary outcomes, such as whether a customer will make a purchase or not, whether a patient will have a certain disease or not.

## PMF of Bernoulli distribution



Three example of bernoulli distribution.

■  $P(x=0) = 0.2$  and  $P(x=1) = 0.8$

■  $P(x=0) = 0.8$  and  $P(x=1) = 0.2$

■  $P(x=0) = 0.5$  and  $P(x=1) = 0.5$

## Binomial Distribution

↳ Binomial Distribution is a Discrete Probability distribution that Number of success in a fixed number of independent Bernoulli trials with two possible outcome (often called "success" or failure), where the probability of success is constant for each trial. The binomial distribution is characterized by two parameters the no of trials n and the probability of success p.

Note → Binomial distribution ek discrete probability distribution hai jo baatka hai ki agar koi experiment n baar repeat karo aur 1 baar do possible ho (success ya failure) ho to kitni baar success milne ki probability hogi;

The Probability of anyone watching this lecture in future and liking it is 0.5. What is the probability that:

① NO-one Out of 3 People will like it.

$$\begin{array}{l} \text{L} \\ \frac{1}{8} \end{array}$$

YYY

YYN

YNY

YNW

NYY

NYN

NNY

NNW

② 1 Out of 3 People will like it.

$$\begin{array}{l} \text{L} \\ \frac{3}{8} \end{array}$$

③ 2 Out of 3 People will like it.

$$\begin{array}{l} \text{L} \\ 3/8 \end{array}$$

Topic 2

Binomial distribution ^ Parameter :

 $n \rightarrow$  no of trials $P \rightarrow$  probability of success.

Prob of Binomial distribution :

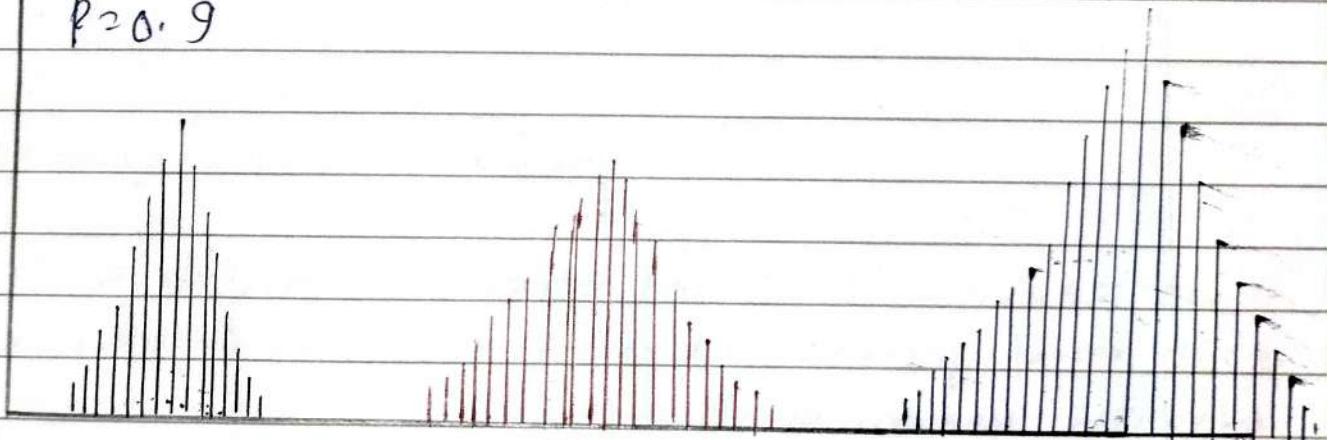
$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

Binomial distribution with Different Probability of success

$$P=0.2$$

$$P=0.6$$

$$P=0.9$$



### Note

Jab  $P=0.5$  distribution Perfect symmetric hoti hai - bechma centred.

Jab  $P \neq 0.5$ , distribution right skewed (tail right side pe hoti hai)

Jab  $P > 0.5$ , distribution left skewed (tail left side on hoti hai)

Certain:

- ① The process consist of n trials.
- ② Only 2 exclusive outcome are possible, a success and a failure.
- ③  $P(\text{success}) = p$  and  $P(\text{failure}) = 1-p$  and it is fixed from trial to trial.
- ④ The trials are independent.

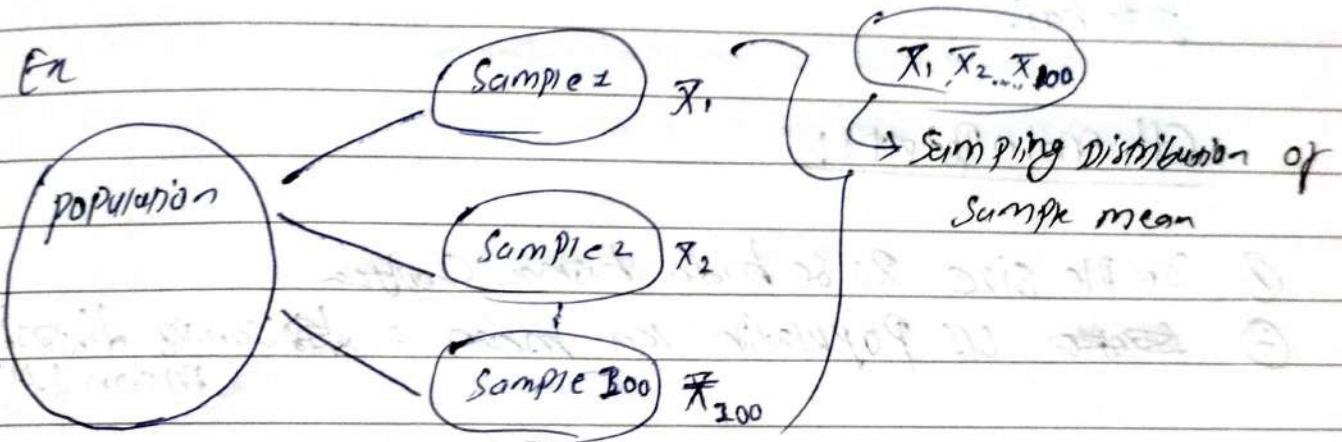
$$P(X=k) = \binom{n}{k} p^k (1-p)^{n-k}$$

Where Binomial distribution used in:

- ① Binary Classification Problems
- ② Hypothesis testing
- ③ Logistic regression
- ④ A/B testing.

## Sampling Distribution

Sampling distribution ka matlab hai, Jab Population se equal size ~~of sample~~ & multiple sample lena & baad, un sab samples ka mean (ya koi aur statistic) ka set jg banta hai, wahi hota hai sampling distribution.



### Why Sampling Distribution is important?

Sampling distribution important hai Kyunki Yeh Batata hai ke Sample statistic Population k respect ma kime vary karega, Jiss se hum confidence intervals, hypothesis tests, aur prediction kar sakte hain.

mean median  
mode

## Central Limit theorem

→ Agar hum kisi Bhi population (chahiye wo normal ho ya  
na ho) se bahat sare random samples take lein. ~~consequently~~  
aur hum har sample ka mean nikale de, ~~and~~ toh  
Un sub sample means ka distribution (graph) normal ho  
jata hai.

### Clt main points:

- ① Sample size 30 se bada huna chahiya
- ② ~~US~~ US Population ka mean = ~~sample~~ sample distribution mean ( $\mu$ )
- ③ Population std =  $\sigma$   
Sample size =  $n$

$$\text{sample distribution} = \frac{\sigma}{\sqrt{n}}$$

$$\text{Sample mean} = \text{Population mean} = \mu$$

$$\text{Sample std} = \frac{\sigma}{\sqrt{n}}$$

→ The Clt is important in statistic and ml because it allows us to make probabilistic inferences about population based on sample of data. for ex: we can use Clt to construct confidence intervals, perform hypothesis tests, and make predictions about the mean based on the sample.

## SOME TERMS

**Population** → A Population is the entire group or set of individuals, objects, or events that a researcher want to study or draw conclusion about it.

**Sample** → A Sample is a subset of the population that is selected for study.

**Parameter** → Parameter ek number hota hai jo population ki property batata hai like, mean, std, var etc.

**Statistic** → ~~sample~~ statistic ek number hota hai jo sample ki property batata hai like  $\bar{x}$ ,  $s$ ,  $var$  etc.

Population → parameter → Mostly unknown

sample → statistic → calculate kar sakte ho

**Point estimate** → ek single Numerical value jo population parameter ka best guess batata hai, sample ke basis pe.

In other words:

Sample ke statistic ek Numerical value hota hai : Jab ham us value ko use population ka parameter ka aandara lagane ke liye ~~keya~~ keya jata hoi, tb use point estimate kaheta hoi.

## Confidence interval (CI)

Confidence interval ek range of values hai jisme Population Parameter (like mean, proportion) hone ki Probability hyn hote hai - based on sample data

- Sample ke data se Confidence interval nikalte hain, aur uss range se Population ka aasli value guess karte hain.

→ Confidence interval is created for Parameters and NOT for statistics. Statistics helps us get the Confidence interval for a parameter.

## Confidence level

It is usually expressed as a percentage like 95%, indicates how sure we are that the true values lies within the interval.

Confidence interval = Point estimate  $\pm$  Margin of error

Ways to calculate confidence interval:

Z producer

t - producer

### ① Zproducer (sigma known)

#### Assumption

- ① Random sampling  $\rightarrow$  The data must be randomly selected to avoid bias
- ② Known population standard deviation
- ③ Normal Distribution or large sample ( $n \geq 30$ )  $\rightarrow$  The population should be normal, or the sample should be large enough for the central limit theorem to apply.

formula of CI using Z producer:

$$\boxed{CI = \bar{x} \pm z_{\alpha/2} \times \frac{\sigma}{\sqrt{n}}}$$

$\bar{x}$  = Sample mean

$\sigma$  = Population standard deviation (known)

$n$  = sample size

$z_{\alpha/2}$   $\rightarrow$  critical value (from z-table)

90% CL  $\rightarrow 1.645$

95% CL  $\rightarrow 1.96$

99% CL  $\rightarrow 2.575$

Z can calculate by  
Area to left =  $\frac{(CI + I)}{2}$

Put this value  
in place of  
 $z_{\alpha/2}$

- (1) The average test scores in chemistry class is normally distributed with std of 6.5. 100 scores with sample mean of 82 were selected at random.
- find a 90% confidence interval for the population mean exam score.
  - find the value of margin of error

Soln,

$$\sigma = 6.5, n = 100, \bar{x} = 82$$

$$\text{Area to left} = \frac{CL + I}{2} = \frac{0.90 + 1}{2} = \boxed{0.95}$$

$\frac{Z_{\alpha/2}}{2} \rightarrow 1.645$

$$CI = \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}$$

$$= 82 \pm 1.645 \left( \frac{6.5}{\sqrt{100}} \right)$$

$$= 82 \pm 1.06925$$

$$\boxed{(80.93075, 83.06925)}$$

$$\text{Margin of error} = Z_{\alpha/2} \times \frac{\sigma}{\sqrt{n}}$$

$$= 1.645 \left( \frac{6.5}{\sqrt{100}} \right)$$

$$= \boxed{1.06925}$$

## Interpreting Confidence Interval?

A confidence interval is a range of values within a population parameter such as the population mean, is estimated to lie within a certain level of confidence. The confidence interval provides an indication of the precision and uncertainty associated with the estimate. To interpret the confidence interval values, consider the following points:

### ① Confidence level

- ↳ The confidence level (commonly set at 90%, 95% or 99%) represents the probability that the confidence interval will contain the true population parameter if the sampling and estimation process were repeated multiple times. For ex., a 95% confidence interval means that if you were to draw 100 different samples from the population and calculate the confidence interval for each, approximately 95 of those intervals would contain the true population parameter.

### ② Interval range

- ↳ The width of the confidence interval gives an indication of the precision of the estimate. A narrower confidence interval suggests a more precise estimate of the population parameter, while wider intervals indicate greater uncertainty. The width of the interval depends on the sample size, variability in the data, as the desired level of confidence.

$$95\% \uparrow \rightarrow \text{Width} \uparrow$$

factor that effect confidence interval.

Narrow - chotta  
wide - bada

### Sample size

- ↳ longer sample size → narrow confidence interval
- ↳ more data give more precise estimate of the population parameter

ex

$$n=30 \rightarrow CI(48.5, 51.5)$$

$$n=300 \rightarrow CI(49.5, 50.5)$$

### Confidence level

- ↳ Higher confidence level → wider confidence interval
- ↳ To be more confident, you must allow a bigger range.

ex

90% CL → narrow

99% CL → wider

Note

100% CI =  $-\infty$  to  $+\infty$

### Population variability ( $\sigma$ )

- ↳ Higher variability → wider CI
- ↳ more variation in data = less precision in estimating the mean

$\sigma = 5 \rightarrow$  narrow CI

$\sigma = 20 \rightarrow$  wide CI

Confidence Interval (sigma not known)

## T Procedure (sigma not known)

Assumption:

① Random Sampling → The Data must be randomly selected to avoid Bias.

② Sample Standard deviation → The population std is unknown, and the sample std is used as an estimate.

③ Approximately normal Distribution → The t procedure assumes that the underlying distribution is approximately normally distributed, or sample size large enough to apply (CLT).

④ Independent Observations → The observation in the sample should be independent of each other.

formula of CI using t procedure

Sample std ka mean

$$CI = \bar{X} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$$

$\bar{X}$  = sample mean

$s$  = sample std - ka - mean

$\sqrt{n}$  → size of sample

$t_{\alpha/2}$  No of degree →  $n - 1$

$t_{\alpha/2}$  → CI (Confidence Level)

Kin ek CL K kya chal jaa

aur chart kro t table

ma uska print kya koi

aur kare nikalo.

Note → agar CI nikalne hai using t procedure to find sample mean nikalo fir sample std ka mean nikalo usko sample std k place mai logo use bad t table ma kina confidence level k leya chal jao unhe k leya value dekho aur df ma sample se kam ka value dekho.

- Hypothesis Testing kyu karte hain:

→ Ham data se decision lena chahte hain - lekin har baar jo difference ya pattern data mein dikhta hai, wo real bhi ho sakte hai ya Chained (randomness) se bhi ho saka hain. Toh Hypothesis Testing se hum check karte hain ke:

kya jo result mila hai actual meaningful hai ya sirf luck hai?

### Hypothesis Testing

→ Hypothesis Testing ek statistical method hai jis se hum yeh decide karte hain ke sample data ke basic par population ke bare main jo hum assume kar raha hai wo sanii hai ya galat.

#### Real life example:

Maan to tu ek medicine test kar dahe hai:

Ab tu dekhna chaheta hai ke nayi dawa sach maius benetar kaam kar rahi hai ya baki.

Tu data collect karnega aur Hypothesis Testing se decide karnega ki kya difference statistically significant hai ya sirf randomness chances se hua hai.

## ① Null Hypothesis ( $H_0$ ):

→ In simple terms, the Null Hypothesis is a statement that assume there is no significant effect or relationship between the variable being studied. It serves as the starting point for hypothesis testing and represents the status quo (or the assumption of no effect until proven otherwise). The purpose of hypothesis testing is to gather evidence (data) to either reject or fail the null hypothesis in favor of the alternative hypothesis, which claims there is a significant effect or relationship.

Null Hypothesis → kuch bhi change/relationship lekaa rhi ha  
matlab default assumption : sab kuch normal, koi diff nahi

## ② Alternate Hypothesis ( $H_1$ or $H_a$ ):

→ The alternate Hypothesis, is a statement that contradicts the Null hypothesis and claims there is a significant effect or relationship between the variable being studied. It represents the research hypothesis or the claim the researcher ~~wants~~ wants to support through statistical analysis.

Alternate Hypothesis → Alternate hypothesis will has difference, relationship, or effect bew population ma.

## Important Points

- How to decide what will be null hypothesis and what will be Alternate hypothesis ( Typically the null hypothesis says nothing new is happening).
- We try to gather evidence to reject the null hypothesis.
- It is important to note that failing to reject the null hypothesis doesn't necessarily mean that the null hypothesis is true; it just means that there isn't enough evidence to support the alternative hypothesis.
- Hypothesis test are similar to jury trials . in a sense .  
in a Jury trial ,  $H_0$  is similar to the NOT guilty verdict , and  $H_a$  is the guilty verdict . You assume in a jury that defendant isn't guilty unless the Prosecution can show beyond a reasonable doubt that he or she is guilty . If the jury says the evidence is beyond a reasonable doubt , they reject  $H_0$  , NOT guilty , in favour of  $H_a$  , guilty .

## Steps involved in Hypothesis Testing.

### Rejection Region Approach

1. Formulate a Null and Alternate Approach
  2. Select significant level → This is the probability of rejecting the null hypothesis when it is actually true, usually set (0.05 or 0.01)
- 
- ① Check assumption (example distribution)
  - ② Decide which test is appropriate (z-test, t-test, chi-square, ANOVA)
  - ③ State the relevant test statistic
  - ④ Conduct the test
  - ⑤ Reject or not reject the Null Hypothesis
  - ⑥ Interpretate the result.

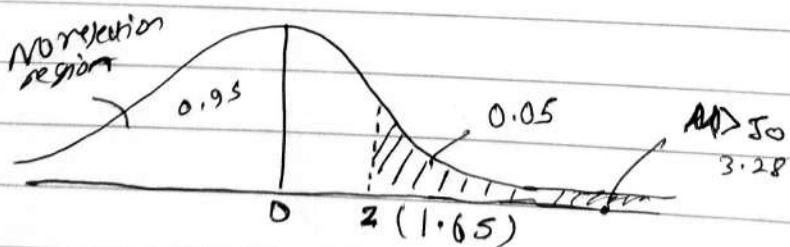
## Performing a Z-test example 1

Suppose a company is evaluating the impact of a new training program on the productivity of its employees. The company has data on the average productivity of its employees before implementing the training ~~new~~ program. The average production was 50 units per day with a known population std of 5 units. After implementing the training program, the company measures the productivity of a random sample of 30 employees. The sample has an average production of 53 units per day. The company wants to know if the new training program has significantly increased ~~productivity~~ Productivity.

So,

Given  $\mu = 50$ ,  $\sigma = 5$

- ①  $H_0: \mu = 50$ ,  $H_a: \mu > 50$
- ②  $\alpha = 0.05 \rightarrow 5\%$
- ③ Normality valid / pop std ( $\sigma$ ) known
- ④ Z-test
- ⑤  $Z$
- ⑥  $Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} = \frac{53 - 50}{5/\sqrt{30}} = \frac{3}{\sqrt{15}} = 3.28$



- ⑦ We can reject the null hypothesis.

$(\mu > 50)$   
will be

## Example 2

Suppose a snack food company claims that their 10gms water packet contain an average weight of sugar per packet. To verify this claim, a consumer ~~watching~~ watchdog organization decide to test a random sample of 40 water packets. The organization want to determine whether the actual average weight differs significantly from the claimed 50gms. The organization collects a random sample of 40 water packets and measure their weights. They find that the sample has an average weight of 49gms, with a known population std of 4gms.

Soln,

$$\mu = 50, n = 40, \bar{x} = 49, \sigma = 4$$

①  $H_0 : \mu = 50, H_a : \mu \neq 50$

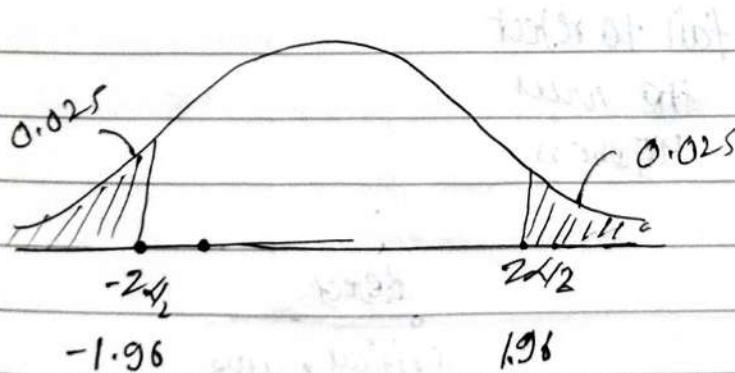
②  $\alpha = 0.05$

③ Normality valid /  $\sigma$  known

④ Z test

⑤ Z

$$Z = \frac{49 - 50}{40/40} = \frac{-10}{4} = -2.5$$

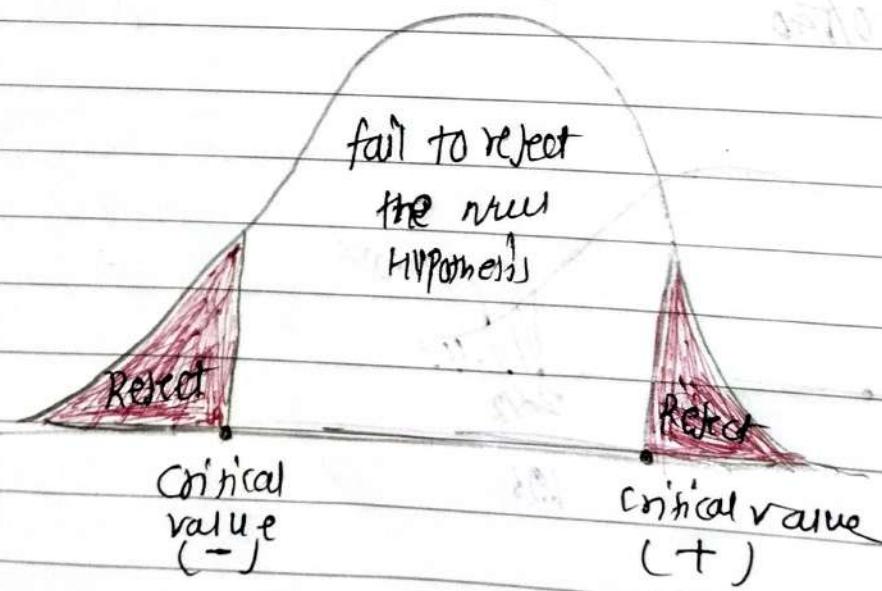


We can reject Null Hypothesis

Significance level  $\rightarrow$  denoted by alpha ( $\alpha$ ), is a predetermined threshold used in hypothesis testing to determine whether the null hypothesis should be rejected or not. It represents the probability of rejecting the null hypothesis when it is actually true, also known as type I error.

Rejection Region  $\rightarrow$  The rejection region is the region of values that corresponds to the rejection of the null hypothesis at some chosen probability level.

$\rightarrow$  Rejection region (aka critical region) is an area under the null hypothesis ( $H_0$ ) to detect by default.



## Problem with rejection region approach.

- ① result depend on the chosen significant level.
- ② Ver hi 'reject' ya 'accept' ka decision data hai, evidence ka strength nahi data hai
- ③ sample size se heavily effect hote hai

### Type I vs Type II error

→ In hypothesis testing, there are two types of error that can occur when making a decision about the null hypothesis.

#### Type I (false Negative)

→ Type I error occur when the sample result lead to rejection of the null hypothesis when it is in fact true.

In other words, it's the mistake of finding a significant effect or relationship when there is none. The probability of committing a type I error is denoted by  $\alpha$  (alpha), which is also known as the significance level. By choosing a significance level we can control the risk of making a Type I error.

#### → Note

Type I error matlab  $\Rightarrow$  Null hypothesis actually true tha, par tumne galat reject kar leya

## Type 2 error (false negative)

→ Type II error occurs when based on the sample result, the Null Hypothesis is not rejected when it is in fact false.

This means that the researcher fails to detect a significant effect or relationship when one actually exists. The probability of committing a type II error is denoted by  $\beta$  (Beta).

Note → Null hypothesis galat tha, par turne 'sahi' naa  
leya (fail to 'reject').

## Truth about the Population

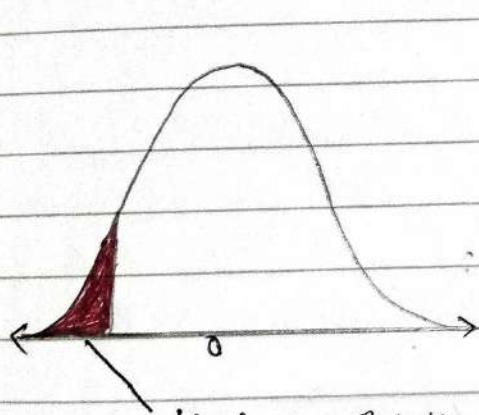
	H <sub>0</sub> true	H <sub>0</sub> false
Decision based on sample		
Reject H <sub>0</sub>	Type I error	Correct decision
Accept H <sub>0</sub>	Correct decision	Type 2 error

## One sided vs two sided test

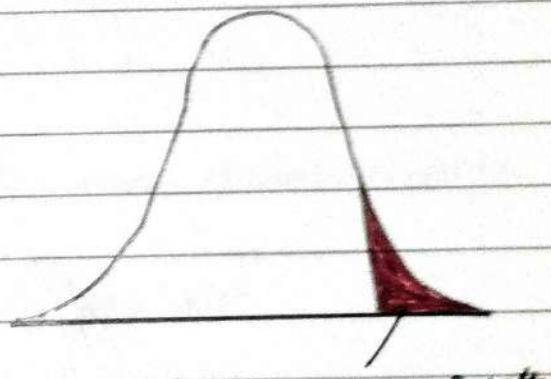
### One-sided (one tailed) test

↳ A one-sided test is used when the researcher is interested in testing the effect in a specific direction. (either greater than or less than the value specified in the null hypothesis). The alternate hypothesis is a one-sided test contains an inequality (either " $>$ " or " $<$ ").

Ex → A researcher want to test whether a new medication increase the avg recovery rate compared to the ~~existing~~ existing medication.



↳ left tail test  
H<sub>a</sub> :  $\mu <$  Value

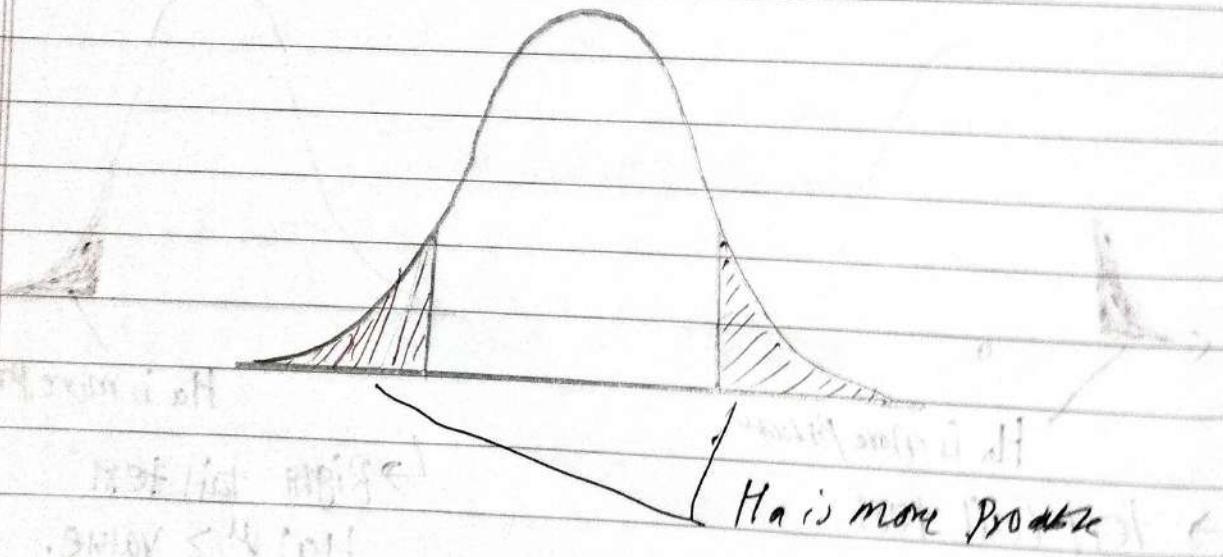


↳ Right tail test  
H<sub>a</sub>:  $\mu >$  value.

## Two-sided [two-tailed] test:

A two-sided test is used when the researcher is interested in testing the effect in both ~~sides~~ direction (ie. whether the value specified in the null hypothesis is different, ~~either~~ greater or lesser). The alternate hypothesis is a two-sided test containing a "not equal to" sign ( $\neq$ )

Ex → A researcher wants to test whether a new medicine has diff avg recovery rate compared to the existing medication.



The main difference between them lies in the directionality of the alternate hypothesis and how the significant level is distributed in critical region.

## P-value Approach

→ P-value is a measure of the strength of the evidence against the null hypothesis that is provided by our sample data.

### Steps in P-value approach

- ① state Hypothesis
  - ② choose test type like Z-test, t-test, chi-square test, ANOVA
  - ③ find the statistic for test:  
on Z-test
- $$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$
- ④ find P-value Using table:
    - Z-table ya T-table se test statistic ke area (probability) dekho.

Us area ko depend karta hai test ka direction par

Right tailed : ~~P-value~~ P-value = 1 - Z-table value

left tailed : P-value = Z-table value

2 tailed : P-value =  $2 \times (1 - \text{table one})$  (agar dobara tall)

- ⑤ Compare with (significance level)  
P-value interpretation

$P < 0.01 \rightarrow$  very strong evidence against  $H_0 \rightarrow$  reject confidently  
 $0.01 \leq P \leq 0.05 \rightarrow$  moderate evidence against  $H_0 \rightarrow$  Reject  $H_0$

$0.05 \leq P \leq 0.10 \rightarrow$  weak evidence against  $H_0 \rightarrow$  fail to reject  $H_0$ , but keep suspicion  
 $P \geq 0.10 \rightarrow$  no evidence against  $H_0 \rightarrow$  fail to reject  $H_0$

P-value in context of z-test

SUPPOSE A COMPANY IS EVALUATING THE IMPACT OF NEW TRAINING PROGRAM ON THE PRODUCTIVITY OF ITS EMPLOYEES. THE COMPANY HAS DATA ON THE AVG PRODUCTIVITY OF ITS EMPLOYEES BEFORE IMPLEMENTING THE TRAINING PROGRAM. THE AVG PRODUCTIVITY WAS 50 UNITS PER DAY. AFTER IMPLEMENTING THE TRAINING PROGRAM, THE COMPANY MEASURED THE PRODUCTIVITY OF A RANDOM SAMPLE OF 30 EMPLOYEES. THE SAMPLE HAS AN AVERAGE PRODUCTIVITY OF 53 UNITS PER DAY AND POPULATION STD IS 4. THE COMPANY WANTS TO KNOW IF THE NEW TRAINING PROGRAM HAS SIGNIFICANTLY INCREASED PRODUCTIVITY.

Z-test

$$\mu = 50 \quad n = 30 \quad \bar{x} = 53 \\ \sigma = 4 \quad \alpha = 0.05$$

$$H_0: \mu = 50$$

$$H_a: \mu > 50$$

$$Z\text{-stat} = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{53 - 50}{4/\sqrt{30}} \approx 4.10$$

P-value

for  $Z \approx 4.11$  (right tailed)

$P\text{-value} \approx 0.0002$  (using  $\Phi(0)$ )

$P < \alpha$ , we can reject  $H_0$