

# Introduction To Statistics

①

## ① Basics To Advance

### (i) Descriptive stats

- { ① Measure of central Tendency }
- { ② Measure of Dispersion }

Summarizing the data.

Histograms, Pdf, Cdf, Probability,  
Permutation, mean, Median, Mode,  
Variance, Standard deviation

### (ii) Inferential stats

- z test  $\Rightarrow$  Python
- t test  $\Rightarrow$  Python

ANOVA  $\rightarrow$  F test

CHI SCUARE

Hypothesis Testing

{ P values }

① Gaussian Distribution

Confidence Intervals

② Lognormal Distribution

z table, t table

③ Binomial Distribution

④ Bernoulli's Distribution

⑤ Pareto Distribution (Power law)

⑥ Standard Normal Distribution

$\Rightarrow$  Python

⑦ Transformation and Standardization

⑧ Q-Q Plot

## What is Statistics?

Statistics is the science of collecting, organizing and analyzing data. { Better Decision Making }

Def. data??

Data: Facts or pieces of information that can be measured.

Eg: The IL of a class

{ 98, 97, 60, 75, 55, 65 }

Ages of students of a class

{ 30, 25, 24, 23, 27, 28 } → Data

## Types of Statistics

① Descriptive stats: It consists of organizing and summarizing data.

② Inferential stats: Techniques where in we used the data that we have measured to form conclusions.

① Classroom of Maths Student (20)

84, 86, 78, 72, 75, 65, 80, 81, 92, 95, 96, 97, ...

Eg: Descriptive stats.

① What is the average marks of the students in the class.

## Inferential stats

e.g:

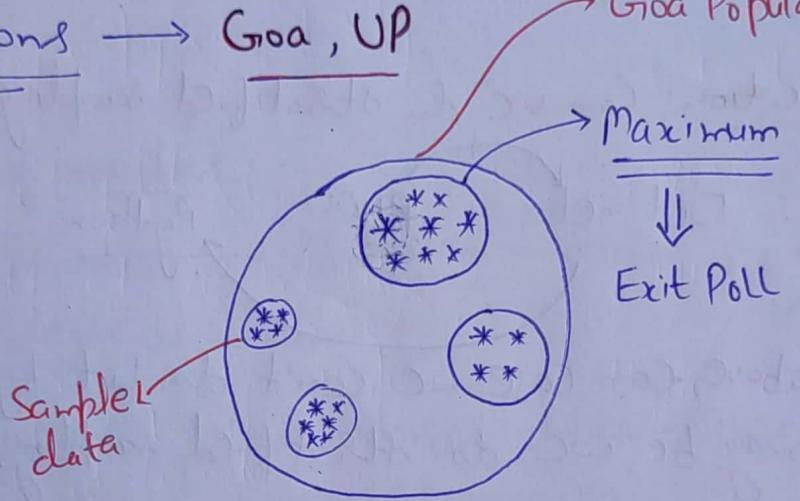
Are the marks of the students of this classroom similar to the age of the Maths classroom in the college?

Sample  $\boxed{1}$   $\xrightarrow{5}$  Population

(2)

## Population And Sample

Elections  $\rightarrow$  Goa, UP

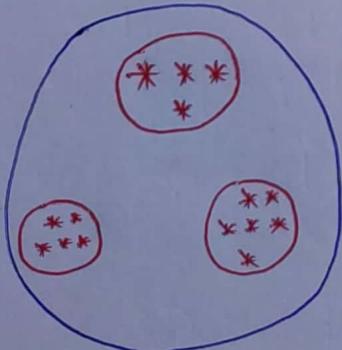


Population ( $N$ )

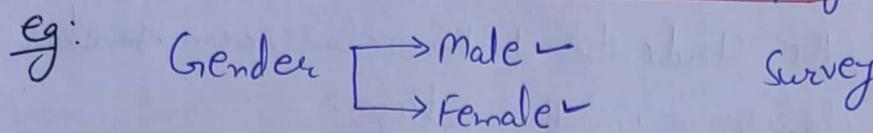
Sample ( $n$ )

## Sampling Techniques

- ① Simple Random Sampling: Every member of the population ( $N$ ) has an equal chance of being selected for your sample ( $n$ ).



② Stratified Sampling: Where the Population ( $N$ ) is split into non-overlapping groups. (strata)

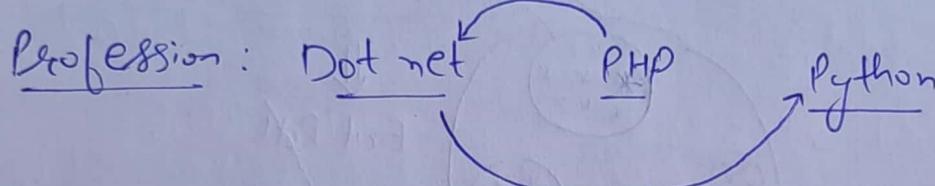


Age-group

(0-10) (10-20) (20-40) (40-100)

90, 100

Interview Question: Can we do stratified sampling based on profession?



In the above case we can't do but by applying some condition it may be we do stratified sampling.

Profession: Doctors, Engineers

In that case we do stratified sampling.

③ Systematic Sampling:

( $N$ ) →  $n$ th individual

Eg: Mall → Survey  
→ 8th person → Survey

④ Convenience Sampling:

Survey ← Only those people  
Data Science ← use case

Eg: Exit Poll  
{Random Sampling}

RBI → Household Survey  
↓  
Survey ⇒ Women

Eg: Drug  $\rightarrow$  Tested  $\Rightarrow$  Stratified + Convenience  
 $\Downarrow$

③

What kind of Sample?

Variables: A variable is a property that can take any value.  
 $\{182, 178, 168, 150, 160, 170\}$

Eg: Height

$$\text{Weight} = \{78, 99, 100, 60, 50\}$$

Two kinds of variables:

- ① Quantitative variable  $\rightarrow$  Measured numerically, {add, subtract, multiply, divide}  
② Qualitative / Categorical variables

L  $\hookrightarrow$  Eg: Gender  $\begin{cases} H \\ F \end{cases}$  {Based on some characteristics we can derive categorical variables.}

Blood group

T-shirt size

A +ve

L

XL

M

S

Eg: IQ

$\frac{0-10}{\Downarrow}$

$\frac{10-50}{\Downarrow}$

$\frac{50-100}{\Downarrow}$

Less IQ

Medium IQ

Good IQ

Quantitative

Discrete variable

Continuous Variables

Eg: Whole number

Eg: Height = 172.5, 162 cm., 163.5 cm.

No. of Bank Accounts

Weight = 100 kg, 99.5, 99.75

Eg: 2, 3, 4, 5, 6, 7,

Rainfall = 1.1, 1.25, 1.35 --

Total children in a family

Eg: 2, 3, 4, 5

Eg: ① What kind of variable is Gender ?? Categorical

② What " " " Marital status? "

③ " " " River length? Continuous

④ " " " Population of the state is? Discrete

⑤ " " " Song length? Continuous

⑥ " " " Blood pressure? Continuous

⑦ PIN Code? {Discrete or Categorical}

## \* Variable Measurement Scales

4 types of measured variable

① Nominal data - {Categorical data} → classes

② Ordinal - Order of the data matters, value does not.

③ Interval - Order matters, value also matters, natural zero is

④ Ratio data  
eg(0 K temp.) not present.



the measurement have  
a true natural zero.

Interview question: Difference b/w  
ordinal and nominal data

Eg: Student marks - 0, 10, 20, 30, 70

Eg:

Student (Marks)

Rank

100

1

96

2

57

4

85

3

44

5

→ Ordinal  
Data

## Frequency Distribution:

(4)

Sample dataset :- Rose, Lilly, Sunflower, Rose, Lilly, Sunflower, Rose, Lilly, Lilly.

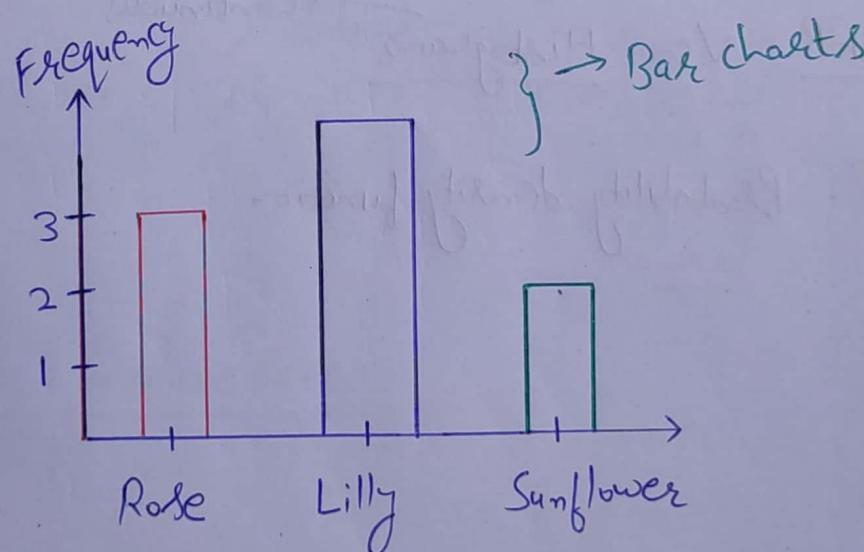
<u>Flower</u>	<u>Frequency</u>	<u>Commulative Frequency</u>
Rose	3	3
Lilly	4	7
Sunflower	2	9

### Interview

When there is discrete variables :- we draw a bar graphs

When " " continuous variables :- we draw a histogram.

### ① Bar Graph

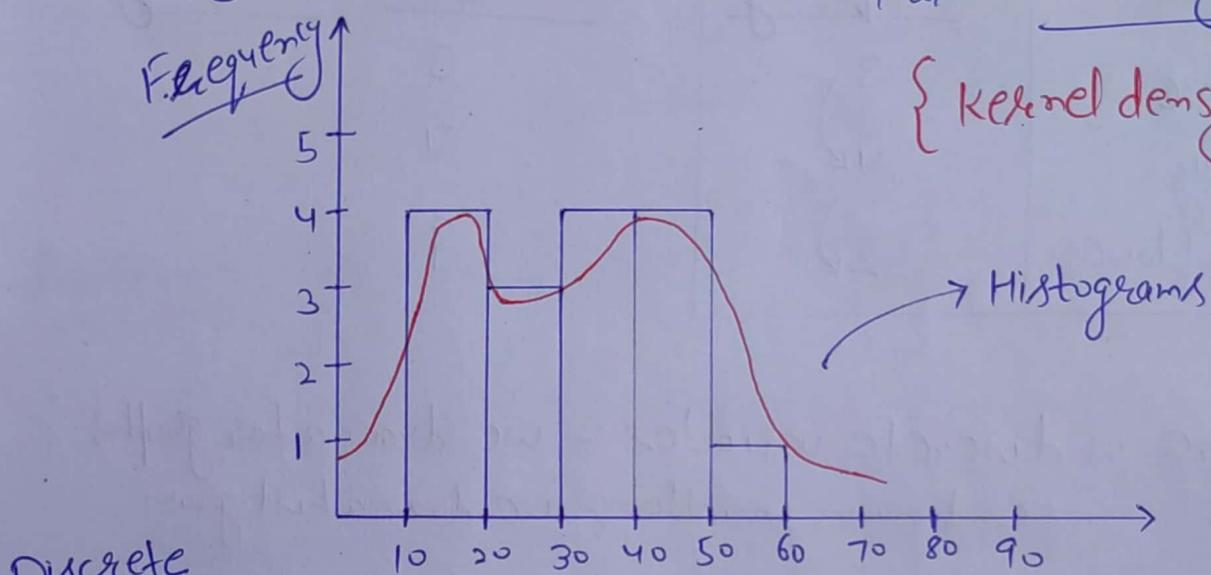


For Summarizing the data we plot a Bar Graph

## ② Histograms :- Continuous

Ages = { 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51 }

↳ Bins = 10



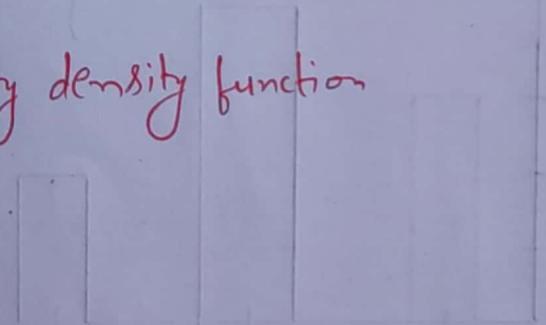
Discrete

BAR

VS Histograms

→ Continuous

Pdf: Probability density function



# Basics To Intermediate Stats

(5)

① Measure of Central Tendency

② Measure of dispersion

③ Gaussian Distribution

④ Z score

⑤ Standard normal Distribution

① Arithmetic Mean for population & sample

Mean (Average)

Population (N)

Population  
data

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

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

$$\text{Population mean} = \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$= \frac{32}{10} = 3.2$$

Sample (n)

$$\bar{x} = \sum_{i=1}^n \frac{x_i}{n}$$
$$= 3.2$$
$$=$$

\* Central Tendency: Central tendency refers to the measure used to measure the "centre" of the distribution of data.

① Mean    ② Median    ③ Mode

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

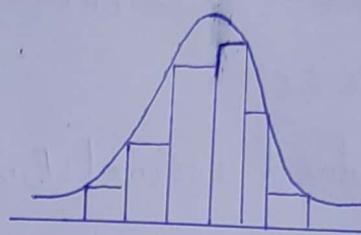
$$\text{Mean} = \frac{32 + 100}{11}$$

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

$$\begin{aligned} \mu &= 3.2 \\ &+ \\ &100 \end{aligned}$$

$$\mu = 12$$

Distribution



② Median:  $\{1, 1, 2, 2, \boxed{3}, 3, 4, 5, 5, 6, \boxed{100}, 112\}$

Steps: ① Sort all the numbers ✓

② Find the central element

→ odd length  
→ even length

$$\checkmark \quad \mu = 12$$

$$\text{median} = 3$$

$$\text{median} = 3.5$$

$$\text{Avg} = \frac{3+4}{2} = 3.5$$

find the average of two middle elements.

$$\begin{array}{l} \mu = 3.2 \\ \mu = 12 \end{array}$$

$$\begin{array}{l} \text{median} = 3 \\ \text{Median} = 3.5 \end{array}$$

\* { Median works well with outliers }

③ Mode: { most frequent elements }

⑥

{ 1, 2, 2, 3, 4, 5, 6, 6, 6, 7, 8, 100, 100, 100, 100 }  
2 3 ↓↓↓↓

[mode = 6] → Measure of central tendency

when to use:

① Outliers → median

② Categorical feature → mode

eg: Type of flower Petal length Petal width Data set  
Rose { Missing → Most frequent } ↓  
Lilly Value occurring element 10% missing data  
Sunflower ↓  
—  
—  
— Categorical variable

Ages of Students

age

25

Mean? ✓

26

⇒ Median?

—

—

32

34

38

Mode ??

## (\*) Measure of Dispersion

① Variance

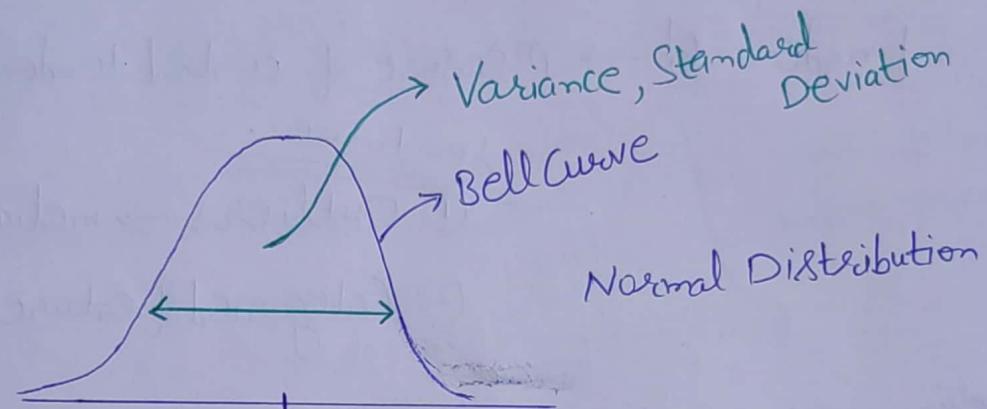
② Standard Deviation

$$\{ \text{Dispersion} \} = \frac{\mu}{\text{Spread}}$$

$\downarrow$

$\{ 1, 1, 2, 2, 4 \}$        $\frac{10}{5} = 2$

$\{ 2, 2, 2, 2, 2 \}$        $\frac{10}{5} = 2$



Mean, Median, Mode

$\downarrow$   
Measure of central  
Tendency

Population ( $N$ )

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

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

Sample ( $n$ )

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

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

Bessel's Correction  
Degree of freedom

# ① Variance ✓

①

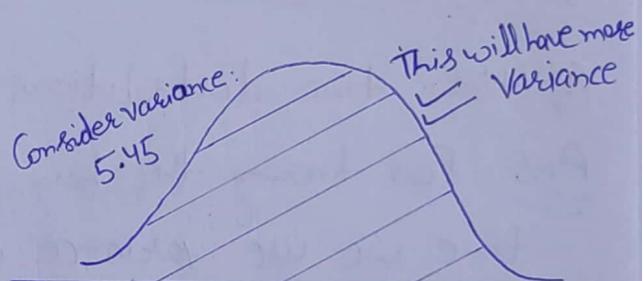
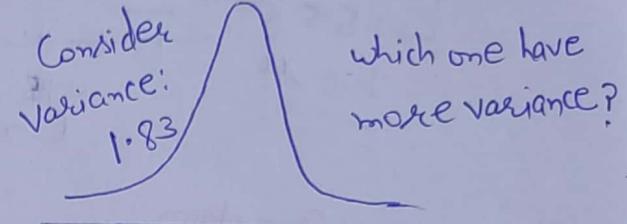
## Population Variance

$$\sigma^2 = \frac{N}{\sum_{i=1}^N} \frac{(x_i - \mu)^2}{N} = \frac{10.84}{6} = 1.81$$

$x$	$\mu$	$x - \mu$	$(x - \mu)^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
2	2.83	+0.83	0.6889
3	2.83	0.17	0.03
4	2.83	1.17	1.37
5	2.83	2.17	4.71
			<u>10.84</u>
	<u><math>\mu = 2.83</math></u>		

## Sample Variance

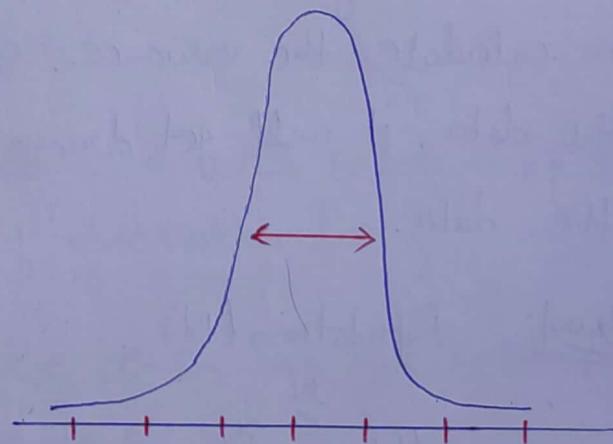
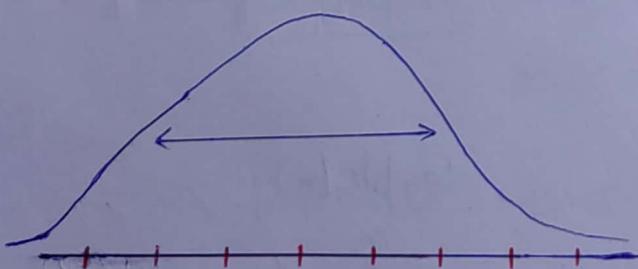
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$



\* more variance value: more spread.

more variance = more spread

less variance = less spread but  
value more peak



graph Variance → Big Number  
S.D. → Big Number  
↓  
means more spread

graph Variance → Small number  
S.D. → Small number  
↓  
means less spread.

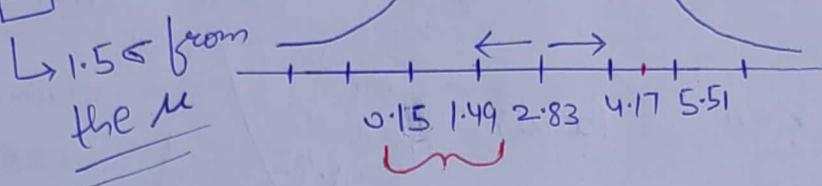
## ② Standard Deviation:

$$S.D = \sigma = \sqrt{\text{variance}} = \sqrt{1.81} = 1.345$$

variance = spread

Let's Consider

$$S.D. = 5$$



$$\begin{array}{r} 2.83 \\ + 1.34 \\ \hline 4.17 \end{array} \quad \begin{array}{r} 2.83 \\ - 1.34 \\ \hline 1.49 \end{array}$$

$$\begin{array}{r} 4.17 & 1.49 \\ + 1.34 & 1.34 \\ \hline 5.51 & 0.15 \end{array}$$

## Interview Question:

① How two distributions are different?

Ans For knowing the how two distributions are different at that time we use variance and standard deviation.

Most important

② Why sample variance is divided by  $n-1$ ?

Ans The reason dividing by  $n-1$  corrects the bias is because we are using the sample mean, instead of the population mean, to calculate the variance. Since the sample mean is based on the data, it will get drawn toward the centre of mass for the data.

Proof: Population ( $N$ )

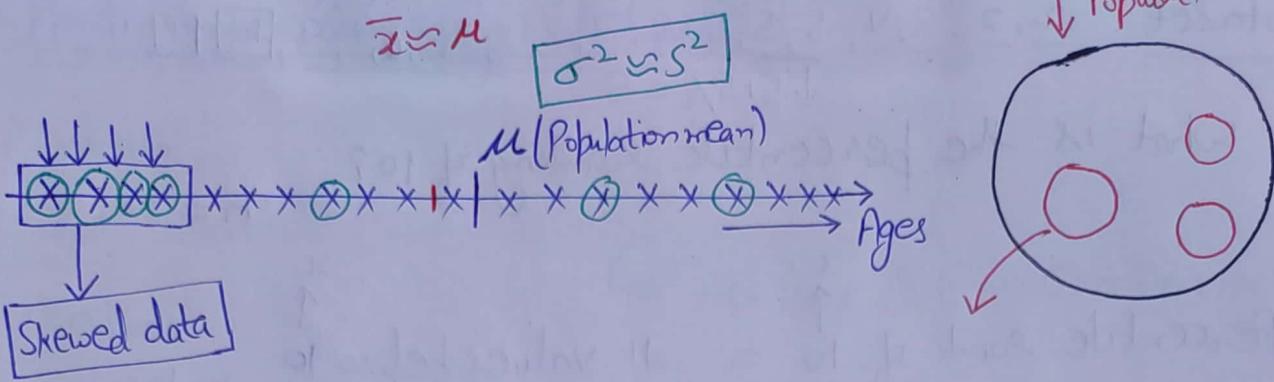
$$\mu = \frac{1}{N} \sum_{i=1}^N x_i$$

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

Sample ( $n$ )

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

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



$n-1, n-2, n-3, n-4, \dots$

## ④ Percentiles and Quartiles

Percentage: 1, 2, 3, 4, 5

% of the numbers that are odd?

$$\% \text{ of odd} = \frac{\# \text{ of numbers that are odd}}{\text{Total numbers}}$$

$$= \frac{3}{5} = 0.6 = 60\%$$

Percentiles: { GATE, CAT, SAT }

→ Definition: A percentile is a value below which a certain percentage of observation lie. }

95 percentile means that the person has got better marks than 95% of the entire student.

Dataset: 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12  
 $\downarrow 5+5/2 = 5$

① What is the percentile ranking of 10?

$$n=20$$

$$\text{Percentile rank of } 10 = \frac{\# \text{ values below } 10}{n} \times 100$$
$$= \frac{16}{20} \times \cancel{100}^5 = \underline{\underline{80 \text{ percentile}}}$$

What is the percentile ranking of 11?

$$\text{Percentile rank of } 11 = \frac{\# \text{ values below } 11}{n} \times 100$$
$$= \frac{17}{20} \times \cancel{100}^5 = \underline{\underline{85 \text{ percentile}}}$$

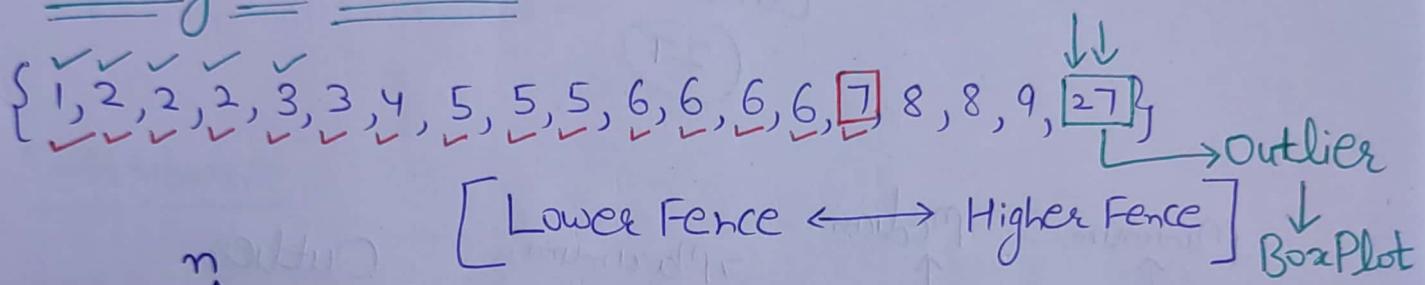
② What value exists at percentile rank of 25?

$$\text{value} = \frac{\text{Percentile}}{100} \times (n+1)$$
$$= \frac{25}{100} \times 21 = 5.25 \Rightarrow \underline{\underline{\text{index position}}} \\ \downarrow \\ \underline{\underline{5 \text{ Answer}}}$$

# \* Five Number Summary And Box Plot

- ① Minimum
  - ② First Quartile (25%) Q1
  - ③ Median
  - ④ Third Quartile (75%) Q3
  - ⑤ Maximum
- }  $\Rightarrow$  Box Plot  $\rightarrow$  Outlier detect

## Removing the Outliers



$$Q1 = \frac{25}{100} \times (19+1) = \frac{25}{100} \times 20 = 5 \Rightarrow \text{index}$$

$\uparrow$        $\downarrow$

$$(25\%) Q1 = 3$$

$$IQR = Q3 - Q1$$

$$\text{Lower Fence} = Q1 - 1.5(IQR)$$

$$\text{Higher Fence} = Q3 + 1.5(IQR)$$

$$Q3 = \frac{75}{100} (19+1)$$

$$= \frac{75}{100} \times 20$$

$$= 15 \Rightarrow \text{index}$$

$$IQR = 4$$

$$(75\%) Q3 = 7$$

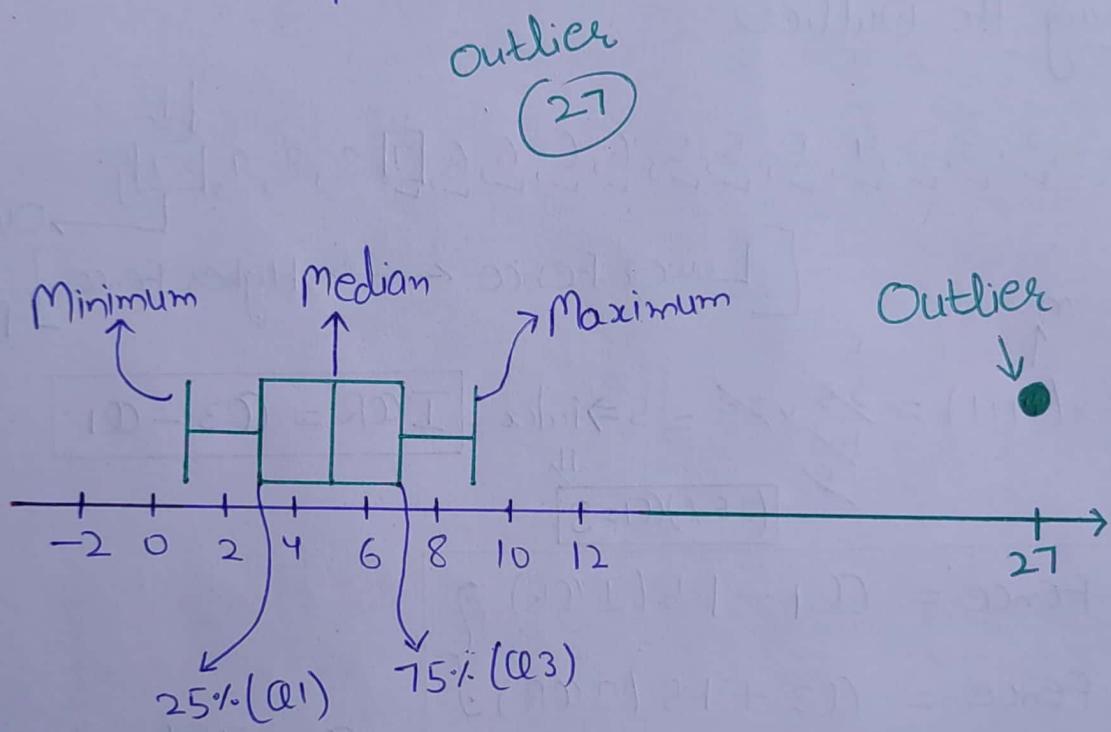
$$L.F. = 3 - 1.5(4) = 3 - 6 = -3$$

$$[-3 \leftrightarrow 13]$$

$$H.F. = 7 + 1.5(4) = 7 + 6 = 13$$

## Box Plot

- ① Minimum = 1
- ② First Quartile (25%) Q1 = 3
- ③ Median = 5
- ④ Third Quartile = 7
- ⑤ Maximum = 9



Interview Question: What is Box Plot and applications?

Ans: A box-plot displays the five-number summary of a set of data. The five-number summary is the minimum, first quartile, median, third quartile and maximum.

We draw a box from the first quartile to the third quartile. A vertical line goes through the box at the median.

- Applications:
- ① It is use to remove the outliers.
  - ② It provides us visualization way where the outlier is present.

# Advance Statistics

(10)

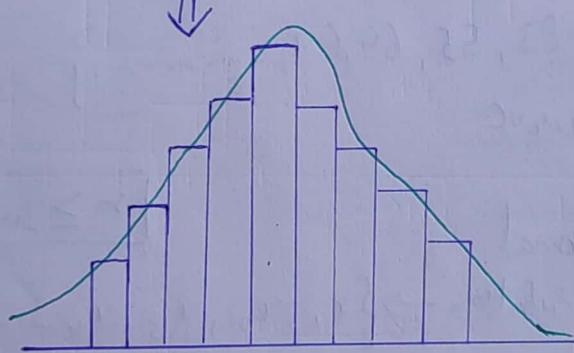
- ① Distributions (Gaussian Distribution)
- ↳ Normal Distribution
  - ↳ Standard Normal Distribution
  - ↳ Z score
  - ↳ Log Normal Distribution
  - ↳ Bernoulli's Distribution
  - ↳ Binomial Distribution

## Practical

- ① Mean, Median, Mode
- ② Variance, Standard Deviation
- ③ Histogram, Pdf, Bar Plot, Violin Plot
- ④ IQR
- ⑤ Log Normal Distribution

## ① Distribution

$$\text{Ages} = \{24, 26, 27, 25, 30, 32, \dots\}$$



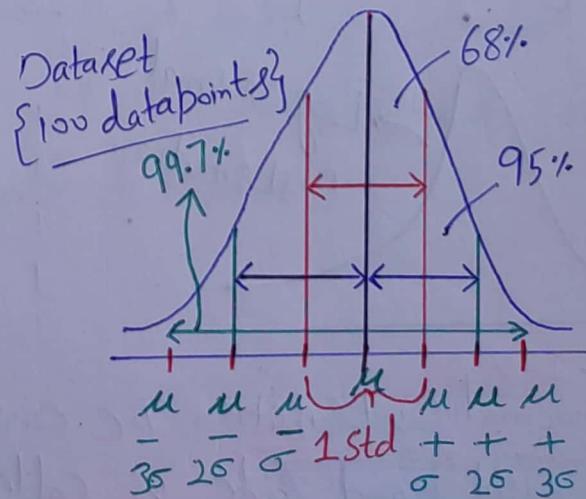
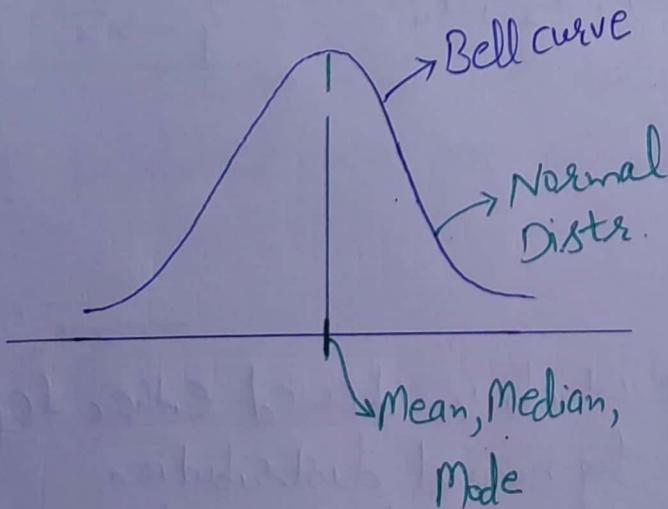
Interview

[3-sigma Rule]  
or  
[Empirical Formula]

68-95-99.7% Rule

Interview

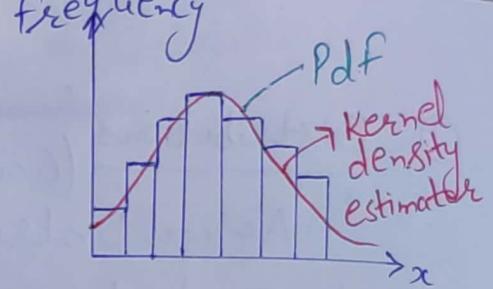
## Normal / Gaussian Distribution



Symmetrical: Left area of the curve  
= Right area of the curve

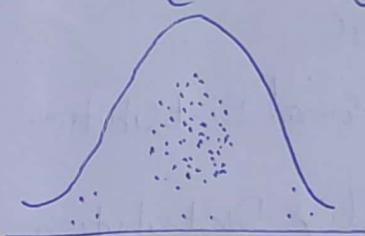
eg: ① Height → Normally Distributed

↓  
Domain Expert → {Doctor}



② Weight    ③ IRIS DATASET ⇒ {sepal length, petal length}

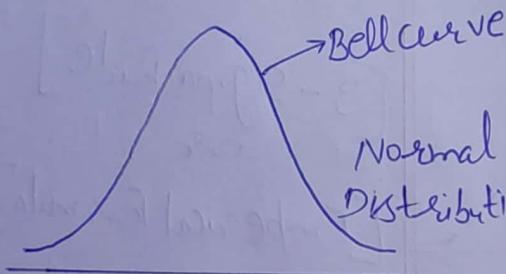
④ Laptop Touchpad →



⑤ inewron Team meetup ⇒ T-shirt ⇒ L, M, XL, XXL

Interview Question: Central limit Theorem

Prob:  $x = \{65, 72, 83, 55, 64, 67, \dots\}$



$n \geq 30 \Rightarrow$  Sample size

$y = \{\dots\}$  weather

$\rightarrow \{x_1, x_2, x_3, x_4, \dots\} \rightarrow \bar{x}_1$

$\rightarrow \{x_1, x_3, x_7, x_5, \dots\} \rightarrow \bar{x}_2$

$\rightarrow \{x_1, x_3, x_4, \dots\} \rightarrow \bar{x}_3$

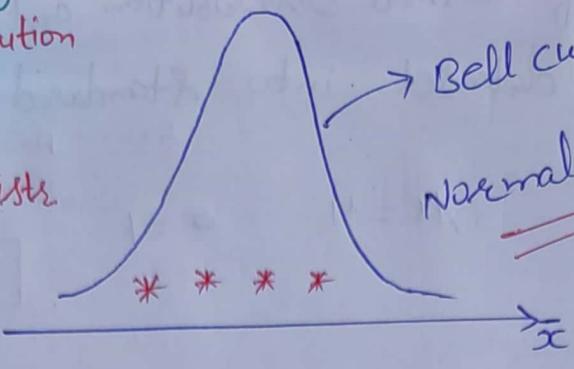
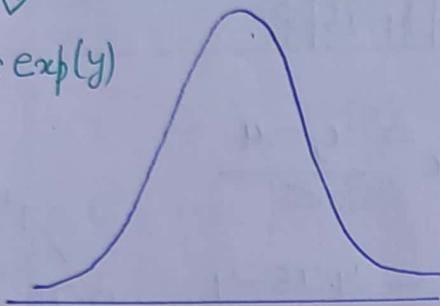
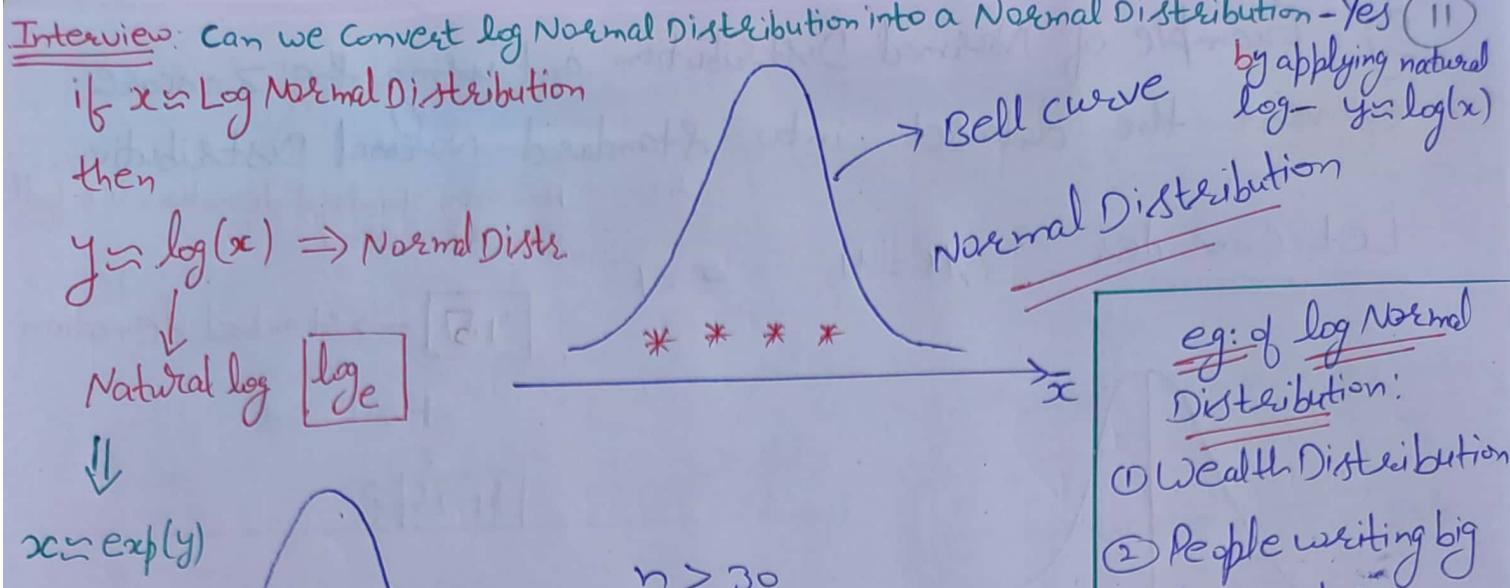
$\vdots$

$\rightarrow \bar{x}_n$

\* Interview

log normal Distribution: The curve which is skewed either left or right hand side is called log normal distribution.

In above example, it is skewed towards the right hand side.



Sample Distribution

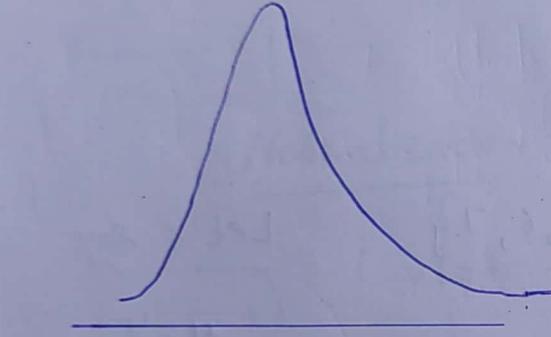
$\downarrow$

$n \geq 30$

Normal Distribution

$\Rightarrow$

$\mu_{\bar{x}}$



### Definition of Central Limit Theorem:

It says whether your distribution is Gaussian distribution or non-Gaussian distribution, if we take the sample where  $n \geq 30$  and we take all the sample mean and we populated we get a normal distribution.

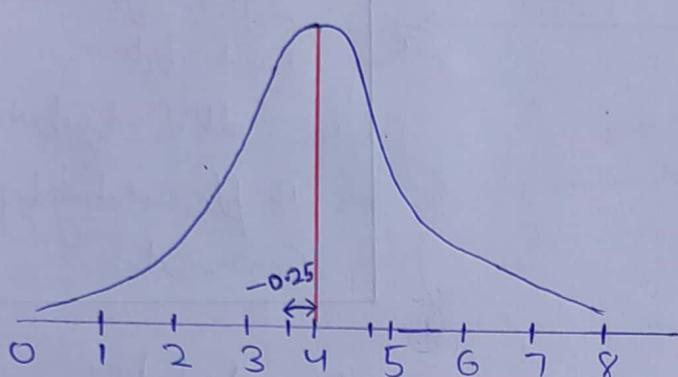
e.g.: Example of Normal Distribution and using the z-score  
convert the dataset into Standard Normal Distribution.

Interview How many standard deviation 4.75 fall from the mean?

Let Consider

$$\mu = 4$$

$$\sigma = 1$$



4.75 → standard deviation

+ 0.5 s.d

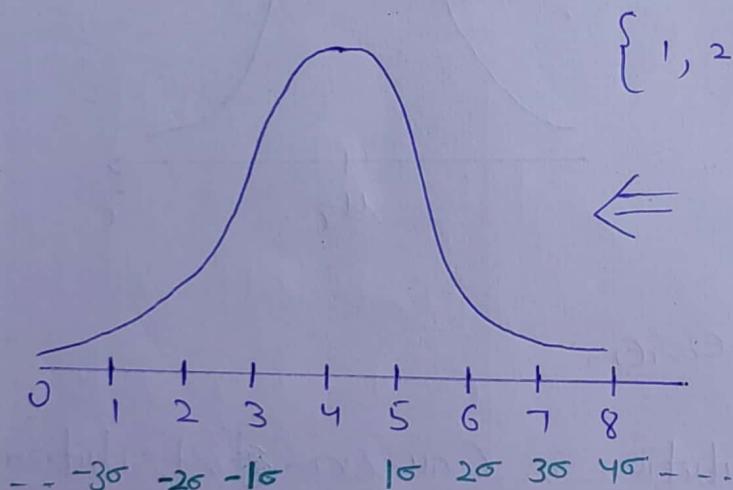
4.75 ??

$$Z_{\text{score}} = \frac{x_i - \mu}{\sigma}$$

$$= \frac{4.75 - 4}{1} = 0.75 \text{ s.d.}$$

3.75

$$Z_{\text{score}} = \frac{3.75 - 4}{1} = -0.25$$



$$\{1, 2, 3, 4, 5, 6, 7\}$$

$$Z_{\text{score}} = \frac{x_i - \mu}{\sigma}$$

Let's say

$$\mu = 4$$

$$\sigma = 1$$

$$\{-3, -2, -1, 0, 1, 2, 3\}$$

$\rightarrow \text{SND} (\mu = 0, \sigma = 1)$

$$Z(1) = \frac{1-4}{1} = -3 \quad Z(3) = \frac{3-4}{1} = -1$$

$$Z(2) = \frac{2-4}{1} = -2$$

$$\{1, 2, 3, 4, 5, 6, 7\}$$

$\downarrow$   
z score

Standard Normal Distribution

$$(\mu = 0, \sigma = 1)$$

$\{-3, -2, -1, 0, 1, 2, 3\}$  ← Satisfying this property

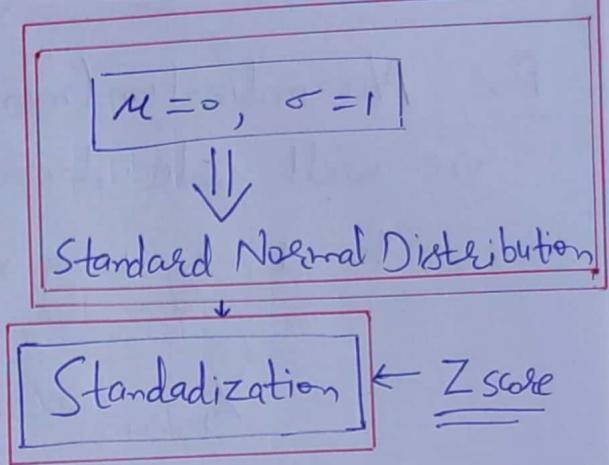
## Practical Application

(12)

DATA SET

Age	Salary	Weight
24	40k	70
25	80k	80
26	60k	55
27	70k	45

Z score



Normalization

$\{ \mu = 0, \sigma = 1 \}$

$\xrightarrow{(-1 \text{ to } 1)}$   $\xrightarrow{(0 \text{ to } 1)}$

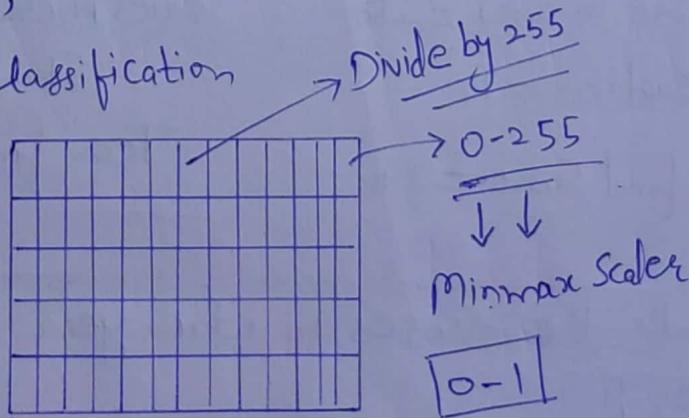
By using Min Max Scaler formula we are able to convert this into  $(0 \text{ to } 1)$ .

Where we use Normalization?

In Deep Learning,

CNN → Image Classification

Pixels  
Normalization →



Interview Question: Difference between Normalization and Standardization?

Ans Normalization (min-max Normalization): In this approach we will scale down the values of the features between 0 to 1 by using min-max scaler formula.

$$X_{\text{norm}} = \frac{X - X_{\text{min}}}{X_{\text{max}} - X_{\text{min}}}$$

Standardization (Z-score Normalization): Here all the features will be transformed in such a way that it will have the properties of standard normal distribution with mean( $\mu$ ) = 0 and standard deviation( $\sigma$ ) = 1 using the Z-score formula.

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

\* Practical Example: { India Vs SA }

① ODI Series (2021)

Series Average Score in 2021 = 250  
Standard Deviation = 10

Series < Team final Score = 240

In 2020

Series Average Score in 2020 = 260  
Standard Deviation = 12

Team final Score = 245

Compare both the scores in which year Team final score was better?

Ans In 2021,

$$Z\text{ score} = \frac{x_i - \mu}{\sigma} = \frac{240 - 250}{10} = \frac{-10}{10} = -1$$

In 2020,

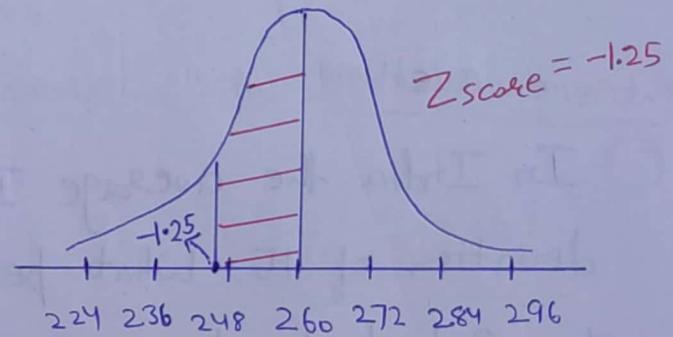
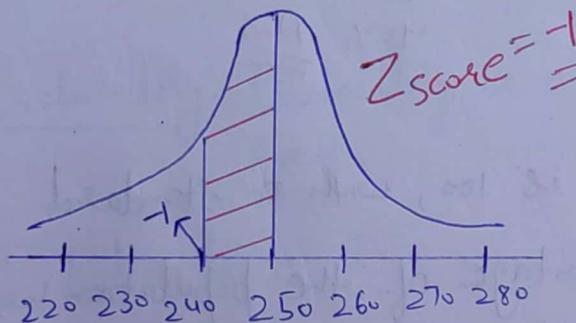
$$Z\text{ score} = \frac{x_i - \mu}{\sigma} = \frac{245 - 260}{12} = \frac{-15}{12} = -1.25$$

In 2021 ✓)

$$\mu = 250 \quad x_i = 240 \quad \sigma = 10$$

In 2020 ✓

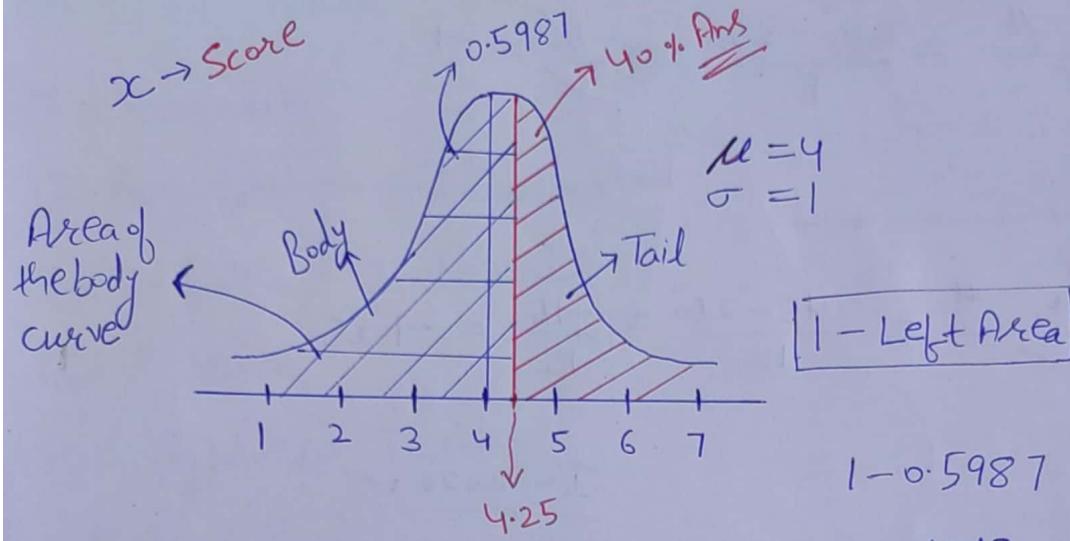
$$\mu = 260 \quad x_i = 245 \quad \sigma = 12$$



Standard deviation is greater in 2020 as compare to in 2021, So in 2020 teams final score is better than the 2021

Important

Interview Question: ① What percentage of scores falls above 4.25?



$$Z = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25$$

$$1 - 0.5987$$

$$= \underline{\underline{0.4013}}$$



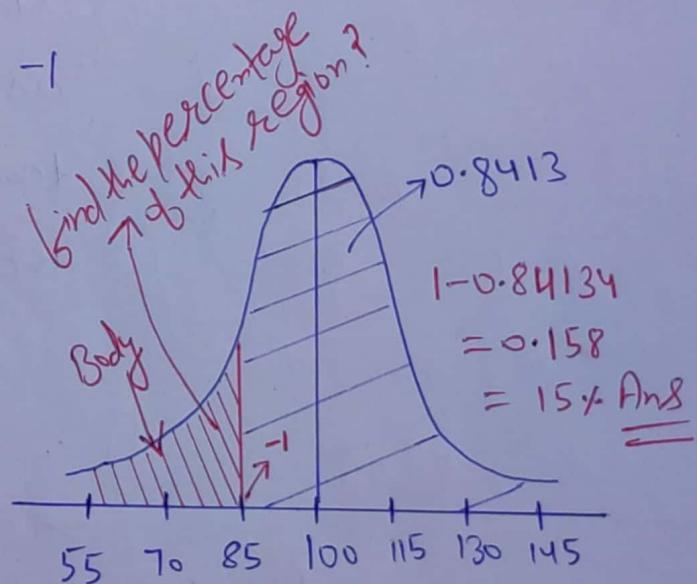
40% Ans

Interview Question

- ② In India the average IQ is 100, with a standard deviation of 15. What percentage of the population would you expect to have an IQ lower than 85?

Ans

$$Z = \frac{85 - 100}{15} = \frac{-15}{15} = -1$$



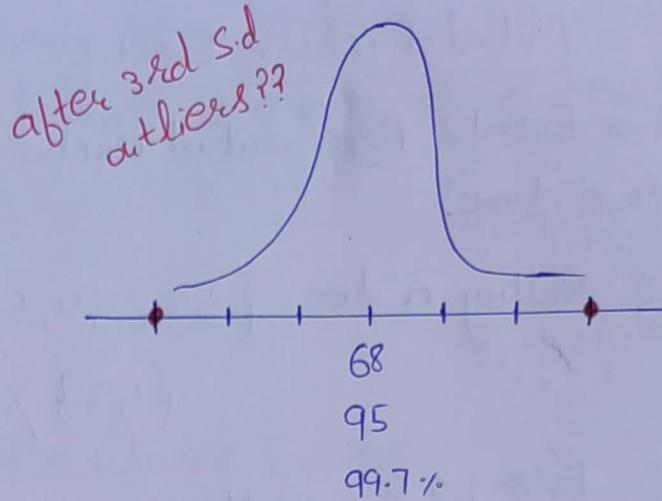
Also came in interview

↓

IQ 90 to 120

# - STATS Today's Agenda

- ① IQR - Python
- ② Probability
- ③ Permutation and Combination
- ④ Confidence Intervals
- ⑤ P value
- ⑥ Hypothesis Testing



$$\boxed{Z\text{ score} = \frac{x_i - \mu}{\sigma}}$$

\* Probability: Probability is a measure of the likelihood of an Event.

e.g.: Roll a dice  $\{1, 2, 3, 4, 5, 6\}$

$$\begin{aligned} \text{Pr}(6) &= \frac{\# \text{ of way an even can occur}}{\# \text{ of possible outcome}} \\ &= \frac{1}{6} \end{aligned}$$

Toss a Coin  $\{H, T\}$

$$\boxed{\text{Pr}(H) = \frac{1}{2}}$$

## ② Addition Rule: (Probability, "or")

### Mutual Exclusive Event

Two Events Are mutual exclusive if they cannot occur at the same time.

Eg: Rolling a dice  $\{1, 2, 3, 4, 5, 6\}$   
 $\{1, 2\} \times$

Tossing a coin  $\{H, T\}$

### Non Mutual Exclusive

Multiple events can occur at the same.

Eg: Deck of Cards  $\{\heartsuit, \diamondsuit\}$

### Example for Mutual Exclusive

① If I Toss a coin, what is the probability of the coin landing on heads or tails?

Ans

### Mutual Exclusive

### Addition Rule

$$Pr(A \text{ or } B) = Pr(A) + Pr(B)$$

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

$$\boxed{Pr(A \text{ or } B) = 1}$$

$\Rightarrow$  Roll a Dice

$$\begin{aligned} Pr(1 \text{ or } 3 \text{ or } 6) &= Pr(1) + Pr(3) + Pr(6) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \\ &= \cancel{\frac{3}{6}}_2 = \frac{1}{2} = 0.5 \end{aligned}$$

## Example for Non Mutual Exclusive

(15)

You are picking a card randomly from a deck. What is the Probability of choosing a card that is queen or a heart?  $\rightarrow (52)$

### Ans Non Mutual Exclusive

$$P(A) = \frac{4}{52} \quad P(V) = \frac{13}{52} \quad P(A \text{ and } V) = \frac{1}{52}$$

Addition Rule for non mutual exclusive Events

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

$$P(Q \text{ or } V) = P(Q) + P(V) - P(Q \text{ and } V)$$

$$= \frac{4}{52} + \frac{13}{52} - \frac{1}{52}$$

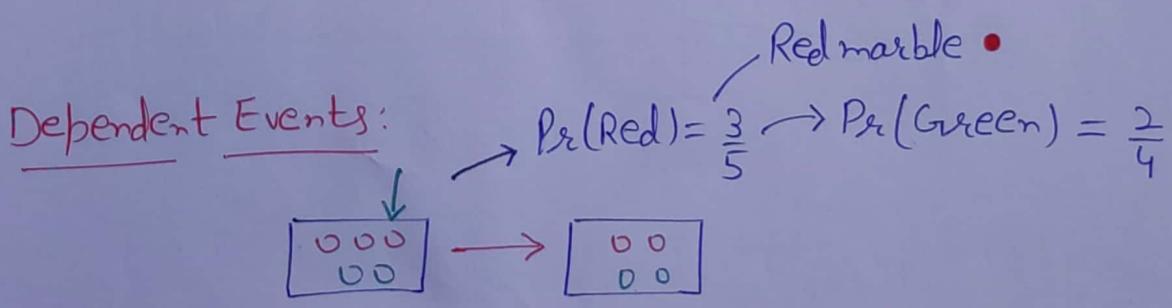
$$= \frac{16}{52} = \frac{4}{13}$$

### ③ Multiplication Rule

{Independent Events}

Eg: Rolling a Dice { $\overset{\uparrow}{1}, \overset{\downarrow}{2}, 3, 4, 5, 6$ }

1, 1, 2, Each and Every are independent



Naive Baye's {Conditional Probability}

Also called (Baye's Theorem)

## Independent Events

✳ What is the probability of rolling a "5" and then a "4" in a dice?

Ans Independent Event

Multiplication Rule

$$Pr(A \text{ and } B) = P(A) * P(B)$$

$$Pr(5 \text{ and } 4) = \frac{1}{6} * \frac{1}{6} = \frac{1}{36}$$

✳ What is the probability of drawing a Queen and then a Ace from a deck of cards?

$$\boxed{\textcircled{Q} \textcircled{Q} \textcircled{Q}} \quad \frac{3}{5} \quad \boxed{\frac{2}{4}}$$

Ans Dependent

$$P(G \text{ and } R) = P(G) * P(R|G)$$

↳ event  
occurred

$$P(A \text{ and } B) = P(A) + P(B|A) \xrightarrow{\text{Conditional Probability}}$$

This is very important  
in Baye's Theorem

$$P(Q \text{ and } A) = P(Q) * P(A|Q)$$

$$= \frac{4}{52} * \frac{4}{51}$$

# \* Permutation and Combination

(16)

## Permutation

School trip {chocolate factory} → Dairy, 5 star, Milky bar, Eclairs,  
Student {Assignment} Gem, Silk

Student ↗

$$\frac{6}{(5)} \times \frac{5}{(4)} \times \frac{4}{(3)} = \underline{\underline{120}}$$

Dairy, Gem, Milky

Milky, Gem, Dairy

## Permutation

$${}^n P_r = \frac{n!}{(n-r)!} = \frac{6!}{(6-3)!}$$

$$= \frac{6 \times 5 \times 4 \times 3!}{3!}$$

$$= \underline{\underline{120}}$$

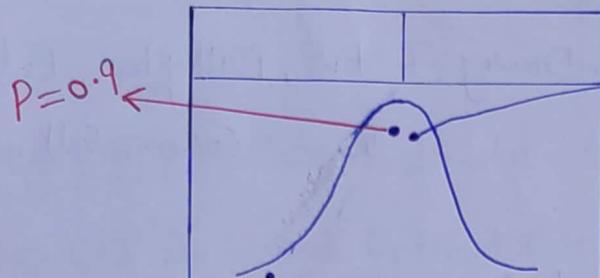
## \* Combination

Dairy	Gem	Eclair
—	—	—
—	—	—
—	—	—

$${}^n C_r = \frac{n!}{r!(n-r)!} = \frac{6!}{3!(6-3)!}$$

$$= \frac{6 \times 5 \times 4 \times 3!}{3 \times 2 \times 1 \times 3!} = \underline{\underline{20}}$$

Interview  
① P value { Many people get confused but its very easy }  
(Significance value)



P = 0.01  
↳ 1 time

Every 100 time I touch the mouse pad 80 times I touch this specific region.

Significance value derived from C.I.

Note: Significance value  $\neq$  P value

P value is also called Significance value: \*P value means you are just laying probability.

\*Significance says what should be within your C.I.

Definition of P value (Significance value): It is the probability for the Null Hypothesis to be true.

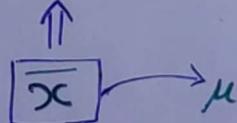
Interview

② Confidence Interval: The range of values that we observe in our sample and for which we expect to find the value that accurately reflects the population.

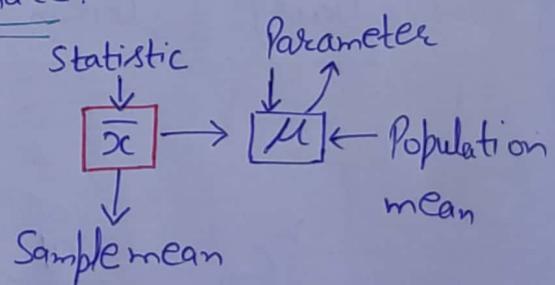
For finding the confidence interval its necessary to understand the concept of Point Estimate.

Point Estimate: The value of any statistic that estimates the value of a parameter is called Point Estimate.

Point Estimate



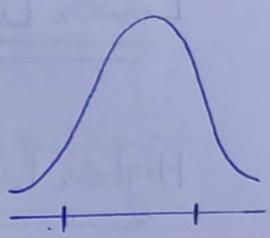
$$\begin{cases} \bar{x} \geq \mu \\ \bar{x} \leq \mu \end{cases}$$



$$\boxed{\text{Point Estimate}} \pm \boxed{\text{Margin of Error}} = \boxed{\text{Parameter} \Rightarrow \text{Population mean}}$$

Lower C.I = Point Estimate - Margin of Error

Higher C.I = Point Estimate + Margin of Error



$$\text{Margin of Error} = Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \Rightarrow \text{Standard Error}$$

$\rightarrow \text{Population Sd}$

$\alpha = \text{Significance value}$

Interview Question: Average size of sharks in the sea with C.I. 95%.  
For solving this question you allowed to assume anything.

And There are various ways of solving this question:  
By using  
①  $\downarrow$  Central limit Theorem

② By using Confidence interval with Population S.D and without Population S.D

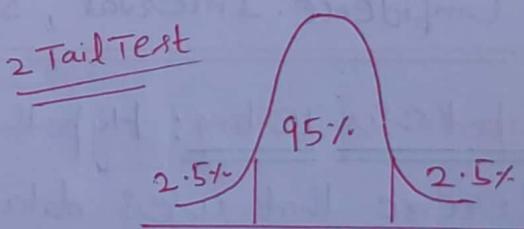
Now I am going to solve this question with the help of Confidence interval where Population S.D. is given.

Let consider, Population Standard Deviation is:

$$\sigma = 100$$

$$\text{Sample}(n) = 30$$

$$\text{Sample mean } (\bar{x}) = 500$$



$$\text{C.I} = \text{Point Estimate} \pm \text{Margin Error}$$

$$= \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \leftarrow \begin{matrix} \text{This is the formula when} \\ \text{your population S.D. is given} \end{matrix}$$

$$= 500 \pm Z_{0.05/2} \frac{100}{\sqrt{n}}$$

$$\begin{aligned} Z_{0.025} &= 1 - 0.025 \\ &= 0.976 \\ &\downarrow \\ &1.96 \end{aligned}$$

$$\underline{\text{Lower Limit}} = 500 - 1.96 \times \frac{100}{\sqrt{30}} = \underline{\underline{386}}$$

$$\underline{\text{Higher Limit}} = 500 + 1.96 \times \frac{100}{\sqrt{30}} = \underline{\underline{613}}$$

Conclusion, If we have 95% C.I., my population mean will be falling between 386 to 613.

\*

### Most Important

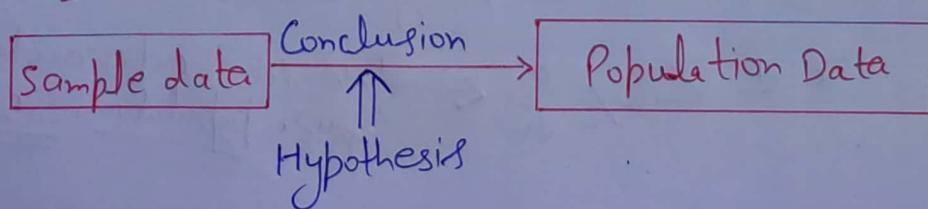
Note : ① Z-test is usually applied when the Population Standard deviation is given. or  $n \geq 30$ .

② When the Population Standard deviation is not given then, we apply t-test. and  $n < 30$  and Sample S.D. is given.

### Interview

③ Hypothesis Testing And Relationship between Hypothesis testing, Confidence Interval , Significance Value , - - - -

Hypothesis Testing: Hypothesis testing is a form of statistical inference that uses data from a sample to draw conclusions about a population parameter or a population probability distribution.



⇒ Inside Hypothesis Testing there are different different experiments (Testings) are available like: Z test, t test, Anova test, Chi square test etc. (18)

Relationship between Hypothesis Testing, C.I. and Significance value:

Hypothesis Testing :-

Coin → Test whether this coin is a fair coin or not by performing 100 tosses.

Sholay coin  $P(H) = 100\%$

$$P(H) = 0.5 \quad P(T) = 0.5$$

100

50 times Head (The coin is fair)

Hypothesis Testing steps:

- ① Null Hypothesis ( $H_0$ ) = Coin is fair.
- ② Alternate Hypothesis ( $H_1$ ) = Coin is not fair or Coin is unfair.
- ③ Experiment :- Toss a coin for 100 times.
- ④ Reject or Accept the Null Hypothesis.

Let's Consider,

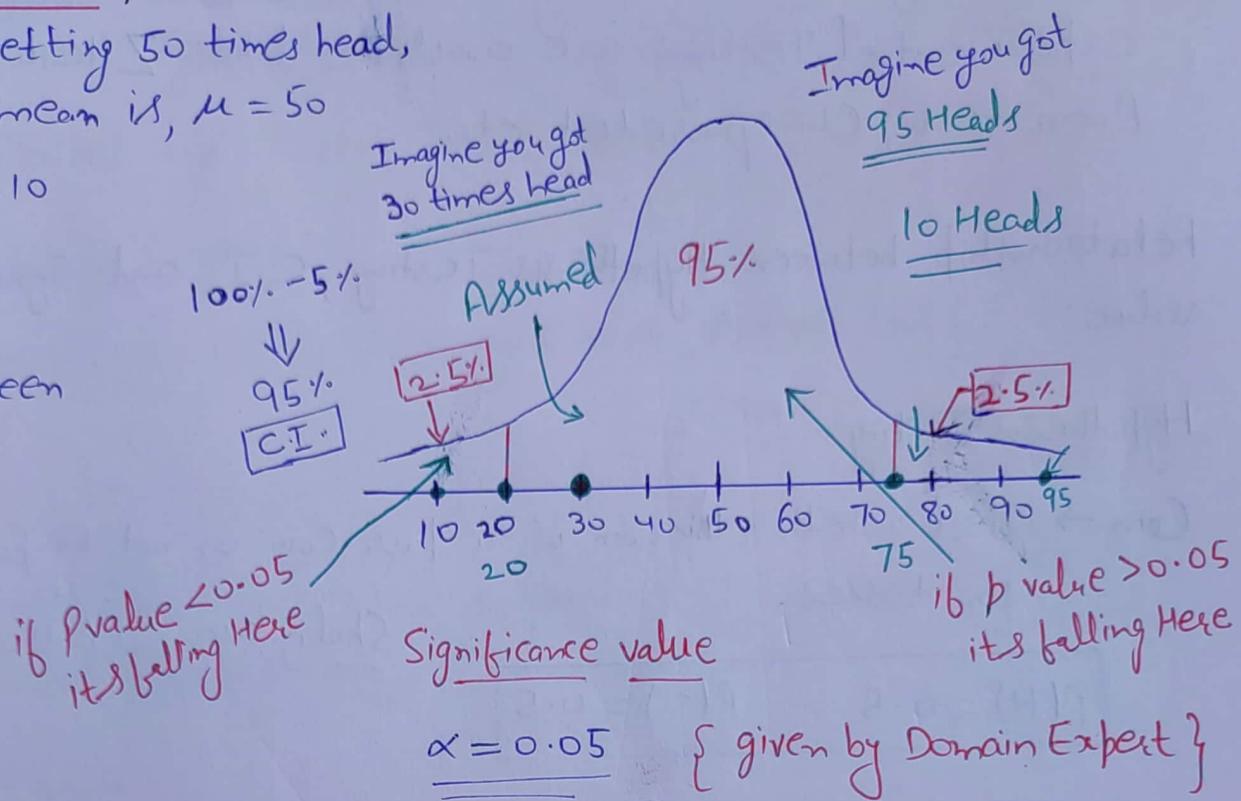
You are getting 50 times head,

Assume, mean is,  $\mu = 50$

S.D.  $\sigma = 10$

Assume,

C.I. between  
20 to 75



Note:

Significance value  $\neq$  P value

- \* P value means you are just saying probability
- \* Significance says what should be within your C.I.

## Today Topics

- ① Type 1 and Type 2 Error
- ② One(1) Tailed and Two(2) Tailed Test
- ③ Confidence Interval
- ④ Z-test, t-test, Chi-square test, ANOVA Test (F test)

### ① Type 1 and Type 2 Error

Null Hypothesis ( $H_0$ ) = Coin is fair

Alternate Hypothesis ( $H_1$ ) = Coin is not fair

#### Reality check

Null Hypothesis is True or Null Hypothesis is False

#### Decision

Null Hypothesis is True or Null Hypothesis is False.

#### Outcome 1:

We reject the Null Hypothesis, when in reality it is false

↓  
Yes

#### Outcome 2:

We reject the Null Hypothesis, when in reality it is true.

e.g.: Movies

Type I  
Error

Person - Death Sentence

### Outcome 3:

We Accept the Null Hypothesis when in reality it is false. Type 2 Error

### Outcome 4:

We Accept the Null Hypothesis when in reality it is true.

Good

P	N
T	TP      TN
F	FP      FN

↓ Type 1

→ Type 2

### ② 1 Tail and 2 tail Test

Example: Colleges in Karnataka have an 85% placement rate. A new college was recently opened and it was found that a sample of 150 students had a placement rate of 88%. With a standard deviation 4%. Does this college has a different placement rate?

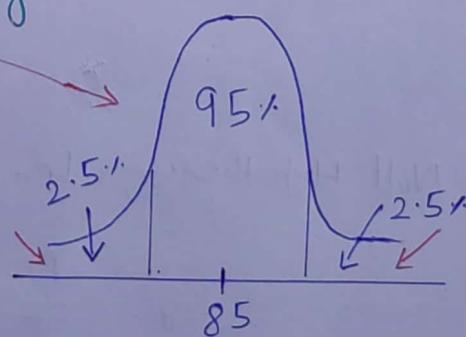
$$\alpha = 0.05 \text{ (Assume)}$$

↳ 85%

85 is our mean.

2 tailed Test

because here two possibilities occur  
placement rate may be greater than 85%  
or may be less than 85%.

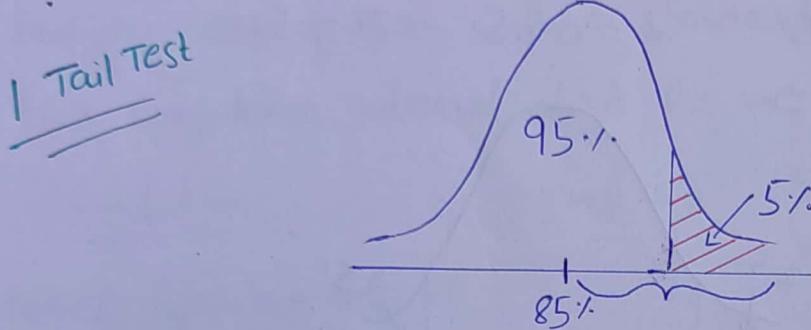


Just change the question - Does this college have a placement rate greater than 85%? (20)

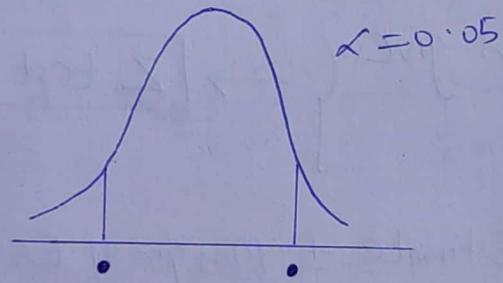
$$\alpha = 0.05 \text{ (Assume)}$$

One tailed test

because it have only one condition or say possibility which is greater.



### ③ Confidence Interval



Point Estimate: The value of any statistic that estimates the value of a parameter.

Inferential stats



$$\boxed{\bar{x}}$$

$$\mu$$

$$\bar{x} = 2.9$$

$$\mu = 3$$

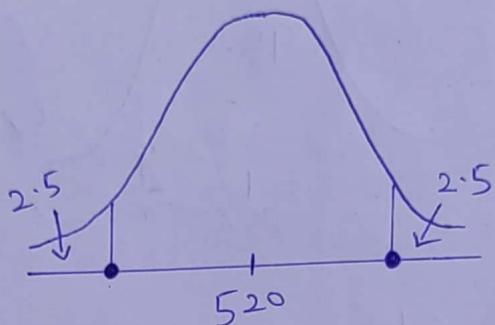
## Confidence Intervals

Point Estimate  $\pm$  Margin of Error

- Q On the quant test of CAT Exam, the standard deviation is known to be 100. A sample of 25 test takers has a mean of 520 score. Construct a 95% C.I about the mean?

Ans  $\sigma = 100$      $n = 25$      $\alpha = 0.05$      $\bar{x} = 520$

$$\alpha' = 1 - 0.95 = 0.05$$



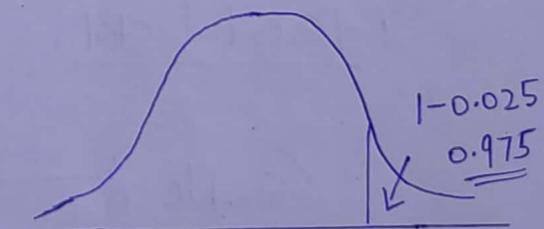
{ ① Population Std is given }  $\rightarrow$  Z test  
 { ②  $n \geq 30$  }

Usually sample is always greater or equal 30, here just for example sample is taken 25.

Point Estimate  $\pm$  Margin of Error

$$\bar{x} \pm Z_{\alpha/2} \left[ \frac{\sigma}{\sqrt{n}} \right] \rightarrow \text{Standard Error}$$

$$Z_{\frac{0.05}{2}} = Z_{0.025}$$

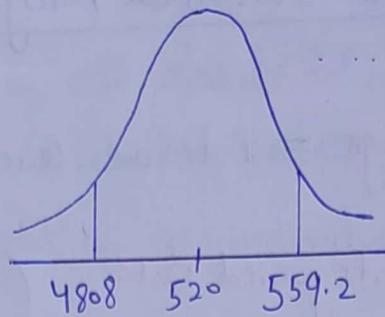


$$\text{Upper bound} = \bar{x} + Z_{\frac{0.05}{2}} \frac{\sigma}{\sqrt{n}}$$

$$\text{Lower bound} = \bar{x} - Z_{\frac{0.05}{2}} \frac{\sigma}{\sqrt{n}}$$

$$\boxed{1.96}$$

$$\begin{aligned}
 \text{Upper} &= 520 + 1.96(20) = 559.2 \\
 \text{Lower} &= 520 - 1.96(20) = 480.8
 \end{aligned}
 \} \quad \underline{\text{Ans}}$$



- Q On the quant test of CAT exam, a sample of 25 test takers has a mean of  $520^\circ$  with a standard deviation of  $80$ . Construct 95% Confidence interval about the mean?

Ans Condition  $n = 25 \quad \bar{x} = 520 \quad S = 80$   
 $\alpha = 0.05$

Here, Population S.D. is  
not given  $\rightarrow t$  - test

Point Estimate  $\pm$  Margin of Error

$$\bar{x} \pm t_{\alpha/2} \left( \frac{S}{\sqrt{n}} \right) \rightarrow \text{Standard Error}$$

$$\text{Upper bound} = \bar{x} + t_{0.05/2} \left( \frac{S}{\sqrt{n}} \right) \quad t_{0.05/2} = \underline{\underline{2.064}}$$

$$\underline{\text{Degree of freedom}} = n-1 = 25-1 = 24$$

$$= 520 + 2.064 \left( \frac{80}{\sqrt{24}} \right)$$

$$= 553.024$$

$$\underline{\text{Lower bound}} = \bar{x} - t_{0.05/2} \left( \frac{S}{\sqrt{n}} \right)$$

$$= 520 - 2.064 \left( \frac{80}{\sqrt{24}} \right)$$

$$= 486.97$$

$$[486.97 \longleftrightarrow 553.024]$$

# Hypothesis Testing And Statistical Analysis

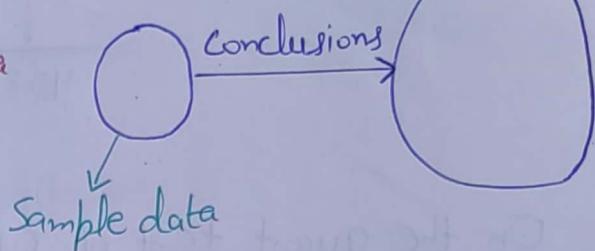
① Z test }  $\Rightarrow$  Average  $\Rightarrow$  Interview Question  
 ② t test }

③ CHI SQUARE  $\Rightarrow$  Categorical Data

④ ANOVA (F test)  $\Rightarrow$  Variance

## Inferential Stats

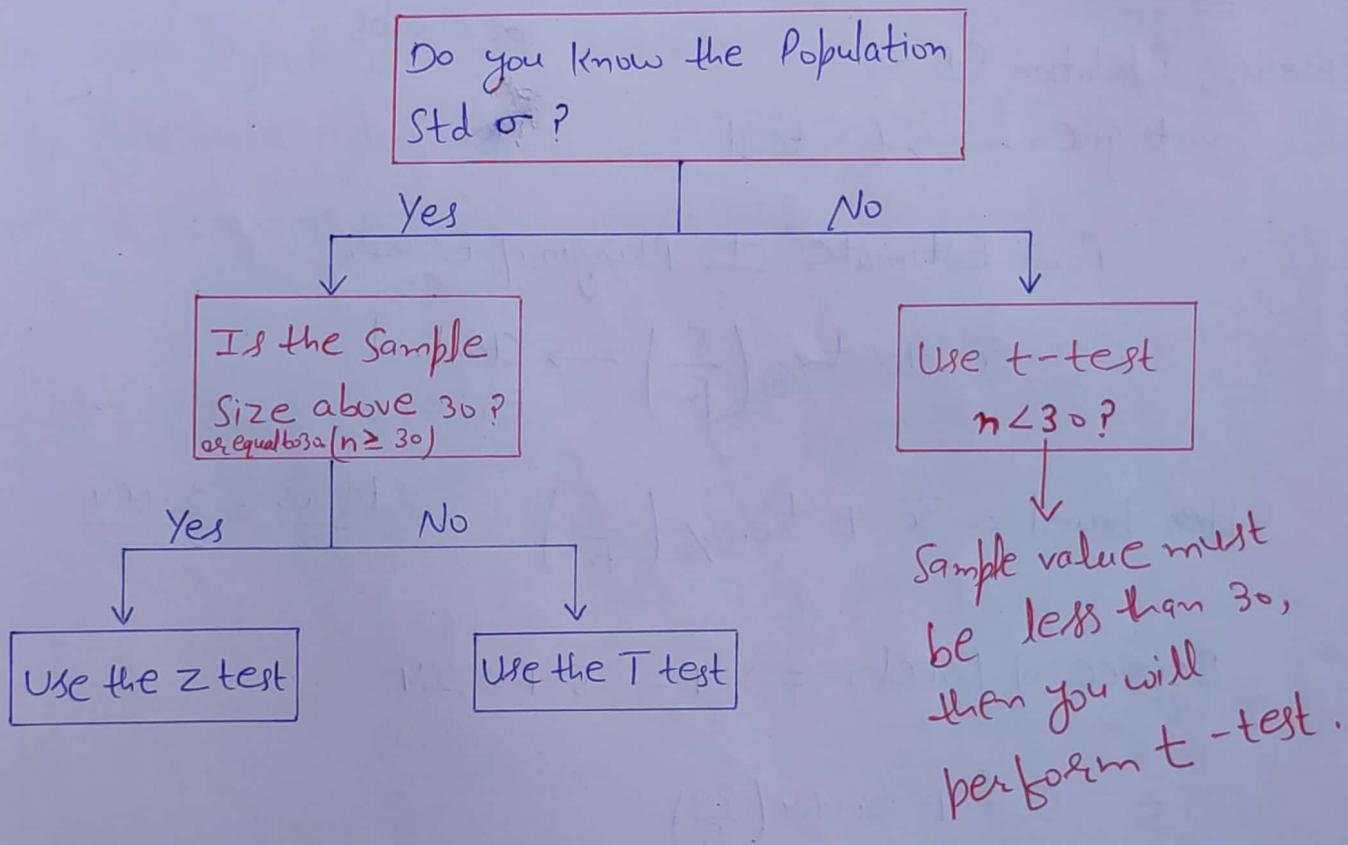
Population  
data  
Conclusion



Z test  $\Rightarrow$  Z table

## Interview Question

When to use T-test Vs Z-test



$n \geq 30$  or Population  $\sigma$  }  $\Rightarrow$  Z-test

$n < 30$  and sample  $\sigma$  }  $\Rightarrow$  T-test

if  $n \geq 30$   
 then perform Z-test

## Z-test

(2)

① The average heights of all residents in a city is 168 cm. with a population std  $\sigma = 3.9$ . A doctor believe the mean to be different. He measured the height of 36 individuals and found the average to 169.5 cm.

(a) State Null And Alternate Hypothesis

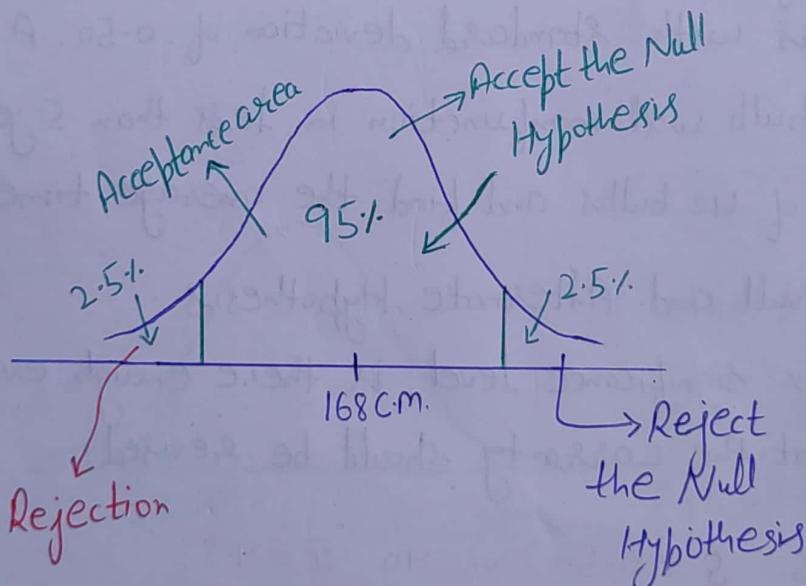
(b) At 95% C.I., is there enough evidence to Reject the Null Hypothesis.

Ans  $\mu = 168 \text{ cm.}$   $\sigma = 3.9$   $n = 36$   $\bar{x} = 169.5 \text{ cm.}$

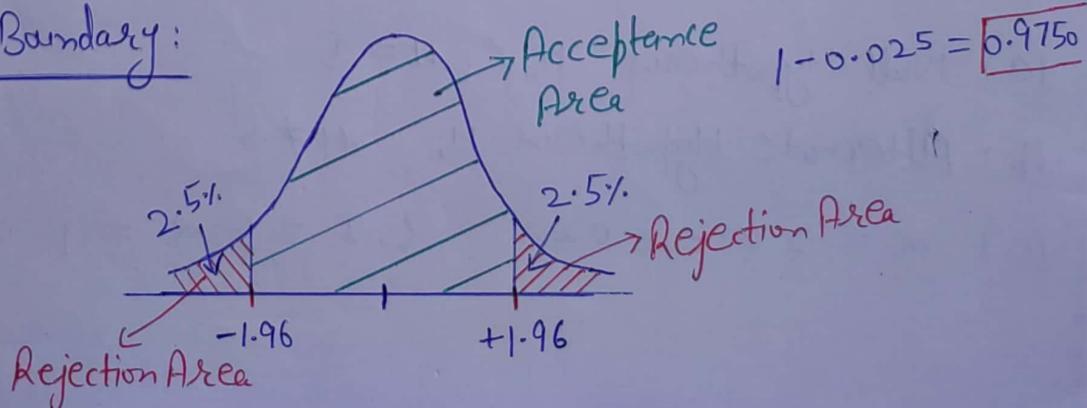
(a) Null Hypothesis  $H_0: \mu = 168 \text{ cm.}$

(b) Alternate Hypothesis  $H_1: \mu \neq 168 \text{ cm.}$  {2 Tail Test}

(c)  $C.I. = 0.95 \Rightarrow 95\% \quad \alpha = 1 - C.I. = 1 - 0.95 = 0.05$



(d) Decision Boundary:



$$1 - 0.025 = 0.9750$$

If Z-test value falls between -1.96 to +1.96 then we fail to reject the Null Hypothesis.

(e) Statistical Analysis

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

|  $n=1$

$$Z_{\text{test}} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} \quad \boxed{\sigma / \sqrt{n}} \Rightarrow \text{Standard Error}$$

$$= \frac{169.5 - 168}{3.9 / \sqrt{36}} = \boxed{2.31}$$

Conclusion: If Z-test is less than -1.96 or greater than +1.96 then Reject the Null Hypothesis.

Here,  $2.31 > 1.96$  {we Reject the Null Hypothesis?}

② A factory manufactures bulbs with a average warranty of 5 years with standard deviation of 0.50. A worker believes that the bulb will malfunction in less than 5 years. He tests a sample of 40 bulbs and find the average time to be 4.8 years.

(a) State Null and Alternate Hypothesis

(b) At a 2% significance level, is there enough evidence to support the idea that the warranty should be revised.

| Tail Test

Ans  $\mu = 5 \quad s = 0.50 \quad n = 40 \quad \bar{x} = 4.8$

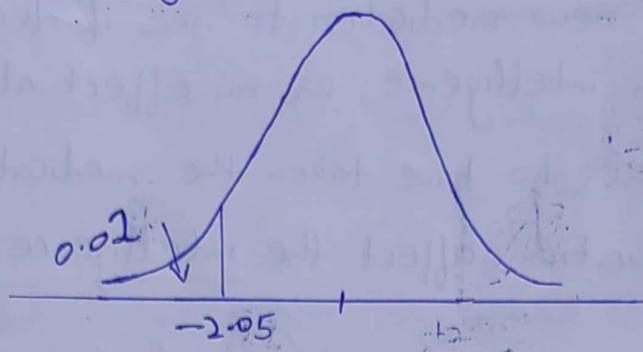
(a) Null Hypothesis  $H_0 \quad \mu = 5$

(b) Alternate Hypothesis  $H_1 \quad \mu \neq 5$

(c)  $\alpha = 2\% \Rightarrow 0.02 \quad C.I. = 1 - \alpha = 1 - 0.02 = 0.98$

(d) Decision Boundary:

(23)

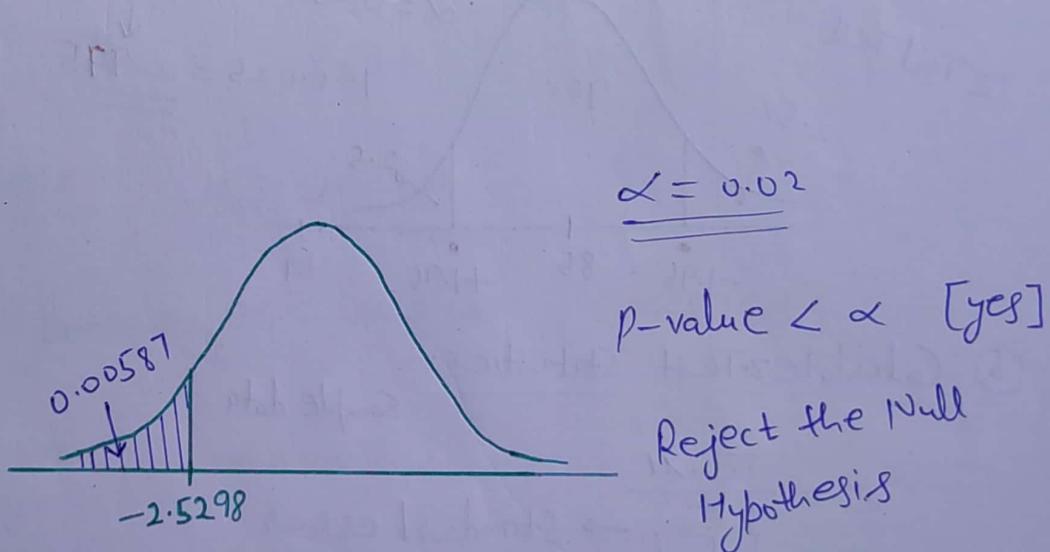


(e) Calculate test statistics (Z-Test)

$$Z = \frac{\bar{x} - \mu}{S/\sqrt{n}} = \frac{4.8 - 5}{0.50/\sqrt{50}} = -2.5298$$

Conclusion: Here,  $-2.52 < -2.05$ . So we reject the Null Hypothesis. Warranty needs to be revised.

P value



Q In the Population, the average IQ is 100 with a SD of 15. Researchers wants to test a new medication to see if there is positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect the intelligence.

$$\underline{\alpha = 0.05} \quad \underline{C.I. = 95\%}$$

Ans ① Define Null Hypothesis

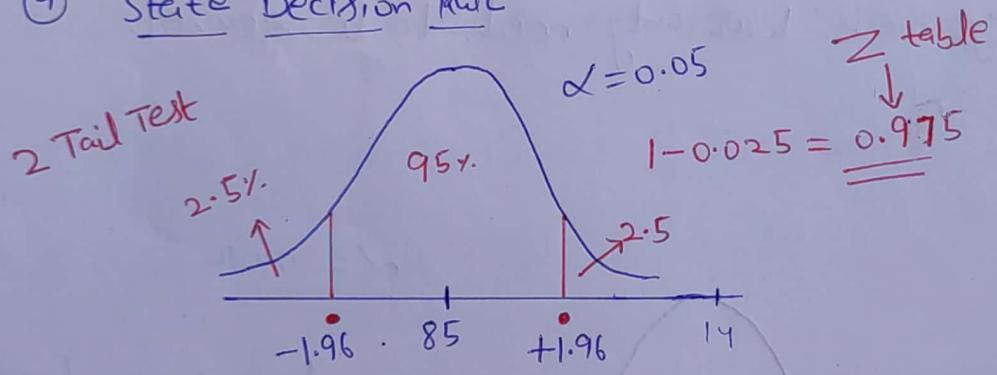
$$H_0 = \mu = 100$$

② Alternate Hypothesis  $H_1, \mu \neq 100$

③ State Alpha

$$\alpha = 0.05$$

④ State Decision Rule



⑤ Calculate Z-test Statistics:

Sample data

$$Z = \frac{\bar{x} - \mu}{\left[ \frac{\sigma}{\sqrt{n}} \right]} \rightarrow \underline{\text{standard error}}$$

$$= \frac{140 - 100}{\frac{15}{\sqrt{30}}} = \frac{40}{15} \times \sqrt{30} = \underline{\underline{14.30}}$$

## ⑥ State our Decision

$$14.60 > 1.96$$

$$\left\{ \begin{array}{l} z = -0.2 \\ z = 14.60 \end{array} \right.$$

If  $z$  is less than  $-1.96$  or greater than  $1.96$ , reject the null hypothesis.

Medication improve the intelligence or decrease?

Improve the  
intelligence.

In the population the average IQ is 100. A team of researchers want to test a new medication to see if it has either a positive or negative effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a standard deviation of 20. Did the medication affect intelligence? C.I. = 95%

Ans  $\mu = 100$ ,  $n = 30$   $S = 20$  C.I. = 95%  $\alpha = 1 - 0.95 = 0.05$

① Null Hypothesis  $H_0: \mu = 100$

Alternate Hypothesis  $H_1: \mu \neq 100$  {2 Tail Test}

② Significance value  $\alpha = 0.05$

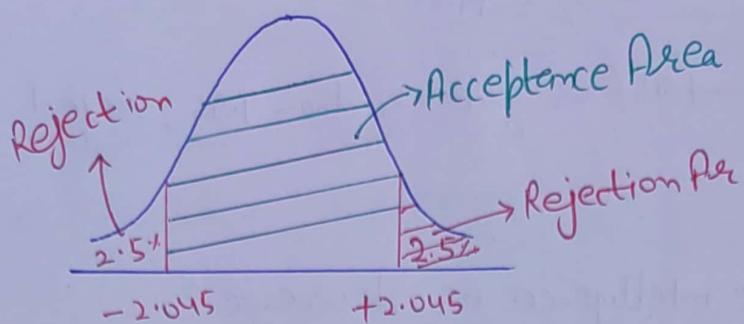
③ Degree of freedom =  $n - 1$

$$= 30 - 1$$

$$= 29$$

#### ④ Decision Boundary

dof = 29



If t test is less than -2.045 or greater than 2.045, then we reject the Null Hypothesis.

#### ⑤ Calculate the t test statistics

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = \frac{140 - 100}{\frac{20}{\sqrt{30}}} = \frac{40}{3.65} = 10.96$$

#### Conclusion:

Decision Rule: If t is less than -2.0452 or greater than 2.0452, Reject the Null Hypothesis.

$$t = 10.96 > 2.0452 \Rightarrow \text{Reject the Null Hypothesis}$$

Q A factory has a machine that fills 80ml of Baby medicines in a bottle. An employee believes the average amount of baby medicine is not 80ml. Using 40 samples, he measures the average amount of baby medicine dispersed by the machine to be 78ml with a standard deviation of 2.5.

(a) State Null and Alternate Hypothesis

(b) At 95% C.I., is there enough evidence to support Machine is working properly or not.

Ans Step I

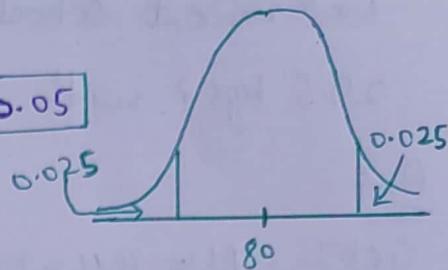
$$\mu = 80 \text{ ml}, n = 40, \bar{x} = 78, S = 2.5$$

(25)

Null Hypothesis  $\mu = 80$   
 Alternative Hypothesis  $\mu \neq 80$

Step 2 C.I. = 0.95

$$S.V(\alpha) = 1 - 0.95 = 0.05$$



Step 3

$$n = 40$$

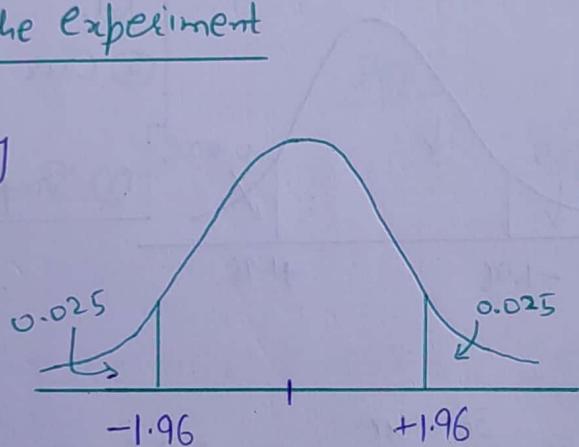
$$S = 2.5$$

Z test

Let's perform the experiment

Decision Boundary

$$1 - 0.025 = 0.975$$



Step 4 Calculate test statistics (Z-test)

$$Z = \frac{\bar{x} - \mu}{\frac{S}{\sqrt{n}}} = \frac{78 - 80}{\frac{2.5}{\sqrt{40}}} = -5.05$$

→ Standard error

Conclusions:

Decision Rule: If  $Z = -5.05$  is less than  $-1.96$  or greater  $+1.96$ ,  
 Reject the Null Hypothesis with 95% C.I.

Reject the Null Hypothesis { There is some fault in the  
 machine.

Q A Complain was registered, the boys in a Government school are underfed. Average weight of boys of age 10 is 32 kgs. with S.D. = 9 kgs. A sample of 25 boys were selected from the Government school and the average weight was found to be 29.5 kgs? with C.I. = 95%. Check it is True or False.

Ans

Step I  $H_0: \mu = 32$

$H_1: \mu \neq 32$

Conditions for Z-test

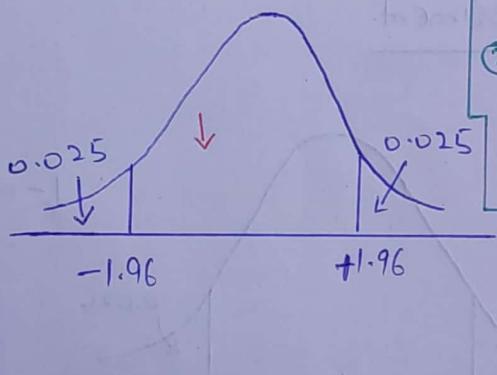
① We know the Population Sd. **OR**

② We do not know the Population Sd but our Sample is large  $\geq 30$

1 Tail Test

② C.I. = 95%  $\alpha = 1 - 0.95 = 0.05$

③ Z-test



Conditions for T-test

① We do not know the Population Sd.

② Our Sample size is small  $n < 30$

③ Sample Sd is given.

$$Z\text{-score} = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{29.5 - 32}{9 / \sqrt{25}} = -1.39$$

Conclusion:  $-1.39 > -1.96$ . So, we Accept the Null Hypothesis

95% C.I. We fail to Reject Null Hypothesis

The boys are fed well.

Q The average weight of all residents in a town XYZ is 168 Pounds. A nutritionist believes the true mean to be different. She measured the weight of 36 individuals and found the mean to be 169.5 pounds with a standard deviation of 3.9.

(a) Null { Alternate Hypothesis }

(b) 95%. Is there enough evidence to discard the null Hypothesis?

Ans  $\bar{x} = 169.5$      $s = 3.9$      $n = 36$      $\mu = 168$     C.I. = 0.95

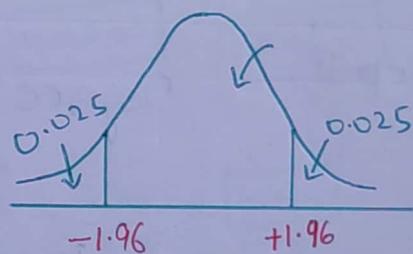
Step 1

$$H_0: \mu = 168$$

$$H_1: \mu \neq 168$$

Step 2: C.I. = 0.95     $\alpha = 0.05$

Step 3:

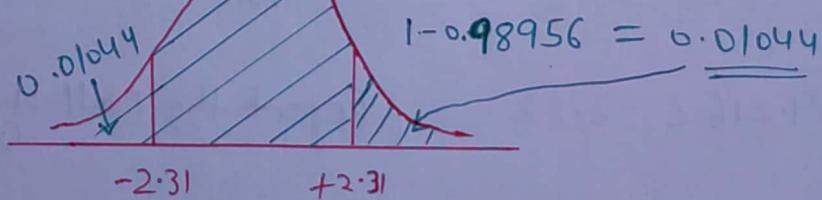


Step 4:  $Z\text{-score} = \frac{169.5 - 168}{3.9/\sqrt{36}} = \boxed{2.31} \Rightarrow 2.3076$

$2.31 > 1.96$  { Reject the Null Hypothesis }

P-value

2 Tail Test



$$P\text{-value} = 0.01044 + 0.01044 = 0.02088$$

$0.02088 < 0.05$

{ Reject the Null Hypothesis }

Q A company manufactures bike batteries with an average life span of 2 years or more years. An engineer believes this value to be less. Using 10 samples, he measures the average life span to be 1.8 years with a standard deviation of 0.15.

(a) State the Null and Alternate Hypothesis

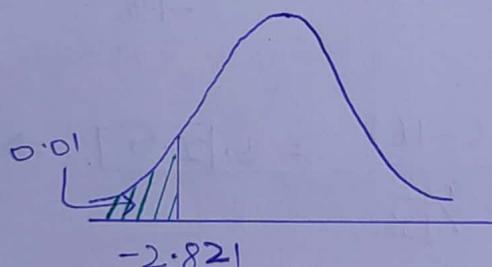
(b) At a 99% C.I., is there enough evidence to discard the  $H_0$ ?

Ans  $\mu = 2$   $n = 10$   $\bar{x} = 1.8$   $s = 0.15$   $C.I. = 0.99$   $\alpha = 0.01$

Step I  $H_0: \mu \geq 2$   
 $H_1: \mu < 2$

1 Tail Test

Step II Degree of freedom =  $n - 1$   
 $= 10 - 1 = 9$



Step III Calculate t test statistics

$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{1.8 - 2}{0.15/\sqrt{10}} = -4.216$$

Step IV  $-4.216 < -2.82$  {Reject the Null Hypothesis}

Conclusion: The average life of the battery is less than 2 years.

## Z test with Proportions

A tech company believes that the percentage of residents in town XYZ that owns a cell phone is 70%. A marketing manager believes that this value to be different. He conducts a survey of 200 individuals and found that 130 responded Yes owning a cell phone?

- (a) State Null and Alternate Hypothesis      2 Tail Test  
 (b) At a 95% C.I., is there enough evidence to reject the Null Hypothesis?

Ans Step 1

$$\text{Null Hypothesis: } P_0 = 0.70 \quad \left. \right\}$$

$$\text{Alternate Hypothesis: } P_0 \neq 0.70 \quad \left. \right\}$$

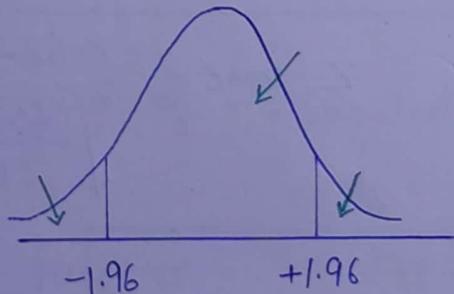
$$q_0 = 1 - P_0 = 0.30$$

$$n = 200 \quad x = 130$$

$$\hat{P} = \frac{130}{200} = 0.65$$

Step 2    C.I = 0.95     $\alpha = 0.05$

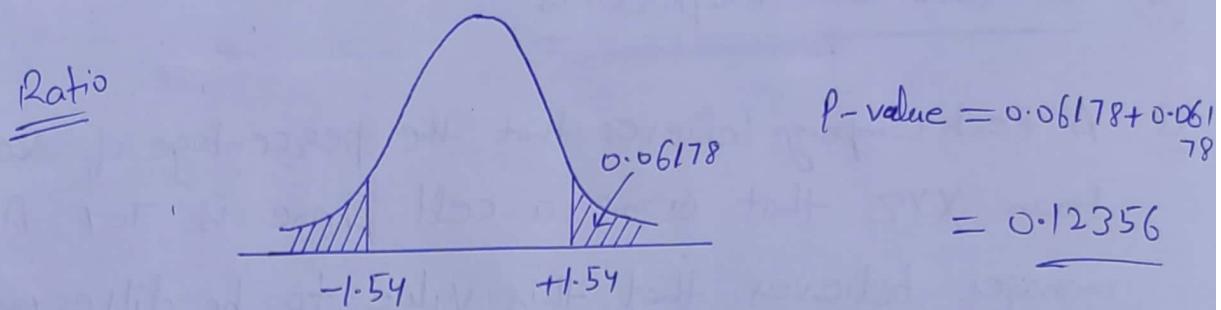
Step 3:



$$Z \text{ test} = \frac{\hat{P} - P_0}{\sqrt{\frac{P_0 q_0}{n}}}$$

$$= \frac{0.65 - 0.70}{\sqrt{\frac{0.7 \times 0.3}{200}}} \approx -1.54$$

Conclusion:  $-1.54 > -1.96$  Fail to Reject the Null Hypothesis.



$P\text{-value} > \text{Significance value} - \text{Fail to Reject the Null Hypothesis.}$

- Q A car company believes that the percentage of residents in city ABC that owns a vehicle is 60% or less. A sales manager disagrees with this. He conducts a hypothesis testing surveying 250 residents and found that 170 responded yes to owning a vehicle.

(a) State the Null and Alternate Hypothesis

(b) At a 10% significance level, is there enough evidence to support the idea that vehicle ownership in city ABC is 60% or less?

$$\underline{\text{Ans}} \quad H_0: P_0 \leq 0.60$$

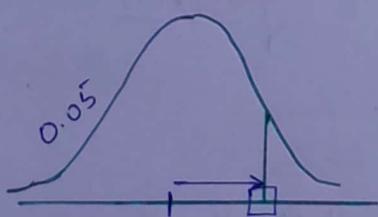
$$H_1: P_0 > 0.60$$

$$\hat{P} = \frac{170}{250} = 0.68$$

$$q_0 = 1 - P_0 = 0.40$$

$$Z\text{-score} = \frac{0.68 - 0.60}{\sqrt{\frac{0.6 \times 0.4}{250}}}$$

$$= \frac{0.08}{0.0309} = 2.588$$



Reject the Null Hypothesis

## Agenda

(28)

- ① CHI SQUARE
- ② Covariance
- ③ Pearson Correlation Coefficient
- ④ Spearman Rank Correlation
- ⑤ Practical Implementation
- ⑥ Z-test, t-test, Chi-square test
- ⑦ F test (ANOVA)

Why CHI SQUARE TEST Performed?

Interview

- ① CHI Square Test claims about population proportions.  
It is a non parametric test that is performed on categorical (nominal or ordinal) data.
- ② In the 2000 Indian Census, the age of the individual in a small town were found to be the following:

Less than 18	18-35	>35
20%	30%	50%

In 2010, age of  $n=500$  individuals were sampled. Below are results

$\leq 18$	18-35	$>35$
121	288	91

Using  $\alpha = 0.05$ , would you conclude the population distribution of ages has changed in the last 10 years?

Ans

$< 18$	$18 - 35$	$> 35$
20%	30%	50%

{Population}  $\underline{\underline{2000}}$

Expected

$< 18$	$18 - 35$	$> 35$
121	288	91
$500 \times 0.2$	$500 \times 0.3$	$500 \times 0.5$

[n = 500]

Observed

Expected

100                  150                  250

(Chi Square Table)

$< 18$	$18 - 35$	$> 35$
121	288	91
100	150	250

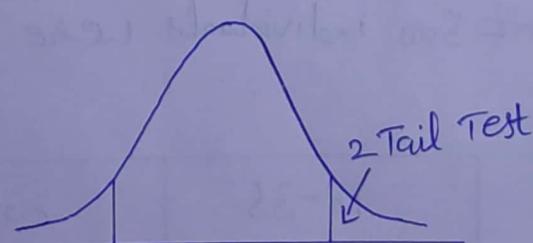
①  $H_0$  = The data meets the distribution 2000 census

$H_1$  = The data does not meet

②  $\alpha = 0.05$  (95% C.I.)  $df = 2$ ,  $\alpha = 0.05$

③ Degree of freedom =  $n-1 = 3-1 = 2$

④ Decision Boundary



If  $\chi^2$  is greater than 5.99, then Reject the Null Hypothesis  $H_0$ .

## ⑤ Calculate Test Statistics

(29)

$$\chi^2 = \sum \frac{(f_o - f_e)^2}{f_e}$$

$$= \frac{(121-100)^2}{100} + \frac{288-150}{150} + \frac{(91-250)^2}{250}$$

$$= 232.494$$

$$\underline{\underline{\chi^2 = 232.494 > 5.99}}$$

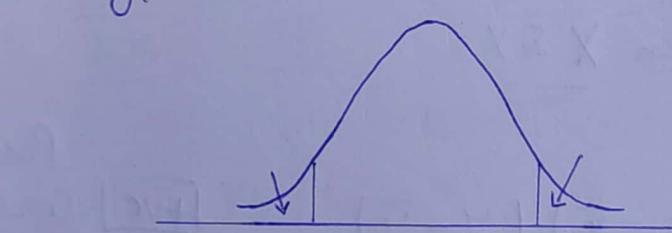
{ Reject the Null Hypothesis }

$$\underline{\underline{\alpha = 0.05}}$$

$$\begin{array}{c} 0.01 \\ \hline 0.10 \end{array} \quad \text{Domain}$$

$0.11 > 0.05$   
 $\downarrow$   
 Accept the Null Hypothesis

$$\underline{\underline{0.002 < 0.05}}$$



{ P-value < Significance value }  
 $\downarrow$   
 Reject the Null Hypothesis

OR

{ Accept the Null Hypothesis }



$\underline{\underline{P = 0.11 > 0.05}}$   
 $\downarrow$   
 Accept

$$P = \underline{\underline{0.002 < 0.05}}$$

$\downarrow$   
 Reject the Null Hypothesis

## ② Covariance:

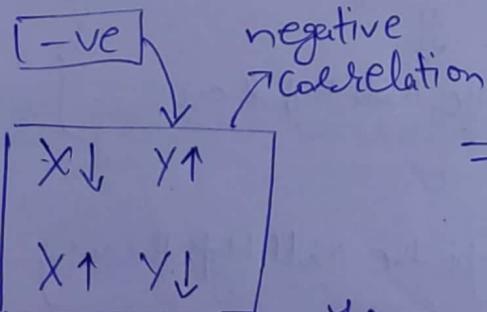
$\boxed{X}$	$\boxed{Y}$	
Weight	Height	
50	160	$\left\{ \begin{array}{ c c } \hline X \uparrow & Y \uparrow \\ \hline X \downarrow & Y \downarrow \\ \hline \end{array} \right\}$
60	170	
70	180	
75	181	

No. of hours Study	Play
2	6
3	4
4	3

Quantity relationship between  $\underline{X \text{ & } Y}$

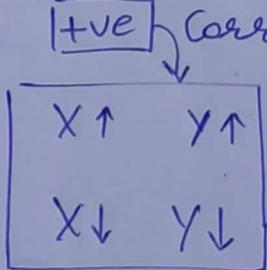
Covariance

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

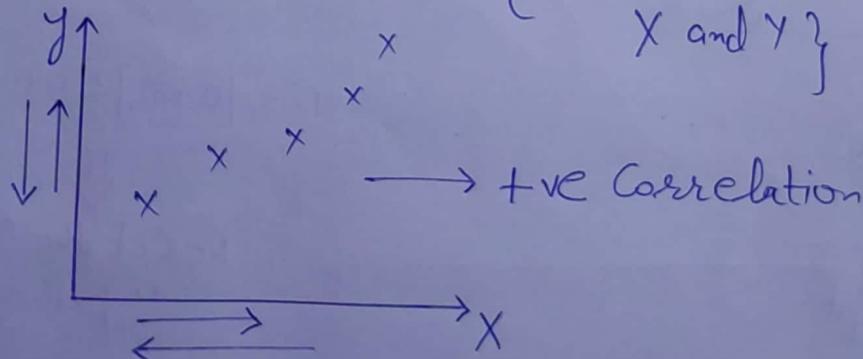


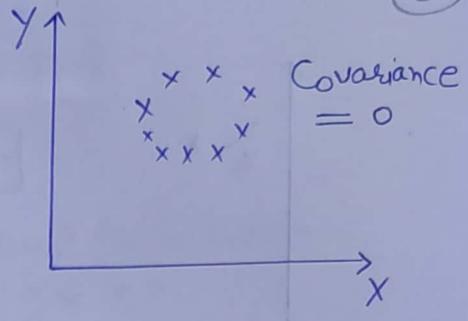
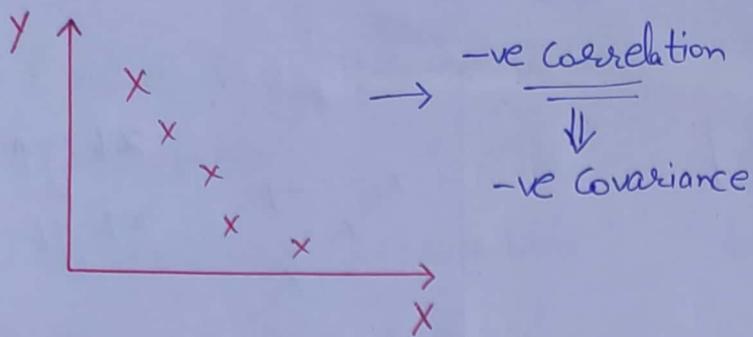
= +ve or -ve

Positive



o { There is no relation b/w  
 $X$  and  $Y$  }





Disadvantage of Covariance

① Positive or Negative ✓

$$\begin{array}{c} +100 \\ \hline -2000 \\ -200 \end{array} \quad \begin{array}{c} +1000 \\ \hline \end{array} \quad \{ \text{Direction} \}$$

$$\begin{array}{cc}
 X & Y \\
 10 & 4 \\
 8 & 6 \\
 7 & 8 \\
 \hline
 \frac{6}{7.75} & \frac{10}{7}
 \end{array}
 \quad \begin{aligned}
 \text{Cov}(x,y) &= \underline{\underline{-ve}} \\
 &= [(10-7.75)(4-7) + (8-7.75)(6-7) + (7-7.75)(8-7) \\
 &\quad + (6-7.75)(10-7)] \\
 &= \underline{\underline{-3.25}}
 \end{aligned}$$

$$\begin{bmatrix}
 X\uparrow & Y\downarrow \\
 X\downarrow & Y\uparrow
 \end{bmatrix}$$

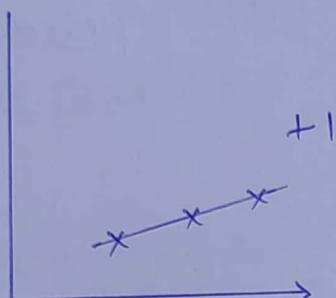
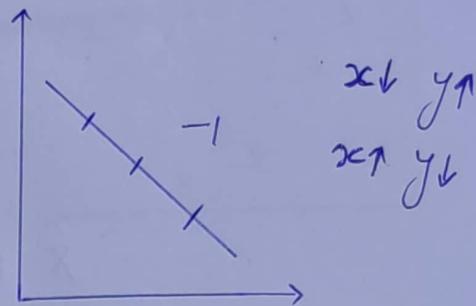
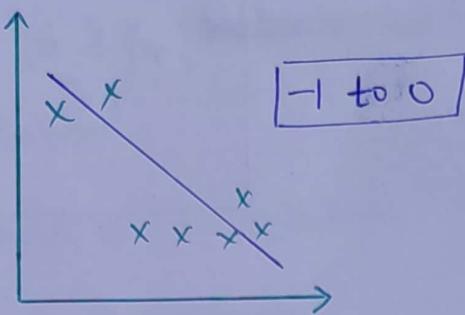
150

③ Pearson Correlation Coefficient:

(-1 to 1) The more towards +1 more positively correlated.

The more towards -1 more negatively correlated.

$$\rho(x,y) = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} \quad \underline{\underline{\{-1 \text{ to } 1\}}}$$



Interview  
why use Pearson Correlation?

for testing whether there is significant relationship b/w two variables

#### ④ Spearman Rank Correlation

$$\text{Spear}(x, y) = \frac{\text{Cov}(R(x), R(y))}{R_{ox} * R_{oy}}$$

x Height	y Weight
170	75
160	62
150	60
145	55
180	85

R(x)	R(y)
2	2
3	3
4	4
5	5
1	1

Interview  
Why we use Spearman Rank Correlation?  
It captures Non Linear Properties

① P value & significance value

② Distribution

③ Bernoulli's Distribution

Log Normal Distribution

④ Binomial Distribution

⑤ Pareto's Distribution {Power law}

⑥ F Test (ANOVA)

① P value & significance value

↳ Derive the P value

Q The average weight of all residents in Bangalore city is 168 pounds with a standard deviation 3.9. We take a sample of 36 individuals and the mean is 169.5 pounds. C.I. = 95%

$$\text{Ans} \quad \mu = 168 \quad \sigma = 3.9 \quad \bar{x} = 169.5 \quad n = 36 \quad \alpha = 0.05$$

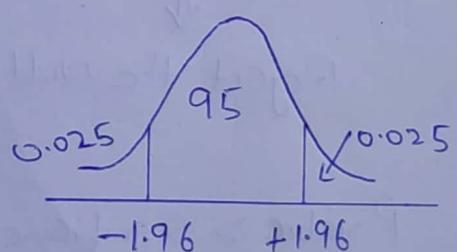
$$① H_0: \mu = 168$$

③ Decision Boundary



$$H_1: \mu \neq 168$$

$$② \alpha = 0.05$$

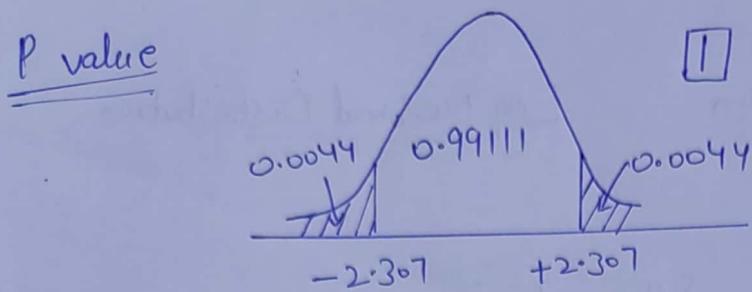


$$1 - 0.025 = 0.9750$$

④ Z test

$$\begin{aligned} Z &= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} \\ &= \frac{1.5}{3.9} \times 6 = \underline{\underline{2.307}} \end{aligned}$$

$$z = 2.307 > 1.96 \quad \begin{cases} \text{Reject the Null} \\ \text{Hypothesis} \end{cases}$$



$$1 - 0.99111 = 0.00889$$

$$\begin{aligned} P \text{ value} &= 0.0044 + 0.0044 \\ &= \underline{\underline{0.0088}} \end{aligned}$$

$$P \text{ value} < 0.05$$

$$0.0088 < 0.05 \rightarrow \begin{cases} \text{Reject the Null Hypothesis} \end{cases}$$

Note:

$P \text{ value} \leq \text{Significance value}$

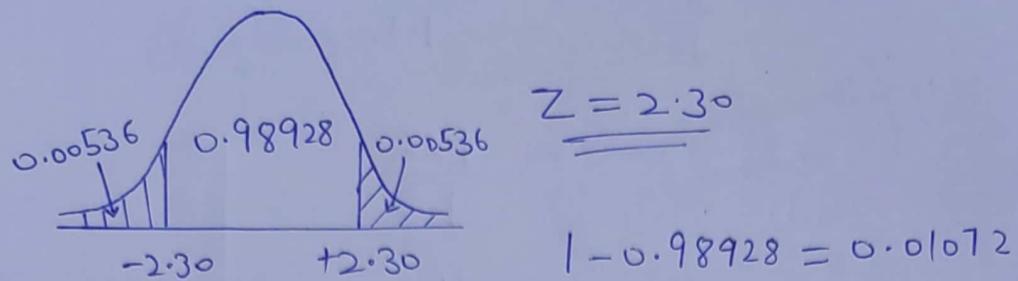


Reject the Null Hypothesis

$P \text{ value} > \text{Significance value}$



Fail to Reject the Null Hypothesis



$2.30 > 1.96 \{ \text{Reject the Null Hypothesis} \}$

$$\begin{aligned} P \text{ value} &= 0.00536 + 0.00536 \\ &= 0.01072 \leq \alpha \Rightarrow \text{Reject the Hypothesis} \end{aligned}$$

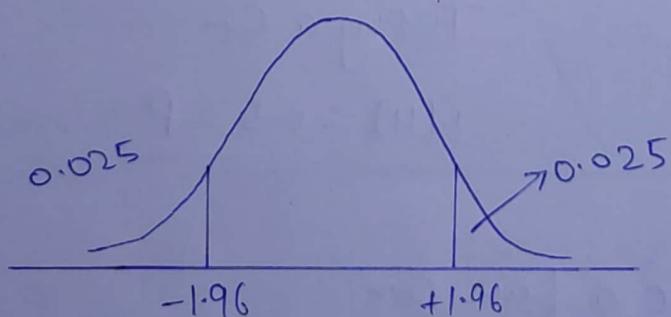
Q Average age of a college is 24 years with a standard deviation 1.5. Sample of 36 students mean is 25 years. with  $\alpha = 0.05$ , C.I = 95%, does the age vary?

Ans ①  $H_0: \mu = 24 \quad \sigma = 1.5 \quad n = 36 \quad \bar{x} = 25 \quad \alpha = 0.05$

$H_1: \mu \neq 24$

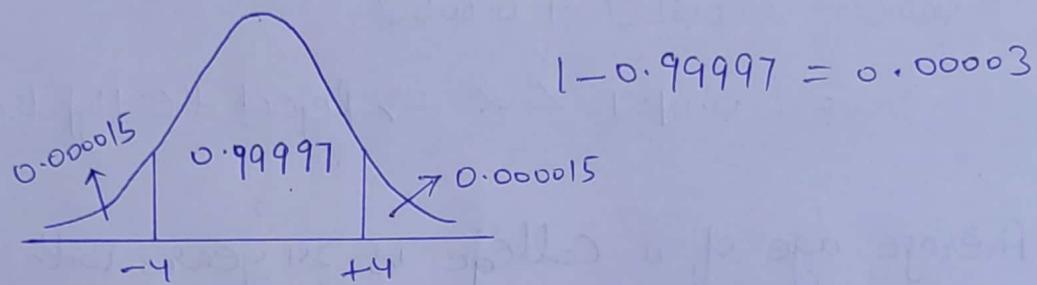
②  $\alpha = 0.05$

③ Decision Boundary



$$\begin{aligned}
 \textcircled{4} \quad Z\text{-score} &= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{25 - 24}{1.5} \times 6 \\
 &= \frac{1 \times 6}{1.5} \\
 &= 4
 \end{aligned}$$

$4 > 1.96$  {Reject Null Hypothesis}



$$\begin{aligned}
 \text{P value} &= 0.000015 + 0.000015 \\
 &= 0.00003
 \end{aligned}$$

$\text{P value} \leq \text{Significance}$  {Reject the Null Hypothesis}

### ③ Bernoulli's Distribution:

2 outcomes      [0 or 1]      Single Trial

$$P = 0.5$$

$$q = 1 - P = 0.5$$

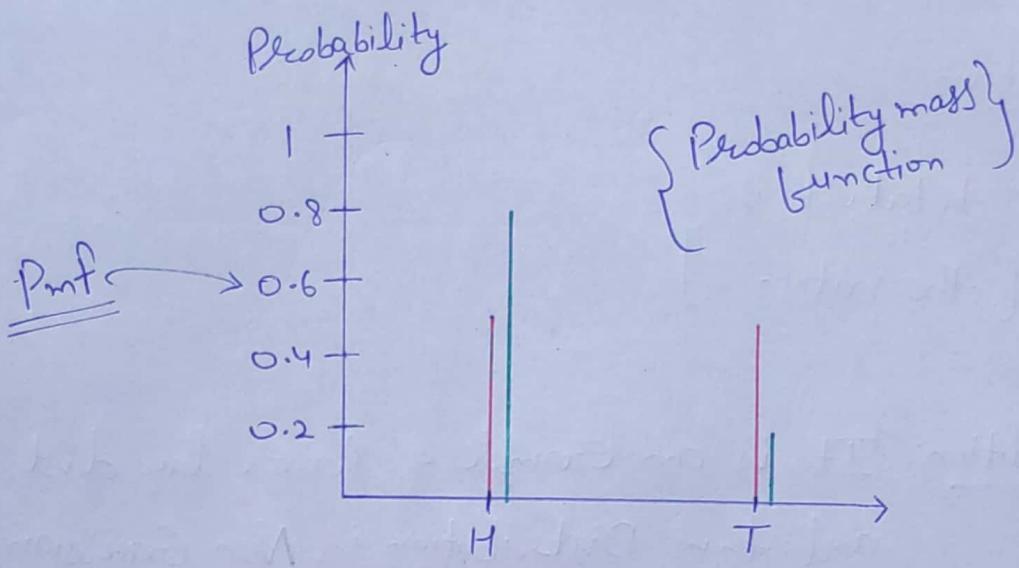
Tossing a Coin

$$\underline{P(H) = 0.5 = P}$$

Do not have a fair coin

$$P(H) = 0.3 = P$$

$$P(T) = q = 1 - P = 1 - 0.3 = \underline{0.7}$$



\* Probability mass function (Pmf): used for categorical variables.

Probability density function (Pdf): used for continuous variables.

#### ④ Binomial Distribution:

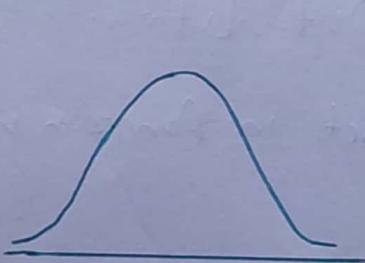
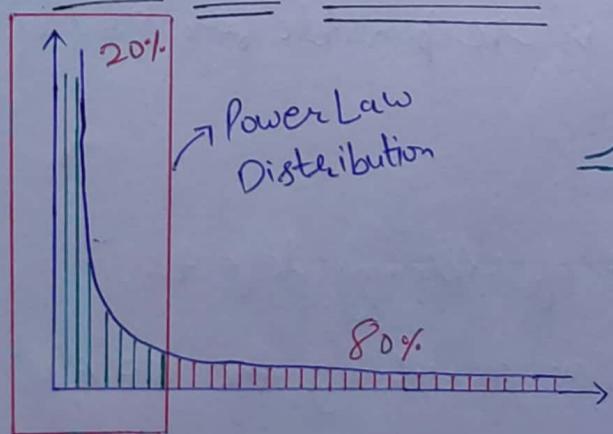
Every Trial  $\rightarrow$  Bernoulli's distribution

Multiple Trial

$$\begin{array}{ll} P(H) = 0.5 & P(H) = 0.6 \\ P(T) = 0.5 & P(T) = 0.4 \end{array} \quad \dots \dots \dots$$

#### Interview

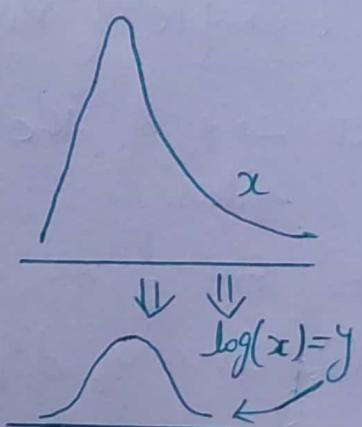
#### ⑤ Power Law Distribution



80-20% Rule

e.g. ① IPL  $\rightarrow$

20% of Team is responsible for winning 80% of match



② 80% of wealth is distributed by 20% of the total Population.

③ Oil Richness

{ 80% of the total oil is }  
with 20% of the nation }

Interview

Pareto Distribution: It is an example of Power law distribution and this Distribution is Non Gaussian.

Interview: Can we Convert a Pareto Distribution into a Normal Distribution? Interview Question

Ans Yes, By using Box Cox transformation we convert a Pareto distribution into Normal Distribution.

Importance of Normal Distribution: If we have any Normal Distributed data, then by using the Machine Learning algorithms like Linear Regression, Logistic Regression we are able to make a Efficient Model.

Example of Pareto Distribution:

- ① 20% of the Product in Amazon is responsible for 80% of sales.
- ② 80-20%
- ③ 20% of defects solves the 80% of upcoming defects.
- ④ 20% of Team is responsible in delivering the 80% of the projects.

Interview: Relationship between Log Normal and Pareto Distribution?