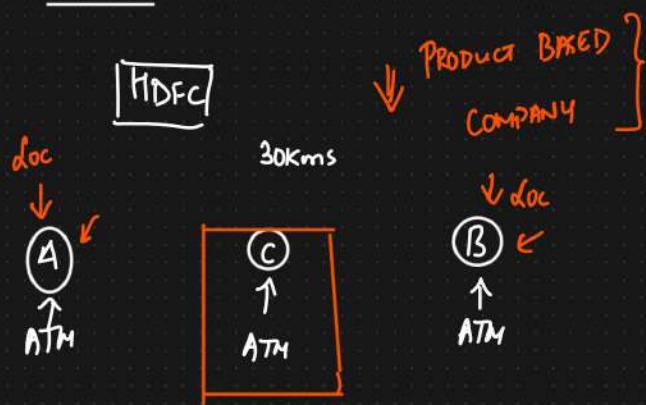


# Statistics

## Usages



X Statistician → 5 years



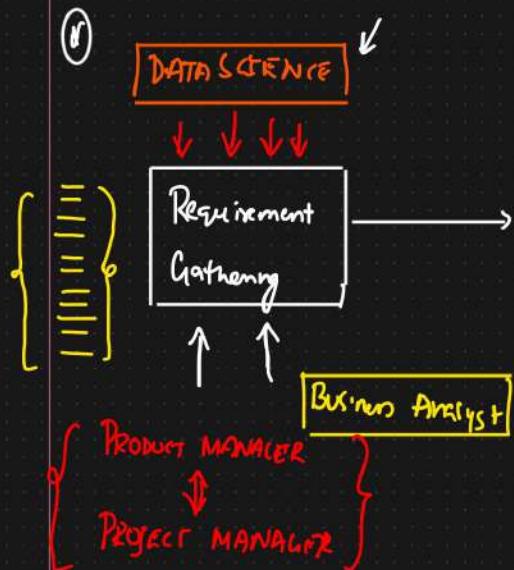
- ① DATA ANALYST
- ② DATA SCIENTIST

blue whale

✓ Amazon

- ② Find the average size of the **shark** throughout the world?
- ③ Amazon Big Billion Day Sale {Intuit} → Which month should you select?

# Statistics {Life Cycle of DATA Science Project}



## DATA ANALYST TEAM

DATA ANALYST /  
DATA SCIENTIST

- ① DATA ANALYST
- ② DATA SCIENTIST
- ③ BIG DATA Engineers
- ④ Cloud Engineers

Domain knowledge

PRODUCT MANAGER  
BA

Apple

PRODUCT BASED

{ Google }

YouTube, GPAY,  
Google Ads, GMAIL

Sales

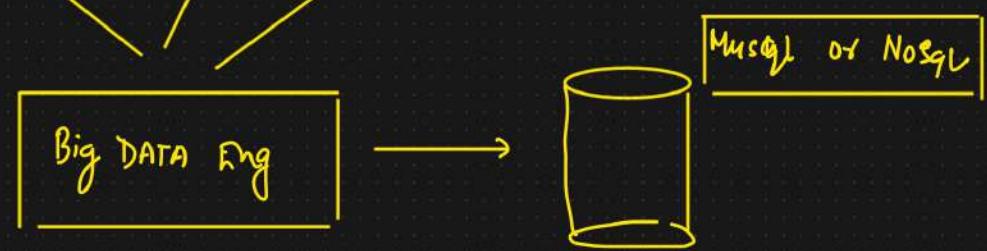
Domain Expertise

Product Manager

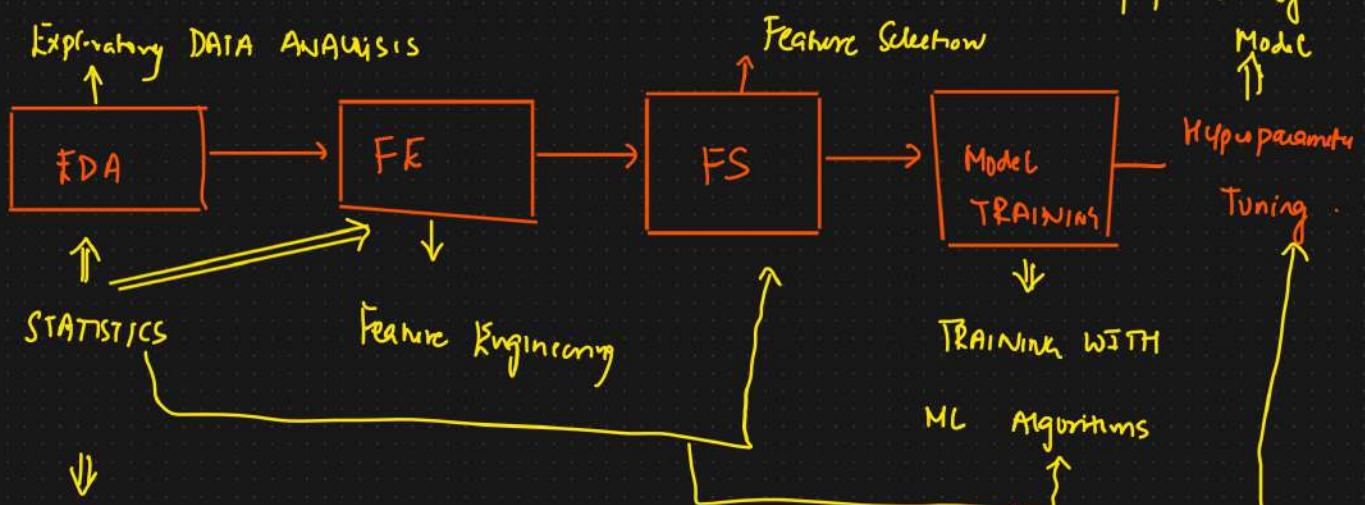
Internal  
DATABASE

3rd party  
API'S

Web Scraping



## Life cycle of DS Project



$$\text{Age} = \{12, 13, 14, 18, 20, 25\} \Rightarrow \text{Average Age} \Rightarrow \text{Measure of Central Tendency}$$

↓

DESCRIPTIVE STATS

Statistics = Defn : Statistics is the science of collecting, organising and analysing the data.

Data : "facts or pieces of information"

Eg: Ages of Students in Classroom

{24, 25, 32, 29, 28}  $\Rightarrow$  Mean, Median, Mode

Standard deviation

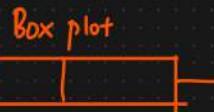
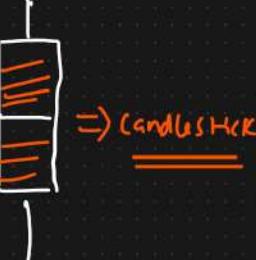
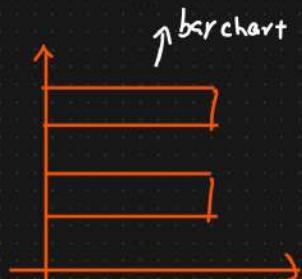
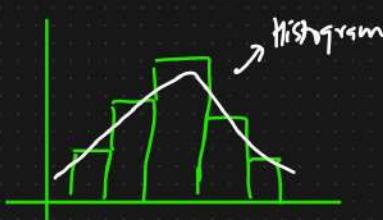
② Weights of Students in Classroom

Descriptive Stats [EDA + FF]

Inferential Stats

① It consists of organising and summarizing the data.

④ It consists of collecting sample data and making conclusion about population data using some experiments



University  $\rightarrow$  500 people

Class A  $\rightarrow$  60 people

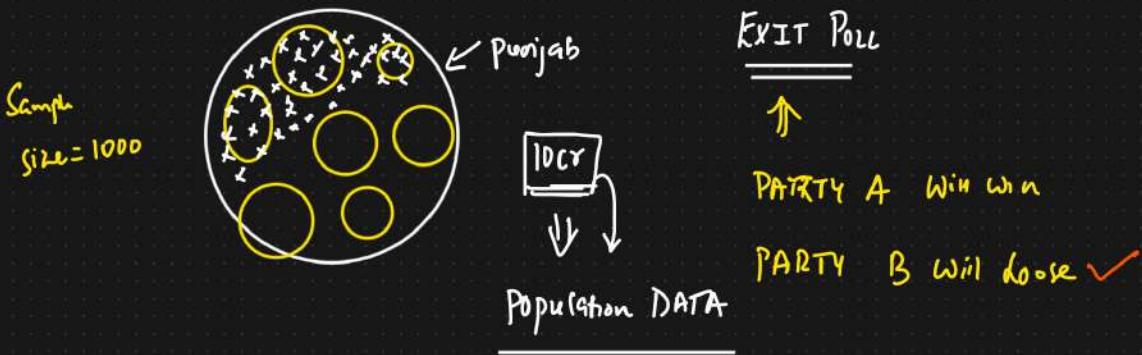
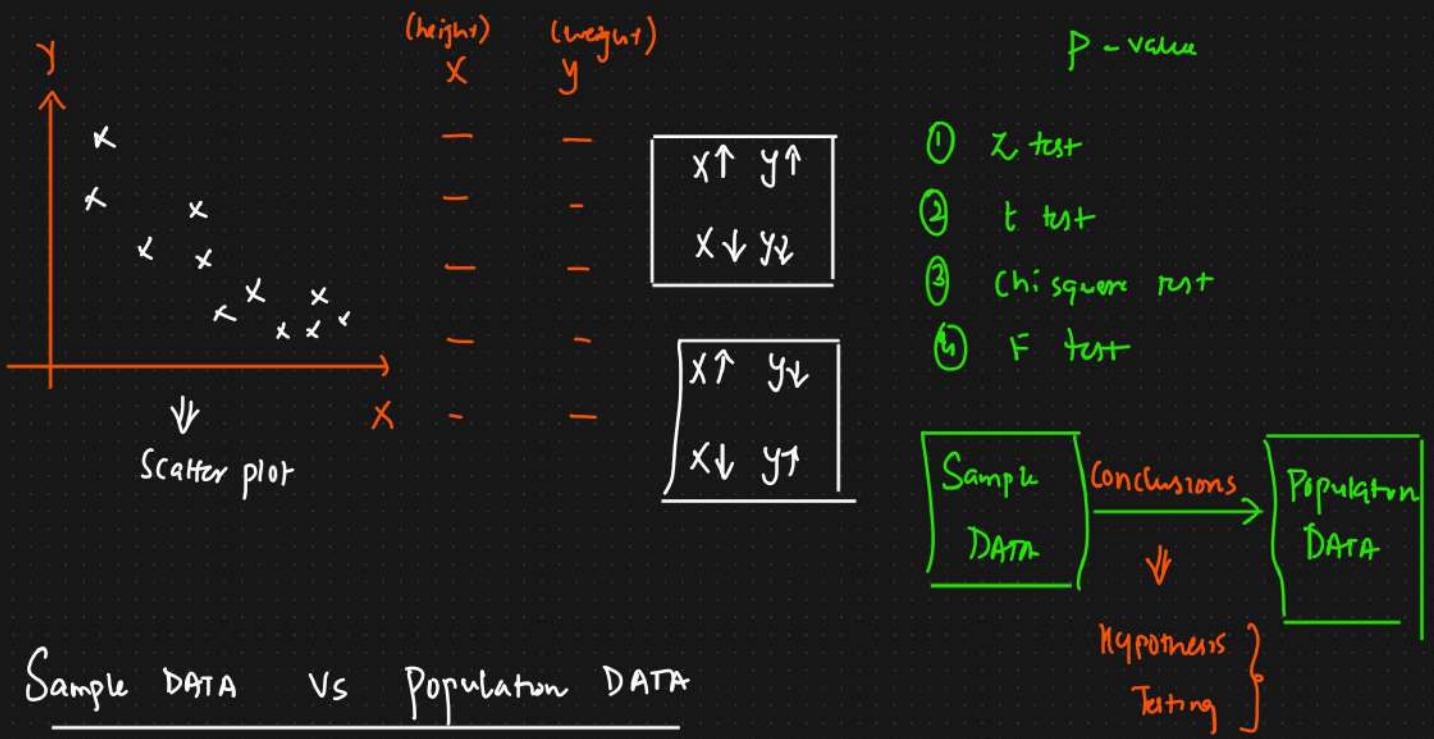
↓  
Sample data  $\Rightarrow$  Age  $\Rightarrow$  Average age of the entire university

[Yes]

↓

Hypothesis Testing

C.I  $\Rightarrow$  Confidence Interval



Eg: let's say there are 20 classrooms in a university and you have collected the age of students in one classroom

Ages  $\{21, 20, 18, 34, 17, 22, 24, 25, 26, 23, 22\}$

Weight  $\{ - \dots - \dots - \dots - \dots \}$

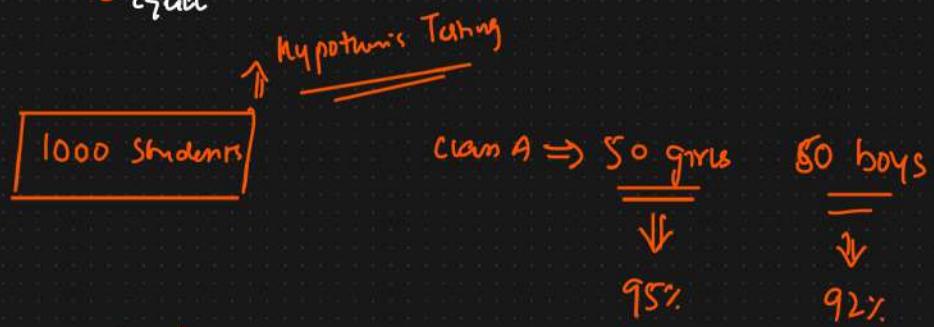
Descriptive Stats : What is the average age of students in the classroom?

Relationship between Age & Gender?

Inferential Stats : Are the average age of the students in the classroom

$\int \uparrow$  less than the average age of the students in the university?

⌈ Greater  
 ↓  
 Equal



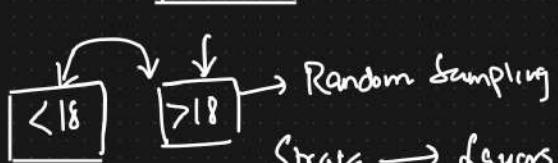
Choose a sample  
↑ Sampling Techniques

Population ( $N$ )      sample ( $n$ )

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



$n=1000$



Strata → Layers → Clusters ⇒ Groups

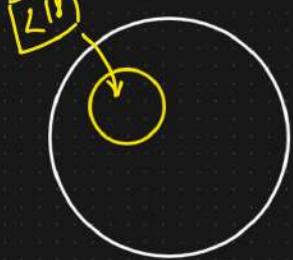
- ② Stratified Sampling

Gender  
→ Male  
→ Female

Education  
Degree  
→ High School  
→ Master  
→ Phd

Blood groups  
E

Population {Exit poll}.



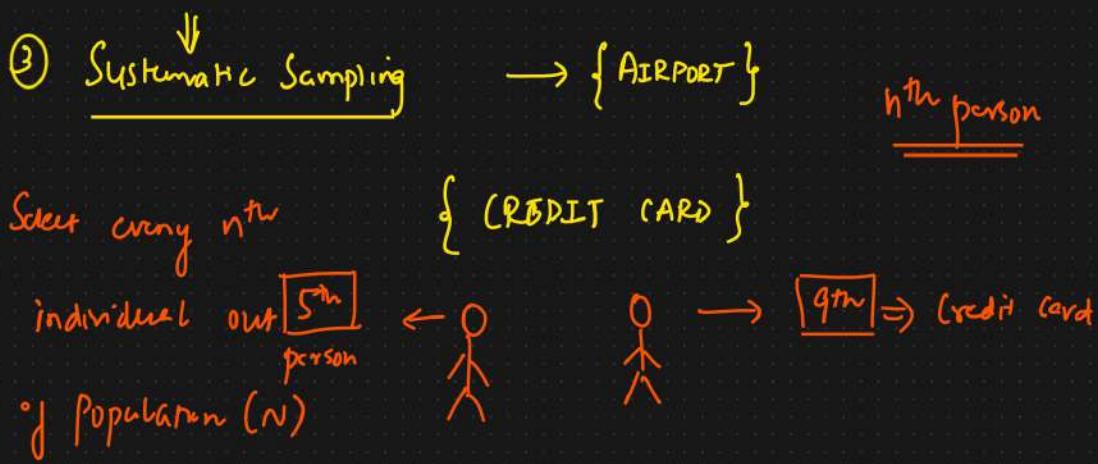
$<18$

↑  
Not Vote

↓  
Random Sampling

↓ Voting

$>18$



Select every  $n^{\text{th}}$  individual out of Population ( $N$ )

④ Convenience Sampling  $\div$  Only those who are interested in the survey  
will only participate

{ DATA SCIENCE SURVEY → General AI Survey }  $\uparrow$  inceron job for a specific  
Fill the Form }

① Survey Regarding New Technology  $\Rightarrow$  Convenience Sampling

② RBI Survey  $\Rightarrow$  Women  $\Rightarrow$  Stratified + Random Sampling  $\rightarrow$  Married Women

③ Credit Card  $\div$  Stratified + Random Sampling

① Variable : A Variable is a property that can take any values

Eg: age = 14      Variables

age = 25       $\text{Age} = [24, 25, 26, 27, 28, 29] \Rightarrow \text{Collection}$

age = 100

Two different types of Variable

① Quantitative Variable  $\rightarrow$  Measured Numerically {Mathematical Operations}.

Eg: Age, weight, height, rainfall(cm), temp, distance

② Qualitative Variables  $\rightarrow$  Categorical Variables {Based on some characteristics they are grouped together}.

Eg: Gender, Types of flowers, Types of Mammals

Quantitative Variable



Eg: Whole number  $\rightarrow$  fixed

Eg: No. of Bank Accounts

$\{1, 2, 3, 4, 5\}$        $\boxed{25} \times$

Eg: Continuous  $\rightarrow$  Decimal values

Eg: Height, weight, ages, Rainfall

Speed

Eg: No. of children  $\div$  Whole numbers

Pincode = fixed

Mammal       $\rightarrow$  Mammal       $\uparrow$  Categorical  
Non Mammal       $\rightarrow$  Non Mammal      variables.

## Assessment

① What kind of variable is Marital Status? Categorical variable

Gender ? Categorical

Length River length? Continuous

Movie duration? Continuous

Pincode ? Discrete

IQ ? Discrete

105.75, 90.5,

Pancard

Pincode Fixed Categorical

360099      }      ↓      Ans no. of Categories

120098  
960097

FE

⇒ It is many?

Categorical

Variables

Continuous



Discrete ←

Continuous

C

Whole number

Bank Account = { 1, 2, 3, 4, 5 }

Pincode = { }

Continuous

Gender Pincode

M —

F —

— —

— —

— —

— —

— —

— —

— —

PAN ←

{ }

⇒ Categorical

## Profile Building

- ① LinkedIn
  - ② GitHub
  - ③ Instagram
  - ④ Resume
  - ⑤ Mock Interviews
- ⇒ Jobs, Opportunities
- Freelancing

# Day 2 - Statistics

## Agenda

- ① Histograms ✓
- ② Measure of Central Tendency ✓ } ← 1.15 hrs
- ③ Measure of Dispersion ✓
- ④ Percentiles And Quartiles }
- ⑤ 5 Number Summary (Box plot) . }

## ① Histogram

$Ages = \{ 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100 \}$

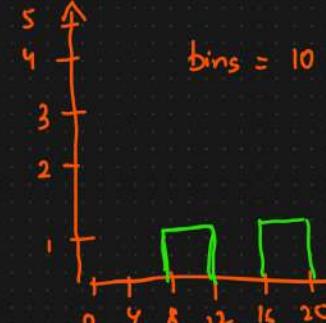
① Sort the Numbers

frequency (count)  
[10, 20, 25, 30, 35, 40]

min = 10

max = 40

② Bins → No. of groups



$$\frac{40}{10} = 4$$

③ Bins size → Size of Bins

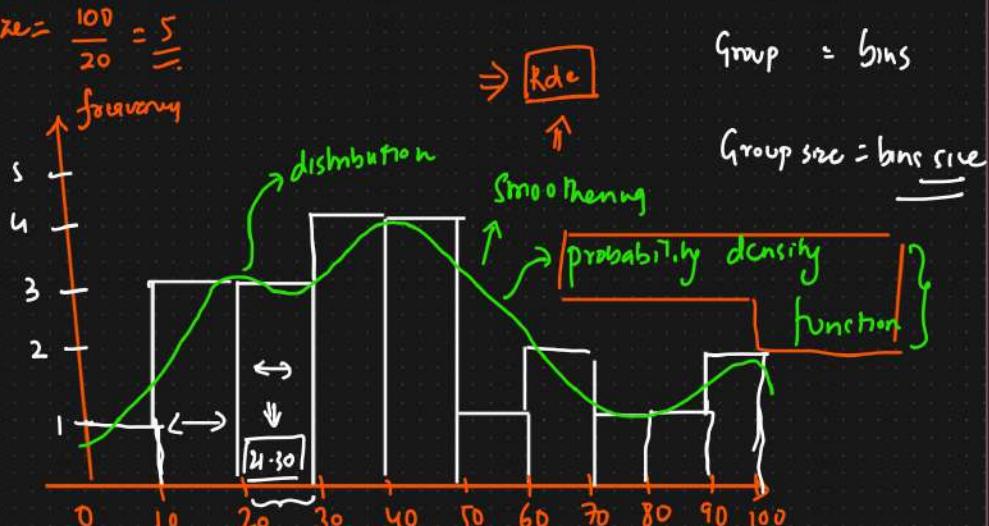
$$\underline{\underline{\text{bins} = 10}}$$

$$\text{bin size} = \frac{100}{20} = 5$$

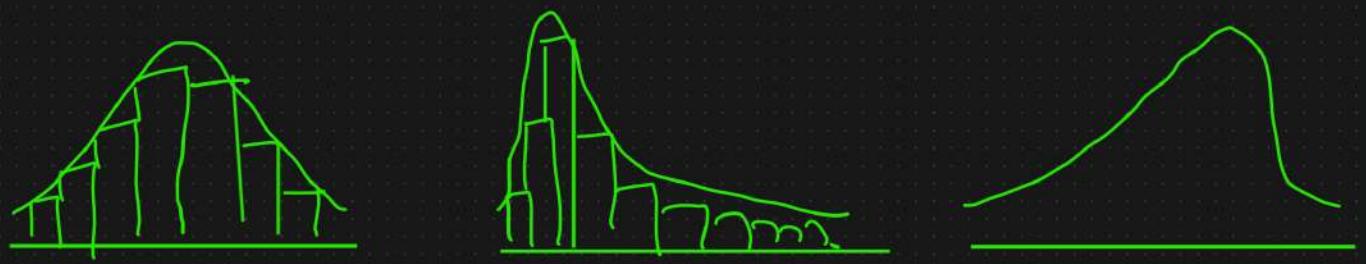
$$\text{bin size} = \frac{100}{10} = 10$$

$$\text{Bin size} = \frac{\text{Max} - \text{Min}}{\text{bins}}$$

$$\text{bins}:$$



$Ages = \{ 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100 \}$



Assignment

Weight = {30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 80, 90, 95}

bins = 10

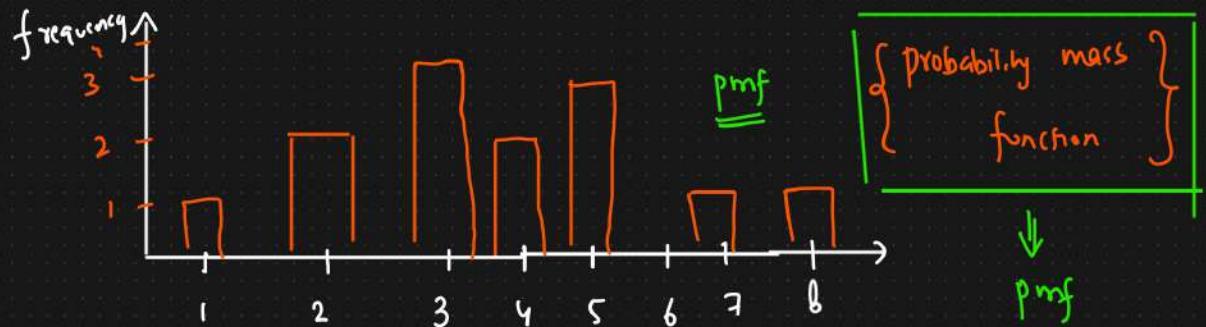
{continuous value}

$$\text{bin size} = \frac{95 - 30}{10} = \frac{65}{10} = 6.5$$

paf

## ② Discrete

No. of Banks accounts = [2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 4, 5]



pdf : probability density function } → continuous

pmf : probability mass function } → discrete

## ④ Measure of Central Tendency

① Mean ✓

{ A measure of CT is a single value that attempts to describe a set of data identifying the central position

② Median

③ Mode.

Mean  $X = \{1, 2, 3, 4, 5\}$  Average / Mean =  $\frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$

Population ( $N$ )

$$N \gg n$$

Sample ( $n$ )

Population mean ( $\mu$ ) =  $\left[ \sum_{i=1}^N \frac{x_i}{N} \right] \quad N \gg n$

Sample mean ( $\bar{x}$ ) =  $\left[ \sum_{i=1}^n \frac{x_i}{n} \right]$

$$N = 6$$

$$N > n$$

$$n = 4$$

$$\{24, 23, 28, 27\} \leftarrow \text{Age}$$

Population Age =  $\{24, 23, 21, 28, 27\}$

Sample Age =  $\{24, 21, 28, 27\}$

$$\begin{matrix} 13 \\ 59 \\ 41 \end{matrix}$$

Population mean ( $\mu$ ) =  $\frac{24+23+21+28+27}{6}$

$$\mu = 17.5$$

Sample mean ( $\bar{x}$ ) =  $\frac{24+21+28+27}{4}$

$$\bar{x} = 13.5$$

[hp.nan]  $\leftarrow$  Null values

### Practical Application (Feature Engineering)

Age	Salary	Family Size
-	-	-
-	-	-
-	-	-
NAN	-	-
-	-	-
-	NAN	-
-	-	NAN
-	NAN	-
NAN	-	-

$\leftarrow$  loss of Info

$$\text{Age} = 29.6$$

$$\downarrow \downarrow \downarrow$$

$$38$$

Mean

NAN  $\rightarrow$  Not A Number

$$\begin{matrix} \text{NULL} \\ \downarrow \\ 10/4 = \frac{1, 2, 3, 4}{\text{NAN}} \end{matrix} \quad \text{mean} \uparrow \text{NAN}$$

Age	Salary
24 ✓	45
28 ✓	50
29 ✓	NAN
NAN X	60
31 ✓	75
36 ✓	80
NAN X	NAN

$$\downarrow \downarrow \downarrow$$

$$85$$

$$\text{Outliers} \leftarrow [80] [200] \leftarrow$$

### ① Median

$$\begin{aligned} \{1, 2, 3, 4, 5\} &= \{1, 2, 3, 4, 5, \boxed{100}\} \\ \bar{x} = 3 &\quad \longrightarrow \bar{x} = 19.16 \end{aligned}$$

Outlier

$$= \frac{1+2+3+4+5+100}{6} = 19.16$$

### Steps to find out median

① Sort the Numbers

② Find the central number

①

{ if the no. of elements are even we find the average of central elements }

② if the no. of elements are odd we find the central elements.

Sorted

$$\{0, 1, 2, 3, 4, \boxed{5, 6}, 7, 8, 100, 120\}$$

Mean =  $\frac{25.6}{10}$

$$\text{median} = \frac{5+6}{2} = 5.5$$

median = 5

③ Mode : {Most frequent occurring elements}

$$\{1, 2, 2, 2, \boxed{3, 3, 3}, 4, 5\}$$

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

$$\boxed{2, 3}$$

## Dataset

Types of flower      (categorical variable)

lily

Sunflower

Rose

NAN ← ROSE

Rose

Sunflower

Rose

NAN ← Rose

$\boxed{40\%}$   $\Rightarrow$  Biased

Under 19

17, 18, 19, 16, 15,  $\boxed{32} \rightarrow$  outlier

## (F) Measure of Dispersion

① Variance ( $\sigma^2$ )  $\leftarrow$  Spread of Data.

② Standard deviation ( $\sigma$ )  $\leftarrow$

$$X = \{1, 2, 3, 4, 5\} \quad \mu = 3$$

### Variance

Population Variance ( $\sigma^2$ )  $\left\{ \begin{array}{l} \text{Degree of freedom} \\ \text{Bands (correction)} \end{array} \right.$

$$\sigma^2 = \frac{\sum_{i=0}^N (x_i - \mu)^2}{N}$$

Sample Variance ( $s^2$ )

$$s^2 = \frac{\sum_{i=0}^n (x_i - \bar{x})^2}{n-1}$$

$$\left\{ 1-100 \right\} \downarrow \quad \left\{ 1-100 \right\} \downarrow \quad \text{Second one.}$$

$= \Downarrow \text{First one.} = \Downarrow$

$\boxed{\text{Assignment}} \leftarrow$

$$\left\{ 1, 2, 3, 4, 5, 6, 7, 8, 9, 10 \right\} \xrightarrow{\text{Variance}} \left\{ 1, 2, 3, 4, 50, 60, 70, 100 \right\} \xrightarrow{\text{Variance}} \frac{1-100}{7}$$

$$\left\{ 1, 2, 3, 4, 5 \right\} \xrightarrow{\text{Variance}} \left\{ 1, 2, 3, 4, 5, 6, 80 \right\} \xrightarrow{\text{Variance}} \frac{21}{7}$$

$$\mu = 3$$

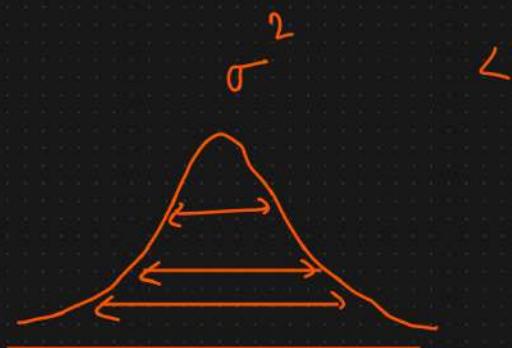
$$\mu = 14.4$$

$$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5} = \frac{10}{5} = 2$$

$$\sigma^2 = \frac{(1-14.4)^2 + (2-14.4)^2 + (3-14.4)^2 + (4-14.4)^2 + (5-14.4)^2 + (6-14.4)^2 + (80-14.4)^2}{7}$$

Variance ↑ Spread ↑↑

$$\sigma^2 = 719.10$$



$$\textcircled{1} \text{ Standard deviation } (\sqrt{\sigma^2}) \Rightarrow \boxed{1.41}$$

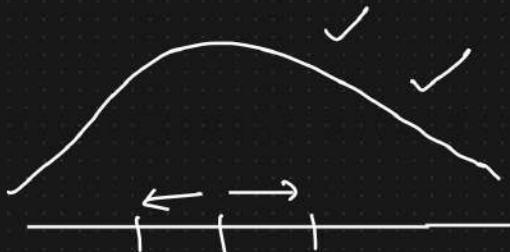
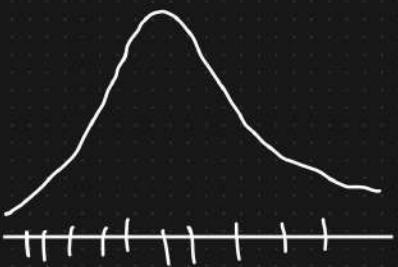
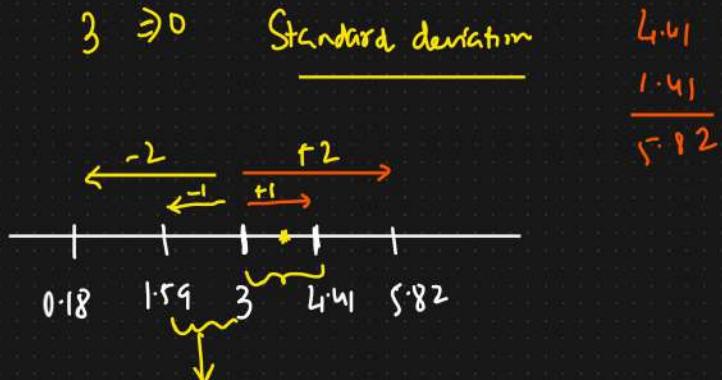
$$\begin{array}{r} 3.00 \\ 1.41 \\ \hline 1.59 \end{array}$$

$$\{ 1, \boxed{2}, 3, \boxed{4}, 5 \}$$

$$M = 3$$

$$\sigma^2 = 2$$

$$\sigma = \sqrt{2} = 1.41$$



## ④ Percentiles And Quartiles

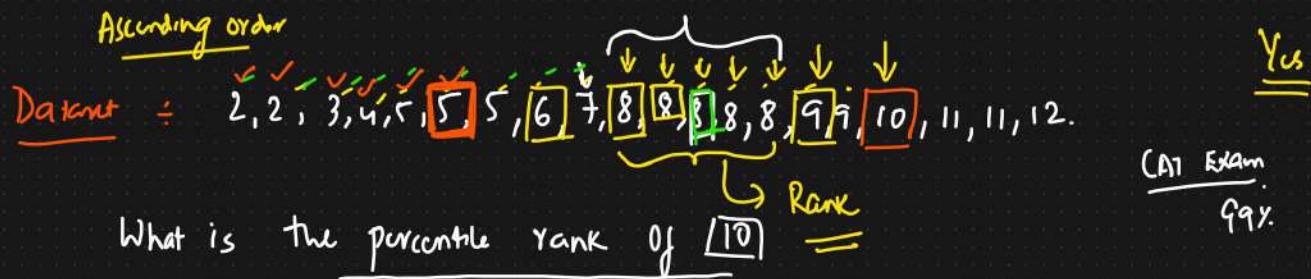
$$\text{Percentage} = \{ 1, 2, 3, 4, 5, 6, 7, 8 \}$$

$$\text{Percentage of Even Number} = \frac{\text{No. of even Numbers}}{\text{Total No. of Number}} = \frac{4}{8} = 0.5 = 50\%$$

Percentiles : CAT, IELTS, SAT, GRE, JEE, NEET  $\Rightarrow$  Percentiles

Defn : A percentile is a value below which a certain percentage of observations lie.

99 percentile = It means the person has got better marks than 99% of the entire students



$$\text{Next item} = 0.8$$

$$\text{Percentile Rank of } x = \frac{\#\text{No. of value below } x}{n} = \frac{16}{20} = 80 \text{ percentile}$$

45 percentile

$$= \frac{14}{20} = 70 \text{ percentile}$$

④ What is the value that exists at 25 percentile

75%

$$\text{Value} = \frac{\text{Percentile}}{100} * \frac{n+1}{(n)}$$

$$= \frac{25}{100} * \frac{20}{20} = 5^{\text{th}} \text{ Index}$$

$$\text{DfP} = 5$$

$$= \frac{95}{100} * 21$$

⑥ 5 number Summary

① Minimum

② First Quartile (25 percentile) (Q1)

③ Median

④ Third Quartile (75 percentile) (Q3).

Box plot  
↑  
⇒ Remove the outliers.

### ⑤ Maximum

$$\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \text{[2]}\} \quad \downarrow \text{outlier} =$$

$\frac{\downarrow}{15} \frac{\downarrow}{16} \frac{\downarrow}{5.25}$



[Lower Fence  $\longleftrightarrow$  Higher Fence]

$$\Downarrow [-3.65 \longleftrightarrow 14.25]$$

$$\leftarrow \text{lower Fence} = Q_1 - 1.5(IQR) \leftarrow$$

$$\text{Higher Fence} = Q_3 + 1.5(IQR) \leftarrow$$

$$IQR = Q_3 - Q_1 = 25$$

$$\Downarrow = =$$

Inter Quartile Range (IQR)

$$Q_1 = \frac{25}{100} \times 21 = 5.25 \quad \text{Index} = 3 =$$

$$Q_3 = \frac{75}{100} \times 21 = 15.75 \quad \text{Index} = \frac{8+7}{2} = 7.5 =$$

$$\text{lower Fence} = 3 - (1.5)(4.5) = \boxed{-3.65}$$

$$\text{Higher Fence} = 7.5 + (1.5)(4.5) = \boxed{14.25}$$

$$\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \text{[2]}\}.$$

-5

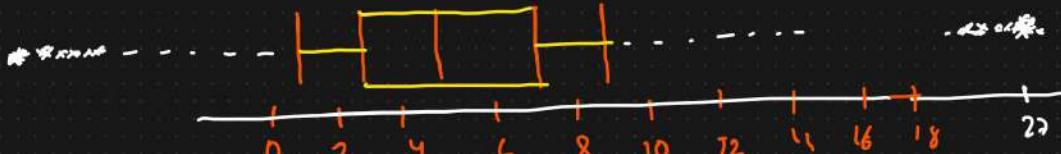
5 Number Summary.

① Minimum = 1 ✓

Box Plot



②  $Q_1 = 3 \checkmark$



③ Median = 5 ✓



④  $Q_3 = 7.5 \checkmark$

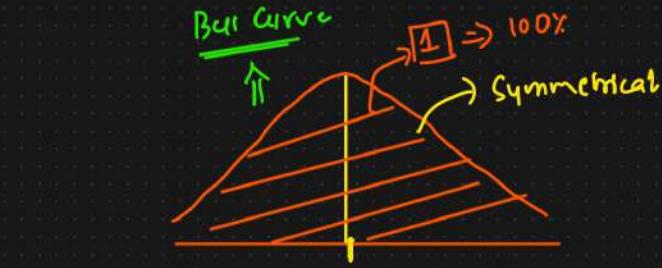
To Treat outliers

⑤ Maximum = 9 ✓

## Day 3 - Stats

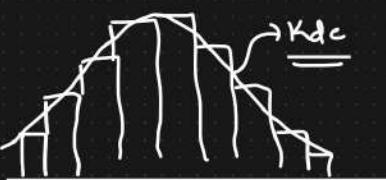
- ① Normal Distribution ✓
- ② Standard Normal Distribution ✓
- ③ Z-score ✓
- ④ Standardization And Normalization ✓.
- ⑤ Gaussian / Normal Distribution



↓  
Age, weight, height  
↑  
Distribution  
⇒ Doctors

Domain Expertise

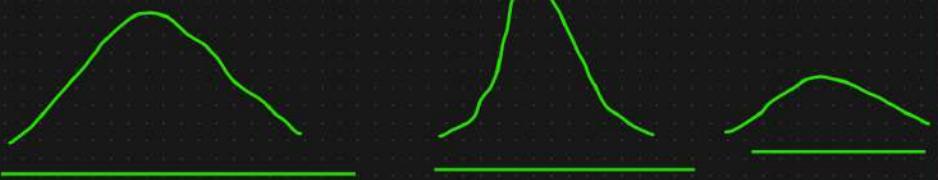
Kernel density estimator



[IRIS DATASET] ←  
↓

Petal length, Sepal length, petal width,  
↓  
Sepal width

Gaussian Distri



- ① [Empirical Rule of Normal Distribution]

Empirical Rule

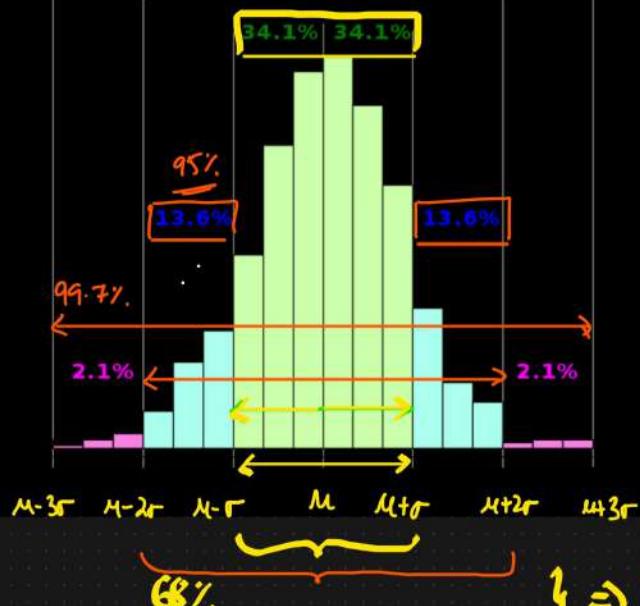


$$68 - 95 - 99.7\%$$



CLEAR

$$\text{Age} = \{$$



Gaussian / Normal Distribution



Assumptions of  
the data

}  $\Rightarrow$  Gaussian / Normal Dist



[Q-Q plot]  $\Rightarrow$  Distribution is Gaussian Or Not?

Standard Normal Distribution

$X \sim \text{Gaussian Distribution } (\mu, \sigma)$

$$X = \{1, 2, 3, 4, 5\}$$



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

$$\mu = 3$$

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

$$\sigma = 1.41$$

= =

$$Z\text{-score} = \frac{x_i - \mu}{\sigma} \quad [n=1]$$

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

$\Rightarrow$  Standard Error  $\Rightarrow$  Inferential stats.

$$Z\text{-score} = \left| \frac{x_i - \mu}{\sigma} \right| \leftarrow \text{Simple}$$

$$X = \{1, 2, 3, 4, 5\}$$

$$\mu = 3 \quad \sigma = 1.414$$

$$= \frac{1-3}{1.414} = -1.414$$

↗  $y = \{-1.414, -0.707, 0, 0.707, 1.414\}$

$$= \frac{2-3}{1.414} =$$

$$\frac{4-3}{1.414} = \frac{1}{1.414} =$$



Why?

[Standardization]  $\Rightarrow [\mu=0 \quad \sigma=1]$

	(years) <u>Age</u>	(kg) <u>Weight</u>	(cm) <u>Height</u>
$\mu=0$	24	72	150
$\sigma=1$	26	78	160
	32	84	165
$\Rightarrow$	33	92	170
	34	87	150
	28	83	180
	29	80	175

[0-1]

Same Scale



Machine Learning

Maths Questions



Algorithm  $\Rightarrow$  Mathematical Model

mathematical

calculation Time ↑↑↑

Height



Feature Scaling

✓ Normalization  $\Rightarrow$

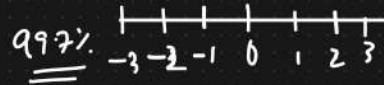
Standardization  $\{\downarrow Z\text{-score}\}$

$$\frac{x_i - \mu}{\sigma}$$



$$\begin{bmatrix} 0 & 1 \end{bmatrix} \quad \begin{bmatrix} -1 & 1 \end{bmatrix}$$

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



$$\begin{bmatrix} 0 & 5 \end{bmatrix}$$

$$\begin{bmatrix} -3 & \leftrightarrow 3 \end{bmatrix}$$

$$\begin{bmatrix} 0 & 4 \end{bmatrix}$$

Normalization [lower scale  $\leftrightarrow$  higher scale]  $\rightarrow$  [Images  $\Rightarrow$  0-255] Standard

① Min Max Scaler [0-1]

$$x_{scaled} = \frac{x - x_{min}}{x_{max} - x_{min}}$$

$$= \frac{1-1}{5-1} = 0$$

$$\frac{4-1}{5-1} = \frac{3}{4}$$

$$\frac{3-1}{5-1} = \frac{2}{4}$$

$$\frac{2-1}{5-1} = \frac{1}{4}$$

$$\frac{1-1}{5-1} = 0$$

1
2
3
4
5

1
2
3
4
5

Imagis  $\Rightarrow$  0-255

$y$	$\downarrow$	$y'$
0		-1.414
0.25		-0.302
0.5		0
0.75		0.302
1		1.414

$\Downarrow$  [RNN] [ANN]

Apply ??



Deep learning



$0-255 \rightarrow 0-1$

Normalization

$B_1$
$B_2$
$B_3$

$P_{HD}$

NASA

① Standardization

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

$$\left\{ \begin{array}{l} X \rightarrow \text{Normal Distribution } (\mu, \sigma) \\ \Downarrow Z\text{-score} \\ Y \rightarrow \text{SND } (\mu=0, \sigma=1) \end{array} \right.$$



Why do we do this  $\rightarrow$  Bring the features in the same scale

Normalization

$[0 - 1]$



① Min Max Scaler

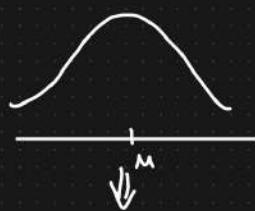


ML

Standardization



ML

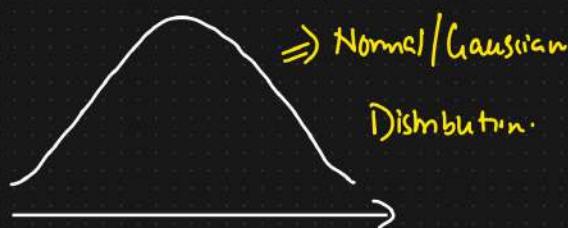


Min Max Scaler

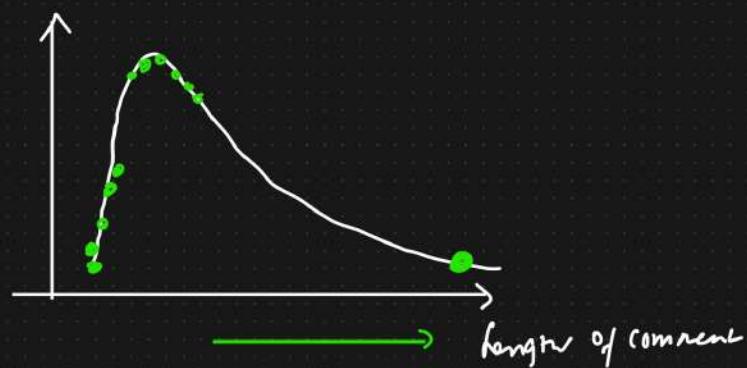
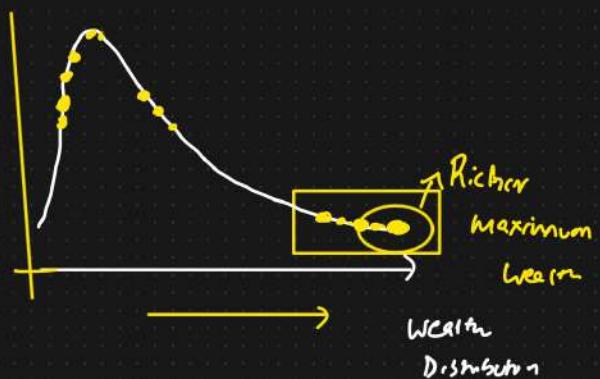


Standardization

① log Normal Distribution



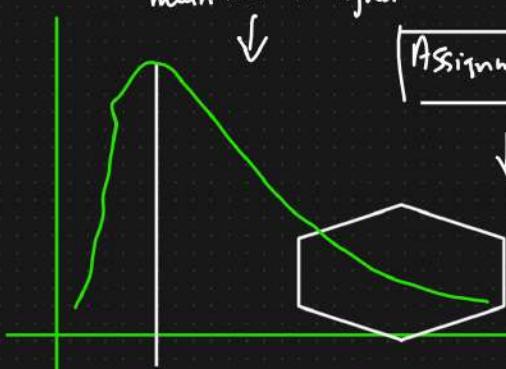
log Normal Distribution



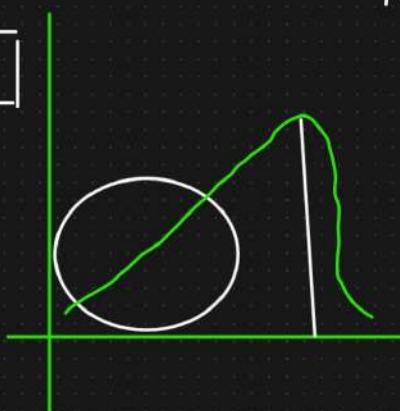
mean will be higher



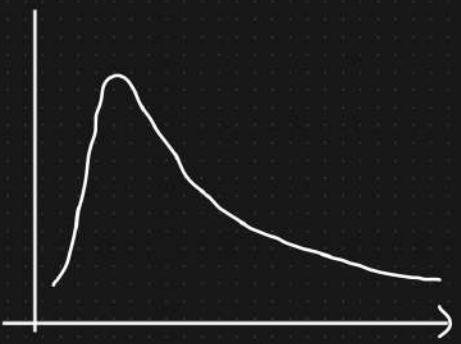
Assignment



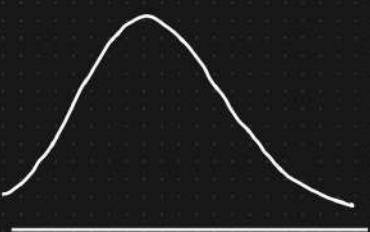
Relation of mean,  
median, mode



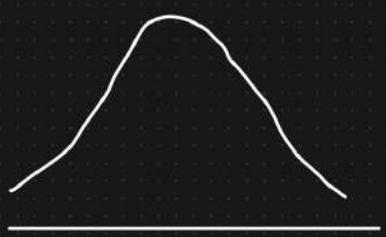
From Ascending order give the relation of mean, median & mode?



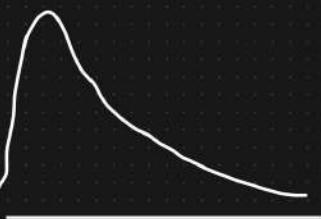
$$Y = \ln(X)$$



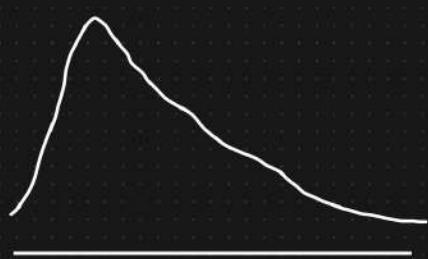
$X \sim \text{log Normal Distribution}$



$$\begin{array}{c} \text{Antilog} \\ \Downarrow \\ \Rightarrow \exp(Y) \\ \Downarrow \\ Y \end{array}$$

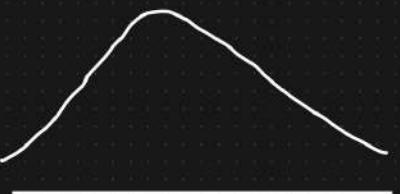


(\*)



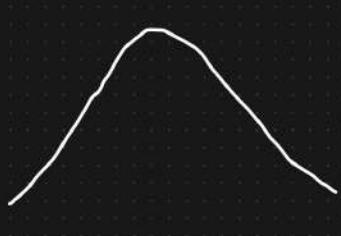
Natural log  
↑  
 $\log_e$

$$\Rightarrow Y = \ln(X)$$

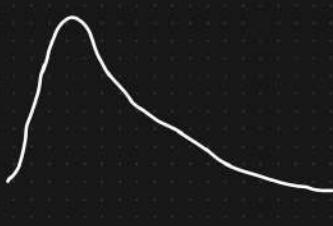


$X \sim \text{log Normal Distribution}$   
 $(\mu, \sigma^2)$

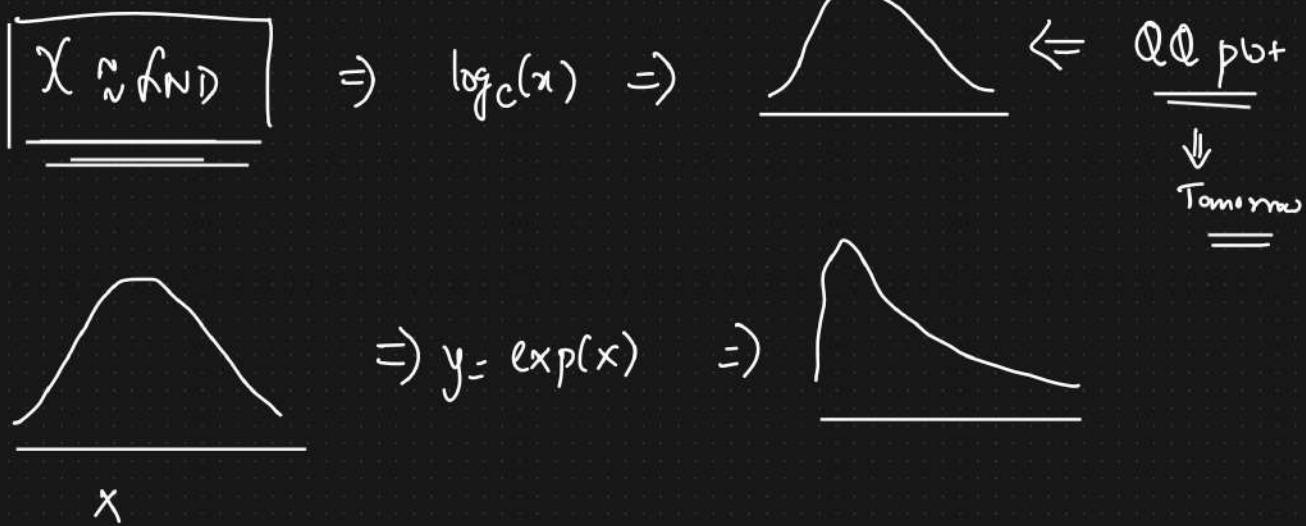
$\Downarrow$   
Inverse



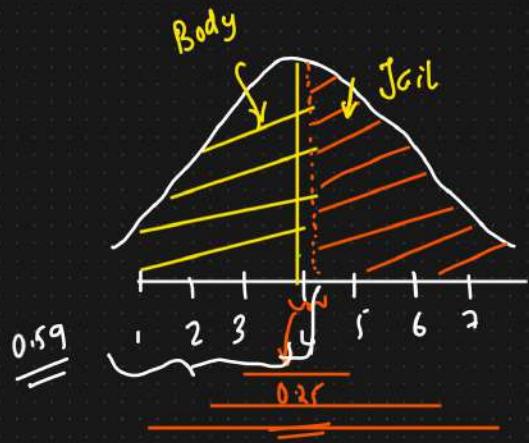
$$\Rightarrow X = \exp(Y)$$



Y



⑥  $X = \{1, 2, 3, 4, 5, 6, 7\}$   $M = 4$   
 $f = 1$



Question: What is the percentage of score

that falls above 4.25?

fall below 3.75?

0.59  $\Rightarrow 59\%$

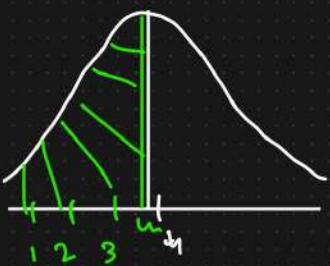
$1 - 0.59 = 0.41 \Rightarrow 41\%$

① Z-score =  $\frac{x - M}{\sigma} = \frac{4.25 - 4}{1} = \boxed{0.25}$

① Z-table (area under the curve)



Z-score =  $\frac{3.75 - 4}{1} = -0.25$   $\approx 40\%$





$$4.25 \quad 4.5 \quad 4.75$$

- ⑥ In India the average IQ is 100 with a Standard Deviation of 15. What is the percentage of population would you expect to have an IQ

Answers

- ① Lower than 85 = 0.1587
- ② Higher than 85 = 0.8413
- ③ Between 85 and 100 = 0.3413.

Assignment

## Day 4-STATS

- ① Central Limit Theorem. ✓.
- ② Probability. ✓
- ③ Permutation And combination ✓
- ④ Covariance, Pearson Correlation, Spearman Rank Correlation. } ✓

⑤ Bernoulli Distribution

⑥ Binomial Distribution

⑦ Power law (Pareto Distribution).

① { Central Limit Theorem }

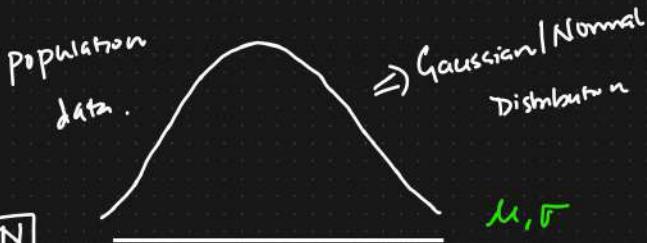
$$\begin{matrix} n < 30 \\ \downarrow \downarrow \\ n > 30 \\ \uparrow \end{matrix}$$

Size of sample

The larger the value the better

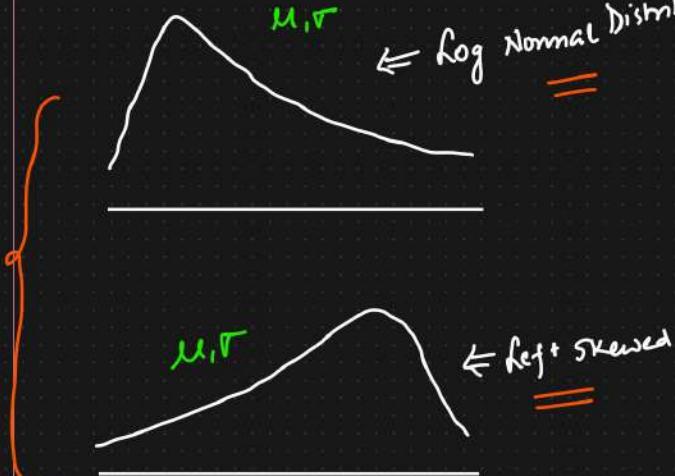
$$\begin{matrix} \uparrow \\ n \\ \rightarrow m \end{matrix}$$

No. of samples



$$\begin{aligned} \rightarrow S_1 &\rightarrow \{x_1, x_2, x_3, \dots, x_n\} \rightarrow \bar{x}_1 = \bar{s}_1 \\ \rightarrow S_2 &\rightarrow \{x_3, x_4, \dots, x_1, \dots, x_n\} \rightarrow \bar{x}_2 = \bar{s}_2 \\ \rightarrow S_3 &\rightarrow \{x_n, x_1, \dots, x_{n-1}\} \rightarrow \bar{x}_3 = \bar{s}_3 \\ &\vdots \\ &\vdots \\ \bar{x}_m &= \bar{s}_m \end{aligned}$$

Sampling with replacement



10 different region

$$n > 30$$

Size of shark through out the world?

Assumptions

$N < 30$

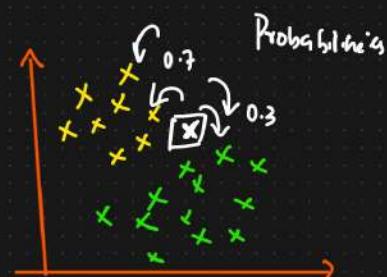
$m \uparrow \uparrow \uparrow$

② Probability: Probability is a measure of the likelihood of an event

Eg: Tossing a fair coin  $P(H) = 0.5$   $P(T) = 0.5$

$\downarrow$   
Shoray  $\rightarrow$  COIN  $P(H) = 1$   
 $\underline{=}$   
unfair coin

Strong  
Basic  
 $\uparrow$



Rolling a Dice

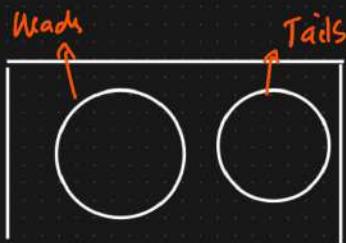
$$P(1) = \frac{1}{6} \quad P(2) = \frac{1}{6} \quad P(3) = \frac{1}{6}$$

① Mutual Exclusive Event

Two events are mutually exclusive if they cannot occur at the same time.

① Tossing a coin

② Rolling a dice



② Non Mutual Exclusive Events

Two events can occur at the same time.



0 0

Bag of  
marbles

④ Picking randomly a card from a deck of cards, two events "heart" and "king" can be selected.

## Mutual Exclusive Event

① What is the probability of coin landing on heads or tails



Addition Rule for mutual exclusive events

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

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

② What is the probability of getting 1 or 6 or 3 while rolling a dice?

$$P(1 \text{ or } 6 \text{ or } 3) = P(1) + P(6) + P(3)$$

$$= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2}$$

## Non Mutual Exclusive Event

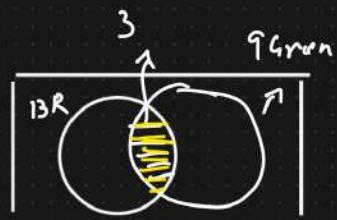
Bag of Marbles : 10 Red, 6 Green, 3 (R&G) R&G

① When picking randomly from a bag of marbles what is the probability of choosing a marble that is red or green?



Non mutual Exclusive

Addition Rule for Non Mutual Exclusive Event



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

$$= \frac{13}{19} + \frac{9}{19} - \frac{3}{19} = \frac{19}{19} = \underline{\underline{1}}$$

Deck of cards  $\rightarrow$  What is the probability of choosing Q or Queen

$$P(Q \text{ or Queen}) = P(Q) + P(\text{Queen}) - P(Q \text{ and Queen}) \\ = \frac{13}{52} + \frac{4}{52} - \frac{1}{52} = \boxed{\frac{16}{52}}$$

### ④ Multiplication Rule

④ Dependent Events : Two events are dependent if they affect one another

Bag of marbles  $\left\{ \begin{array}{l} 0 \ 0 \ 0 \ X \\ 0 \ 0 \ 0 \end{array} \right\}$

$$\Rightarrow P(W) = \frac{4}{7} \xrightarrow{\substack{\text{↑} \\ \text{white} \\ 1 marble}} P(Y) = \frac{3}{6}$$

④ What is the probability of rolling a "5" and then a "3" with a normal 6 sided dice?

Ans)  $P(1) = \frac{1}{6} \quad P(2) = \frac{1}{6} \quad P(3) = \frac{1}{6} \quad P(4) = \frac{1}{6}$  of Independent Events

Multiplication Rule for Independent Events

$$P(A \text{ and } B) = P(A) \times P(B) \\ = \frac{1}{6} * \frac{1}{6} = \boxed{\frac{1}{36}}$$

$P(A \text{ or } B) \Rightarrow$

- Mutual Exclusion
- Non Mutual Exclusion

$$\overbrace{P(K \text{ or } R)}^{\text{common}} = P(A \text{ or } B) = P(A) + P(B) - \boxed{P(A \text{ and } B)} \rightarrow \text{Non Mutual Exclusive.}$$

$$P(A \text{ or } B) = P(A) + P(B) \quad [\text{Mutual Exclusion}]$$

Dependent and Independent Events

Event A      Event B

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

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

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

②

$$\left. \begin{matrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \right\}$$

$\Rightarrow$  Dependent Events

Probability of drawing a "Orange" and then drawing a "Yellow"

marble from the bag?

Ans)

$$\left. \begin{matrix} 0 & 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \right\} \quad P(O) = \frac{4}{7}$$

Orange Marble

$$\xrightarrow{\quad \quad \quad P(Y|O) \quad \quad \quad \text{conditional probability}}$$

$$\left. \begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \right\} \quad \left[ \frac{3}{6} = \frac{1}{2} \right]$$

Naive Bayes

$$P(O \text{ and } Y) = P(O) * \boxed{P(Y|O)}$$

$$= \frac{4}{7} * \frac{3}{6} = \frac{4}{7} * \frac{1}{2} = \frac{2}{7},$$

## ④ Permutation

School of Children

$$\begin{array}{c} 5 \\ \hline \end{array} \times \begin{array}{c} 4 \\ \hline \end{array} \times \begin{array}{c} 3 \\ \hline \end{array}$$

$\left\{ \begin{array}{l} \text{Dairy Milk, Kit Kat, Milky Bar,} \\ \text{Sneakers, 5 star} \end{array} \right\}$

$$= \boxed{\underline{60 \text{ ways}}} \Rightarrow \text{Permutation}$$

With permutation, order matters

$$\begin{array}{c} \{ DM \quad KK \quad MB \} \quad \{ \quad \} \quad \{ \quad \} \\ \{ KK \quad DM \quad MB \} \quad \{ \quad \} \quad \{ \quad \} \\ \{ \quad \} \quad \{ \quad \} \quad \{ \quad \} \quad \{ \quad \} \end{array}$$

Possible Arrangements }.

$n=5$

$$n_p_r = \frac{n!}{(n-r)!} = \frac{5!}{(5-3)!}$$

$$= \frac{5 \times 4 \times 3 \times 2!}{2!} = \boxed{\underline{60}}$$

$n$  = Total No. of Objects

$r$  = # of selection

## ⑤ Combination

Repetition will not occur

$$\{ DM \quad KK \quad MB \}$$

Unique Combination

$$\times \{ MB \quad KK \quad DM \} \leftarrow$$

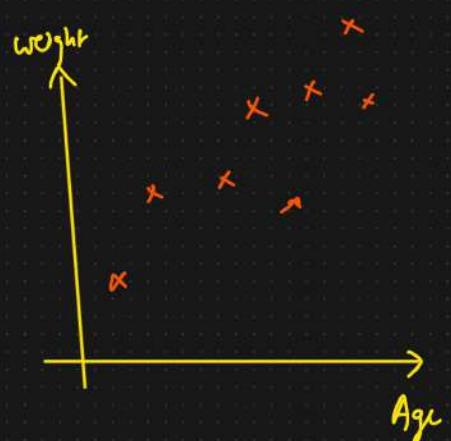
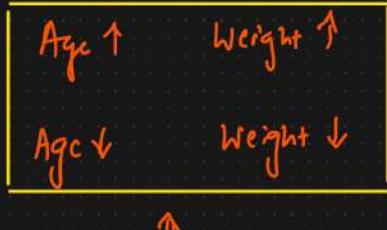
$$n_c_r = \frac{n!}{r!(n-r)!} = \frac{5!}{3!(2)!} = \frac{5 \times 4 \times 3 \times 2 \times 1}{3! \times 2!} = \frac{120}{12} = 10$$

# DREAM 11

① Covariance ✓

X	Y
Age	Weight
12	40
13	45
15	48
17	60
18	62
$\bar{x} = 15$	$\bar{y} = 51$

{Feature Selection}



Quantity the relationship

x & y using mathematical  
question

$$\left[ \text{Cov}(x, y) = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n-1} \right] \quad \sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n-1}$$

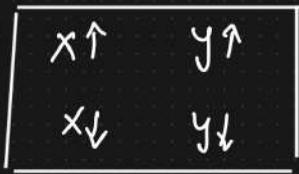
$$\left[ \text{Cov}(x, x) = \frac{\sum (x_i - \bar{x})(x_i - \bar{x})}{n-1} \right]$$

$\Downarrow$

$\boxed{\text{Cov}} \Rightarrow \text{tvc}$        $\text{Cov}(x, x)$        $\Leftarrow$        $\boxed{\sigma^2 = \frac{\sum (x_i - \bar{x})(x_i - \bar{x})}{n-1}}$

$$\left[ \text{Cov}(x, x) = \text{Var}(x) \right] \Leftarrow$$

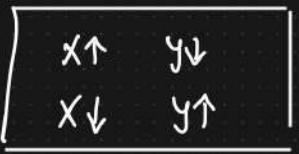
tvc Covariance

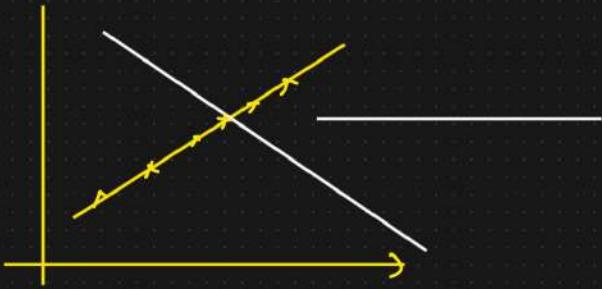


Covariance = 0

No relation  
with x & y

-ve Covariance





$X$	$y$
10	4
8	6
7	8
6	10
$\bar{x} = 7.75$	$\bar{y} = 7$

$$\text{Cov}(x, y) = \underline{\underline{-ve}}$$

$$= \left[ (10 - 7.75)(4 - 7) + (8 - 7.75)(6 - 7) + (7 - 7.75)(8 - 7) + (6 - 7.75)(10 - 7) \right]$$

$$= -3.25$$

$$\begin{pmatrix} x \uparrow & y \downarrow \\ x \downarrow & y \uparrow \end{pmatrix}$$

$\underline{\underline{}}$

① Pearson Correlation Coefficient (-1 to 1)

Scale

-ve Covariance = +ve

+ve Covariance = -ve

$$\begin{matrix} \curvearrowleft & \curvearrowright \\ x & y & z \end{matrix}$$

$$\begin{matrix} \curvearrowleft & \curvearrowright \\ x & y & z \end{matrix} \Rightarrow [0.7] \checkmark$$

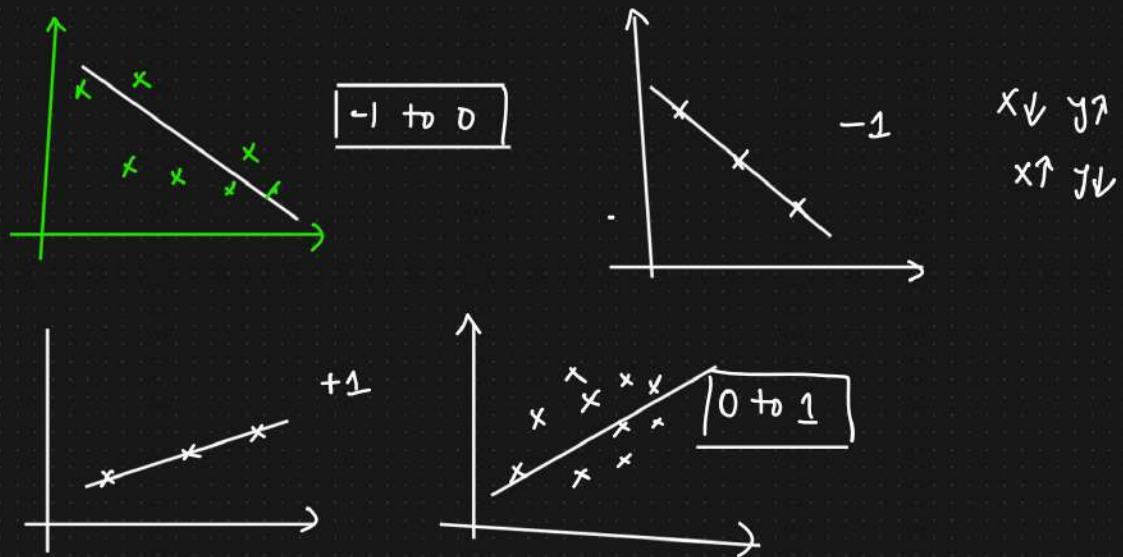
$$[0.5]$$

More the value towards +1

More +ve correlated it is

-1

negative correlated



## (f) Spearman Rank Correlation

$$\gamma_s = \frac{\text{Cov}(R(x), R(y))}{\sigma(R(x)) \sigma(R(y))}$$

Ascending Order

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

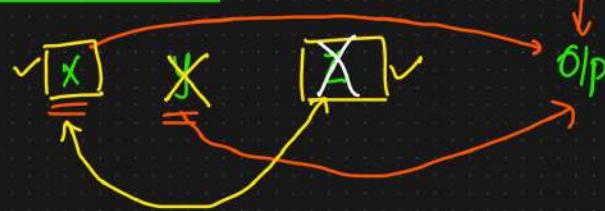
Spearman Rank Correlation

Ascending

Scatter plot showing  $x$  and  $y$  values in ascending order. A green X marks the point (6, 10).

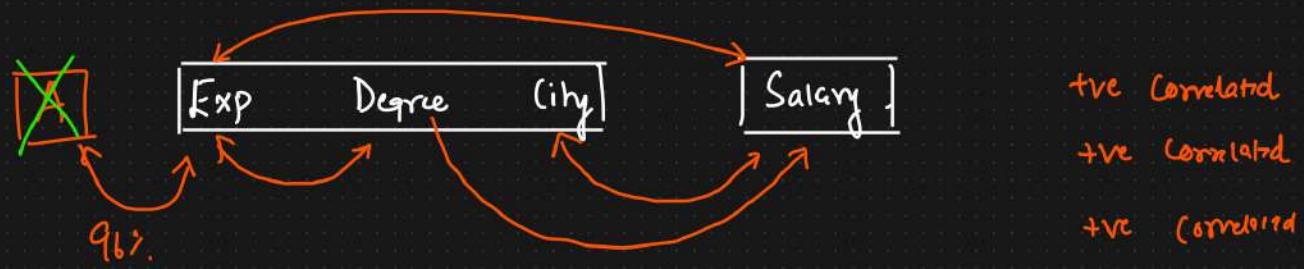
Why this Correlation will be used?

$$\begin{matrix} 0.95 \\ \hline 95\% \end{matrix}$$



$$\begin{matrix} \downarrow \\ \text{tive} \\ -\text{ve} \\ \uparrow \end{matrix} \quad \left. \begin{matrix} \text{Good} \\ \underline{\underline{}} \end{matrix} \right\}$$

$$0.2 \quad 0.01$$



## Inferential Statistics

- ① Hypothesis Testing
- ② p-value
- ③ Confidence Interval
- ④ Significance Value

$\chi^2$  test  
t test  
Chi square test  
Anova test (F-test)

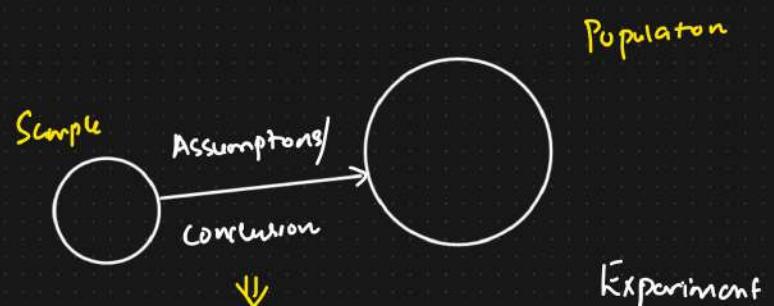
### 3 Distributions

- ① Bernoulli's
- ② Binomial
- ③ Power Law

TRANSFORMATION

## Inferential Stat

### Steps of hypothesis Testing



### Experiment

① Null Hypothesis: Coin is fair  $\Rightarrow$  Accepted  $\rightarrow$  [Coin is fair or not]

② Alternative Hypothesis: Coin is not fair

### ③ Perform Experiments

10  $\rightarrow$  Null Hypothesis is Rejected

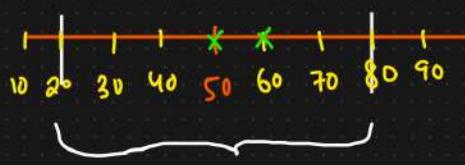
$\rightarrow$  Alternative Hypothesis is Accepted

10 times	75
100 times	60 40
70 30	80 20
50 times Head	<u>Fair</u>
60 times Head	

$$C.I = [20 - 80]$$



Coin is fair



C.I  $\Rightarrow$  Confidence Interval

70 times  $\Rightarrow$  Domain Expert



Confidence Interval

$\hookrightarrow$  We fail to Reject the Null Hypothesis [within C.I]

$\Rightarrow$  Conclusions

$\hookrightarrow$  We Reject the Null Hypothesis [outside C.I]

② Person is Criminal or not {Murder Case}

① Null Hypothesis : Person is not Criminal

② Alternative Hypothesis : Person is Criminal

③ Evidence / Proof : DNA, finger print, weapons, eye witness, foot age



Judge  $\Rightarrow$

Conclusions

Vaccines  $\Rightarrow$  Medical  $\Rightarrow$  critical

$\Rightarrow$  Demand Experiment

Confidence Interval : (CI)

=

$[95\%]$

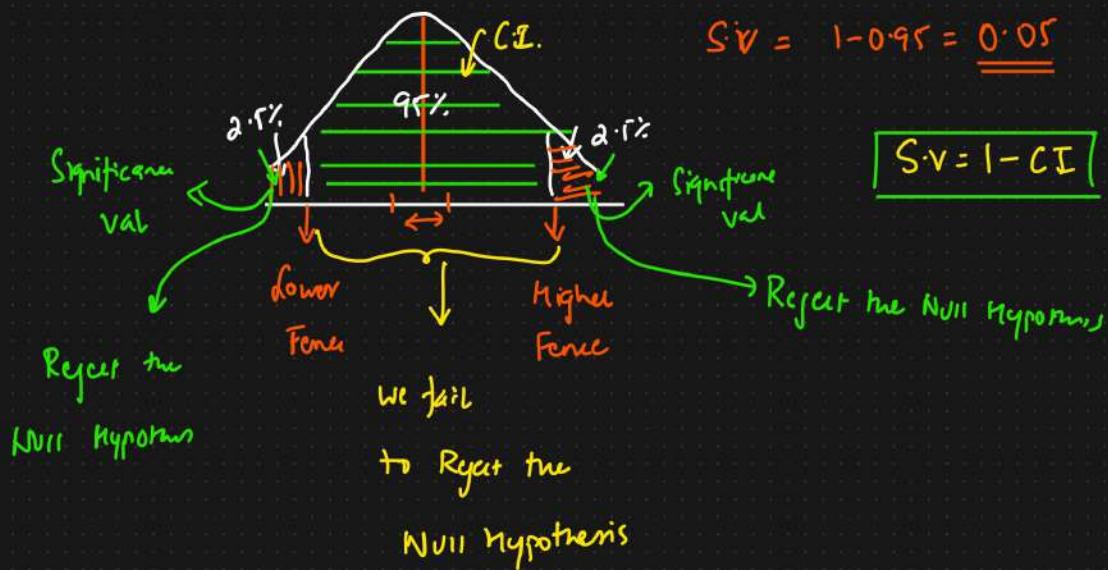
Significance Value

$$C.I = 95\%$$

$$S.V = 1 - C.I$$

$$S.V = 1 - 0.95 = 0.05$$

$$S.V = 1 - C.I$$



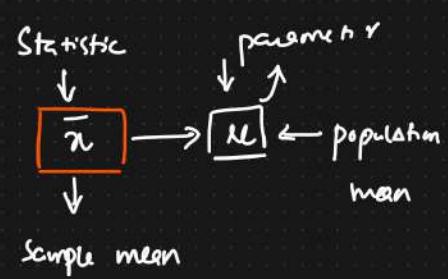
Point Estimate : The value of any statistic that estimates the value of a parameter is called Point Estimate

Point Estimate

$$\uparrow$$

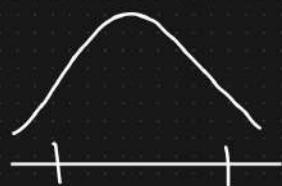
$$\boxed{\bar{x}} \rightarrow \mu$$

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



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

Lower C.I :- Point Estimate - Margin of Error



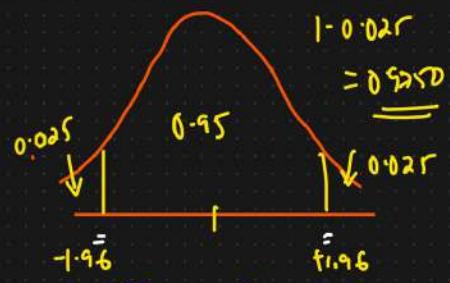
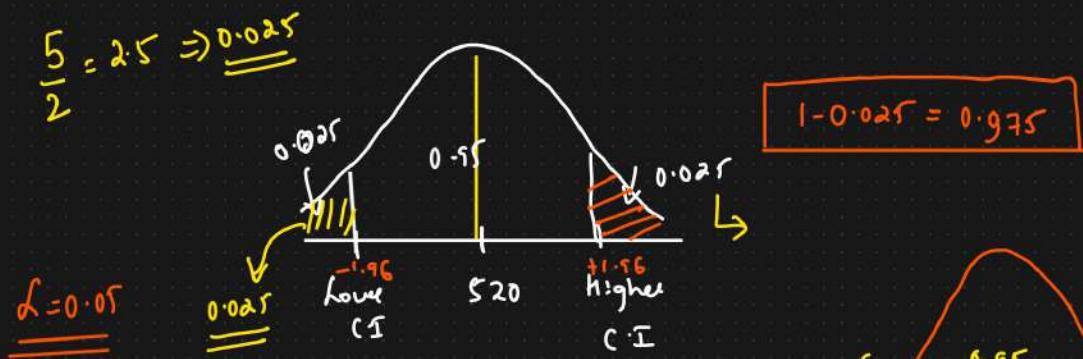
Higher C.I :- Point Estimate + Margin of Error

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

population Sd  
d = Significance value.

- Q) On the quant test of CAT Exam, a sample of 25 test takers has a mean of 520 with a population standard deviation of 100. Construct a 95% C.I about the mean?

Ans)  $n=25 \quad \bar{x}=520 \quad \sigma=100 \quad C.I = 95\% \quad S.V = 1-C.I = 0.05$



Lower C.I = Point Estimate - Margin of Error

$$= 520 - Z_{0.05/2} \frac{\sigma}{\sqrt{n}}$$

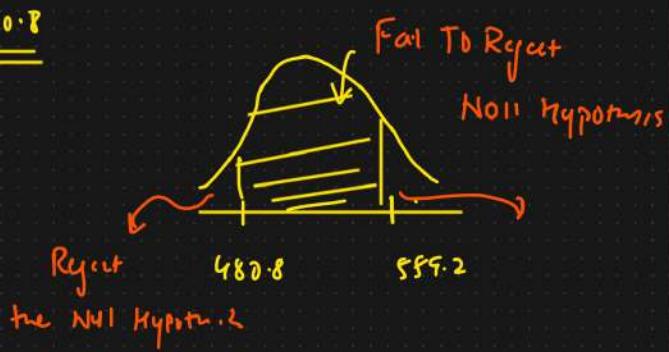
$$= 520 - Z_{0.025} \frac{100}{\sqrt{25}}$$

$$Z_{0.05/2} \Rightarrow \boxed{Z_{0.025}}$$

$$S.V = 1 - 0.95 = 0.05$$

$$= 520 - 1.96 \times 20 = \underline{\underline{480.8}}$$

$$\text{Higher CI} = 520 + 1.96 \times 20 = \underline{\underline{559.2}}$$



①  $\bar{x} = 480$      $S = 85$      $n = 25$     C.I = 90%    Significance

$$= 1 - 0.90 = \boxed{0.10}$$



$$\text{Lower CI} = 480 - Z_{0.10/2} \left[ \frac{85}{5} \right]$$

$$\text{Higher CI} = 480 + 1.64 [17]$$

$$= 480 - Z_{0.05} \left[ \frac{85}{5} \right]$$

$$= 480 + 27.8$$

$$= 480 - 1.64 [17]$$

$$= 507.8$$

$$= 480 - 27.8 = 452.12$$

$$[452.12 \leftrightarrow 507.8]$$

② On the Quant test of CAT exam, a sample of 25 test takers has a mean of 520, with a sample standard deviation of 80.

Construct 95% CI about the mean?

$$\text{Ans) } \bar{x} = 520 \quad S = 80 \quad C.I = 95\% \quad S.V = 1 - 0.95 = 0.05 \quad n = 25$$

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

t test



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

9 9 9

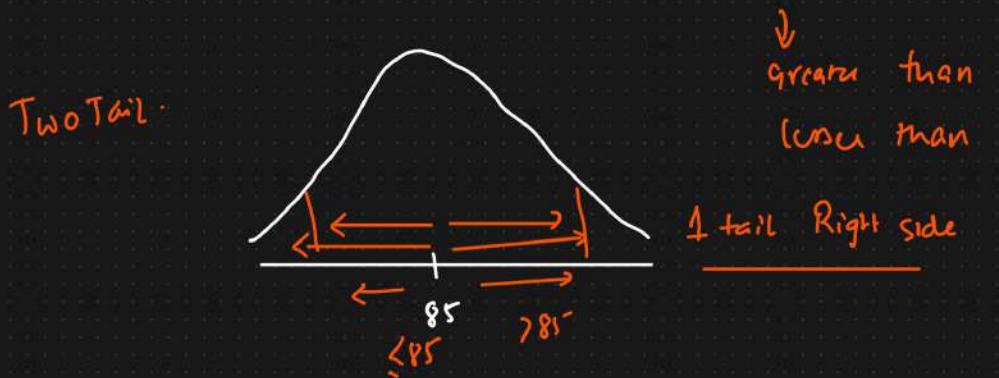
$$\text{Lower C.I} = 520 - t_{0.05/2} \left( \frac{\frac{16}{80}}{81} \right) = 520 - 2.064 \times 16$$

$$\text{Lower C.I} = 486.976$$

$$\text{Higher C.I} = 553.024$$

### ① 1 Tail and 2 Tail Test

- ① Colleges in Town A has 85% placement rate. A new college was recently opened and it was found that a sample of 150 students had a placement rate of  $\underline{88\%}$  with a standard deviation of  $\underline{4\%}$ . Does this college has a different placement rate with 95% C.I?



- ① Z-test }  
② t-test }

④ A factory has a machine that fills 80ml of Baby medicine in a bottle. An employee believes the average amount of baby medicine is not 80ml. Using 40 samples, he measures the average amount dispersed by the machine to be 78ml with a standard deviation of 2.5.

(a) State Null & Alternate Hypotheses

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

Step 1

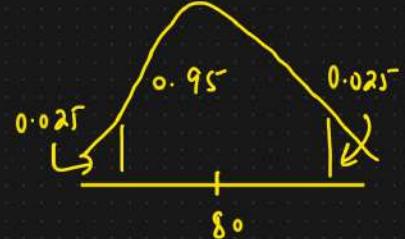
Ans) Null Hypothesis  $\mu = 80$   $\rightarrow$

$$M=80 \text{ ml} \quad n=40 \quad \bar{x}=78 \quad S=2.5$$

Alternate Hypothesis  $\mu \neq 80$   $\rightarrow$

Step 2  $\therefore C.I = 0.95$

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



Step 3  $\therefore$

①  $n > 30$  or population sd  $\} \rightarrow Z \text{ test}$

$$n = 40$$

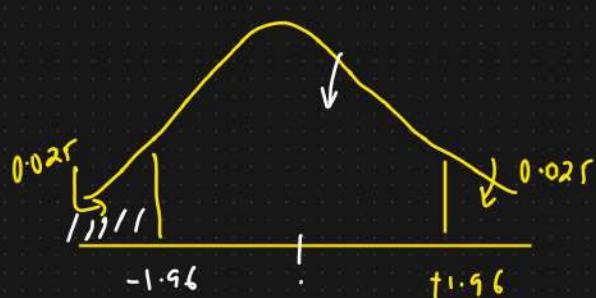
$$S = 2.5$$

②  $n < 30$  and sample std  $\} \rightarrow t \text{ test}$

Z test

Let's perform the Experiment

Decision Boundary



$$1 - 0.025 = 0.975$$

④ Calculate test statistics (z-test)

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

⑤ Conclusions

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

Reject the Null Hypothesis  $\left\{ \begin{array}{l} \text{There is some fault in the} \\ \text{machine.} \end{array} \right.$

⑥ A complaint was registered, the boys in a Government School are underfed. Average weight of the 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 CI: 95% Check if it is True or False.

Ans) Conditions for z-test

$$n=25 \quad \mu=32 \quad \sigma=9 \quad \bar{x}=29.5$$

① We know the population sd. OR

② We do not know the population sd but our sample is large  $n > 30$

## Conditions For T test

- ① We do not know the population std.
- ② Our sample size is small  $n < 30$
- ③ Sample std is given.

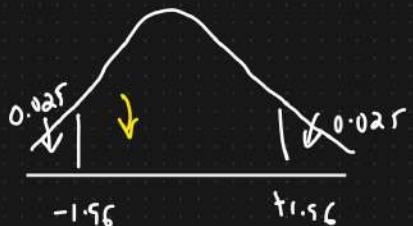
### Step 1

①  $H_0 : \mu = 32$

$H_1 : \mu \neq 32$

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

### ③ Z test



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

Conclusion  $= -1.39 > -1.96$  Accept the Null hypothesis 95% C.I

We fail to Reject the Null hypothesis

The Boys are fed well.

① A factory manufactures cars with a warranty of 5 years <sup>or more</sup> on the engine and transmission. An engineer believes that the engine or transmission will malfunction in less than 5 years. He tests a sample of 40 cars and finds the average time to be 4.8 years with a standard deviation of 0.50. ① State the null & alternate hypotheses

② At a 2% significance level, is there enough evidence to support the idea that the warranty should be revived?  $= -$

Step 1:  $H_0 : \mu \geq 5$

$H_1 : \mu < 5$

Step 2:  $\alpha = 0.02$   $C.I = 0.98$

② In the population the average IQ is 100 with a standard deviation of 15. A team of scientists wants to test a new medication to see if it has a tre or -ve effect, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect intelligence?  $\{ 95\% \}$

## Inferential Statistics

- ① Hypothesis Testing
- ② p-value
- ③ Confidence Interval
- ④ Significance Value

$\chi^2$  test  
t test  
Chi square test  
Anova test (F-test)

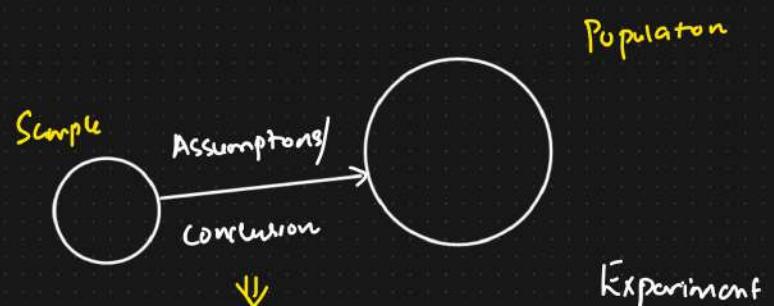
### 3 Distributions

- ① Bernoulli's
- ② Binomial
- ③ Power Law

TRANSFORMATION

## Inferential Stat

### Steps of hypothesis Testing



### Experiment

① Null Hypothesis: Coin is fair  $\Rightarrow$  Accepted  $\rightarrow$  [Coin is fair or not]

② Alternative Hypothesis: Coin is not fair

### ③ Perform Experiments

10  $\rightarrow$  Null Hypothesis is Rejected

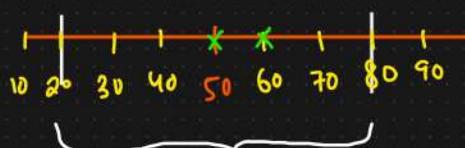
$\rightarrow$  Alternative Hypothesis is Accepted

10 times	75
100 times	60 40
70 30	80 20
50 times Head	<u>Fair</u>
60 times Head	

$$C.I = [20 - 80]$$



Coin is fair



C.I  $\Rightarrow$  Confidence Interval

70 times  $\Rightarrow$  Domain Expert



Confidence Interval

$\hookrightarrow$  We fail to Reject the Null Hypothesis [within C.I]

$\Rightarrow$  Conclusions

$\hookrightarrow$  We Reject the Null Hypothesis [outside C.I]

② Person is Criminal or not {Murder Case}

① Null Hypothesis : Person is not Criminal

② Alternative Hypothesis : Person is Criminal

③ Evidence / Proof : DNA, finger print, weapons, eye witness, foot age



Judge  $\Rightarrow$

Conclusions

Vaccines  $\Rightarrow$  Medical  $\Rightarrow$  critical

$\Rightarrow$  Demand Experiment

Confidence Interval : (CI)

=

$[95\%]$

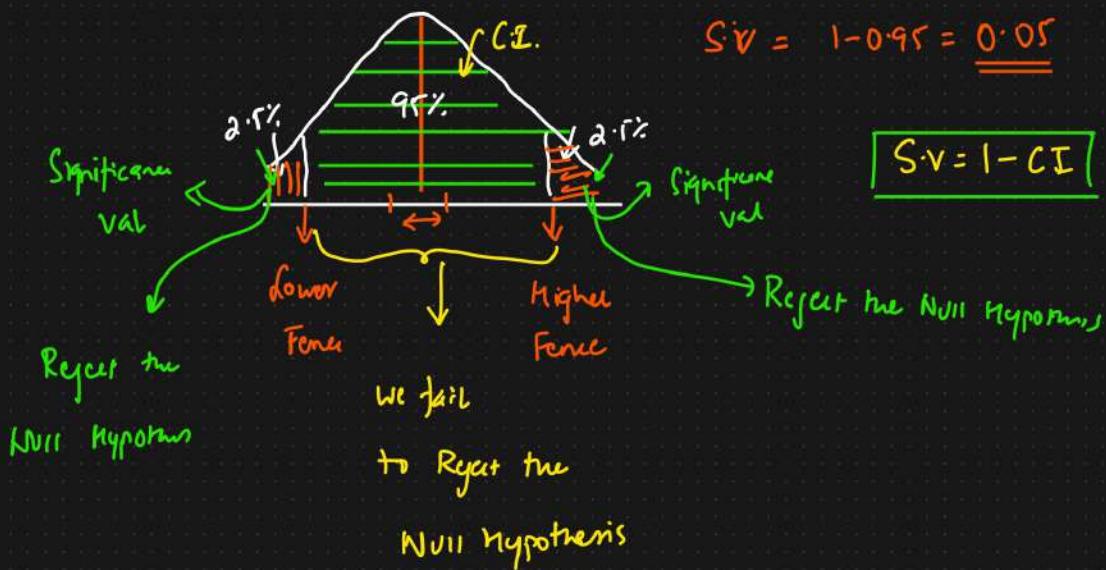
Significance Value

$$C.I = 95\%$$

$$S.V = 1 - C.I$$

$$S.V = 1 - 0.95 = 0.05$$

$$S.V = 1 - C.I$$

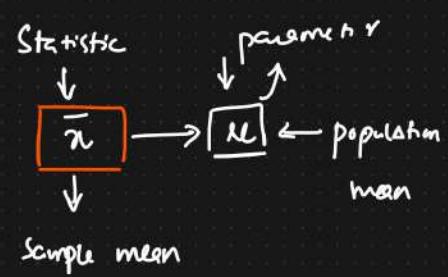


Point Estimate : The value of any statistic that estimates the value of a parameter is called Point Estimate

Point Estimate

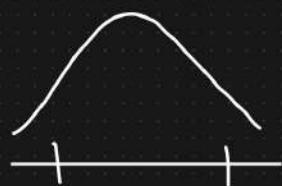
$$\uparrow$$
  
$$\boxed{\bar{x}} \rightarrow \mu$$

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



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

Lower C.I :- Point Estimate - Margin of Error



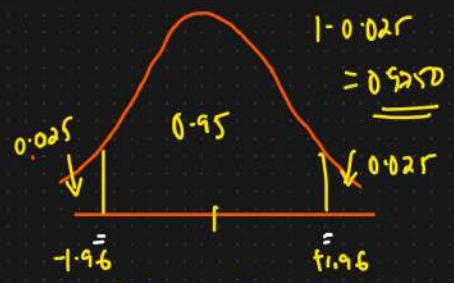
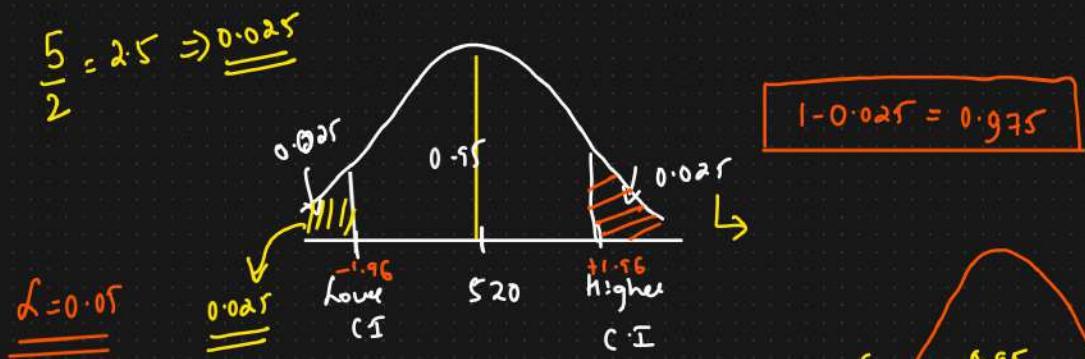
Higher C.I :- Point Estimate + Margin of Error

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

population Sd  
d = Significance value.

- Q) On the quant test of CAT Exam, a sample of 25 test takers has a mean of 520 with a population standard deviation of 100. Construct a 95% C.I about the mean?

Ans)  $n=25 \quad \bar{x}=520 \quad \sigma=100 \quad C.I = 95\% \quad S.V = 1-C.I = 0.05$



Lower C.I = Point Estimate - Margin of Error

$$= 520 - Z_{0.05/2} \frac{\sigma}{\sqrt{n}}$$

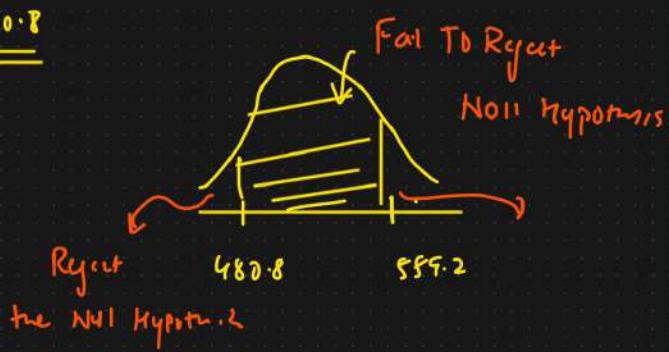
$$= 520 - Z_{0.025} \frac{100}{\sqrt{25}}$$

$$Z_{0.05/2} \Rightarrow \boxed{Z_{0.025}}$$

$$S.V = 1 - 0.95 = 0.05$$

$$= 520 - 1.96 \times 20 = \underline{\underline{480.8}}$$

$$\text{Higher CI} = 520 + 1.96 \times 20 = \underline{\underline{559.2}}$$



①  $\bar{x} = 480$      $S = 85$      $n = 25$     C.I = 90%    Significance

$$= 1 - 0.90 = \boxed{0.10}$$



$$\text{Lower CI} = 480 - Z_{0.10/2} \left[ \frac{85}{5} \right]$$

$$\text{Higher CI} = 480 + 1.64 [17]$$

$$= 480 - Z_{0.05} \left[ \frac{85}{5} \right]$$

$$= 480 + 27.8$$

$$= 480 - 1.64 [17]$$

$$= 507.8$$

$$= 480 - 27.8 = 452.12$$

$$[452.12 \leftrightarrow 507.8]$$

② On the Quant test of CAT exam, a sample of 25 test takers has a mean of 520, with a sample standard deviation of 80.

Construct 95% CI about the mean?

$$\text{Ans) } \bar{x} = 520 \quad S = 80 \quad C.I = 95\% \quad S.V = 1 - 0.95 = 0.05 \quad n = 25$$

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

t test



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

9 9 9

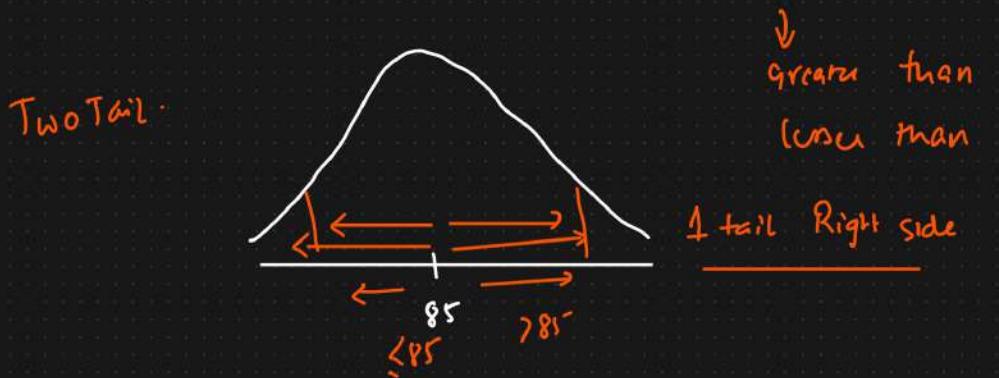
$$\text{Lower C.I} = 520 - t_{0.05/2} \left( \frac{\frac{16}{80}}{81} \right) = 520 - 2.064 \times 16$$

$$\text{Lower C.I} = 486.976$$

$$\text{Higher C.I} = 553.024$$

### ① 1 Tail and 2 Tail Test

- ① Colleges in Town A has 85% placement rate. A new college was recently opened and it was found that a sample of 150 students had a placement rate of  $\underline{88\%}$  with a standard deviation of  $\underline{4\%}$ . Does this college has a different placement rate with 95% C.I?



- ① Z-test }  
② t-test }

④ A factory has a machine that fills 80ml of Baby medicine in a bottle. An employee believes the average amount of baby medicine is not 80ml. Using 40 samples, he measures the average amount dispersed by the machine to be 78ml with a standard deviation of 2.5.

(a) State Null & Alternate Hypotheses

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

Step 1

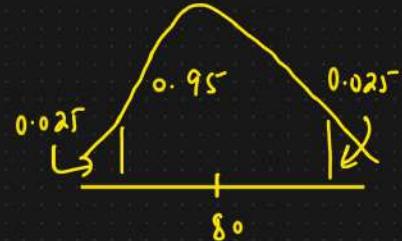
Ans) Null Hypothesis  $\mu = 80$   $\rightarrow$

$$M=80 \text{ ml} \quad n=40 \quad \bar{x}=78 \quad S=2.5$$

Alternate Hypothesis  $\mu \neq 80$   $\rightarrow$

Step 2  $\therefore C.I = 0.95$

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



Step 3  $\therefore$

①  $n > 30$  or population std  $\} \rightarrow Z \text{ test}$

$$n = 40$$

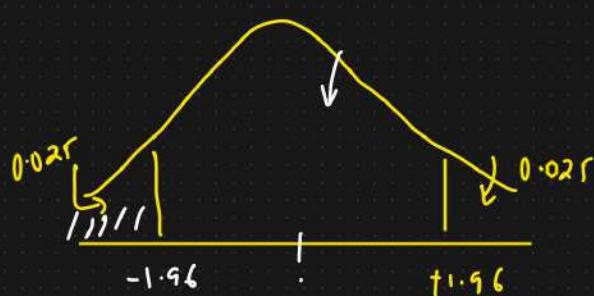
$$S = 2.5$$

②  $n < 30$  and sample std  $\} \rightarrow t \text{ test}$

Z test

Let's perform the Experiment

Decision Boundary



$$1 - 0.025 = 0.975$$

④ Calculate test statistics (z-test)

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

⑤ Conclusions

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

Reject the Null Hypothesis  $\left\{ \begin{array}{l} \text{There is some fault in the} \\ \text{machine.} \end{array} \right.$

⑥ A complaint was registered, the boys in a Government School are underfed. Average weight of the 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 CI: 95% Check if it is True or False.

Ans) Conditions for z-test

$$n=25 \quad \mu=32 \quad \sigma=9 \quad \bar{x}=29.5$$

① We know the population sd. OR

② We do not know the population sd but our sample is large  $n > 30$

## Conditions For T test

- ① We do not know the population std.
- ② Our sample size is small  $n < 30$
- ③ Sample std is given.

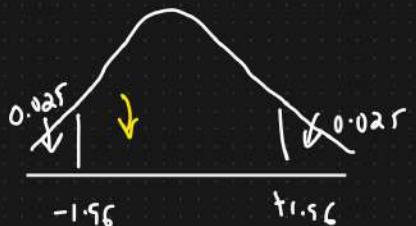
### Step 1

①  $H_0 : \mu = 32$

$H_1 : \mu \neq 32$

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

### 3 Z test



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

Conclusion  $= -1.39 > -1.96$  Accept the Null hypothesis 95% C.I

We fail to Reject the Null Hypothesis

The Boys are fed well.

① A factory manufactures cars with a warranty of 5 years <sup>or more</sup> on the engine and transmission. An engineer believes that the engine or transmission will malfunction in less than 5 years. He tests a sample of 40 cars and finds the average time to be 4.8 years with a standard deviation of 0.50. ① State the null & alternate hypotheses

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

$$\underline{\underline{Z\text{-score}}} = \underline{\underline{-2.5298}}$$

Step 1:  $H_0 : \mu \geq 5$

$H_1 : \mu < 5$

Step 2:  $\alpha = 0.02$   $C.I = 0.98$

② In the population the average IQ is 100 with a standard deviation of 15. A team of scientists wants to test a new medication to see if it has a tve or -ve effect, or no effect at all.

A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affect intelligence?  $\underline{\underline{\alpha = 0.05}}$

$\underline{\underline{Z\text{-score}}} = \underline{\underline{14.96}}$

① A factory manufactures cars with a warranty of 5 years <sup>or more</sup> on the engine and transmission. An engineer believes that the engine or transmission will malfunction in less than 5 years. He tests a sample of 40 cars and finds the average time to be 4.8 years with a standard deviation of 0.50. ① State the null & alternate hypothesis

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

$$\text{Ans) } n = 40 \quad \bar{x} = 4.8 \text{ years} \quad s = 0.50$$

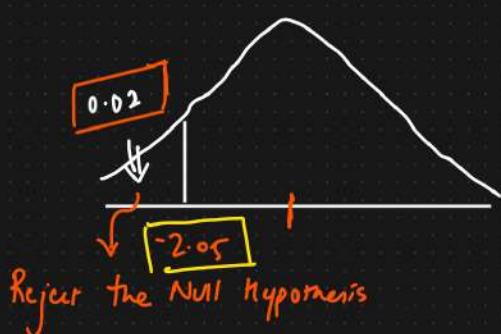
Step 1

$$H_0: \mu \geq 5 \quad \{\text{NULL Hypothesis}\}$$

$$H_1: \mu < 5 \quad \{\text{Alternate Hypothesis}\}$$

Step 2:  $\alpha = 0.02 \quad C.I = 1 - 0.02 = 0.98 = 98\%$

Step 3:



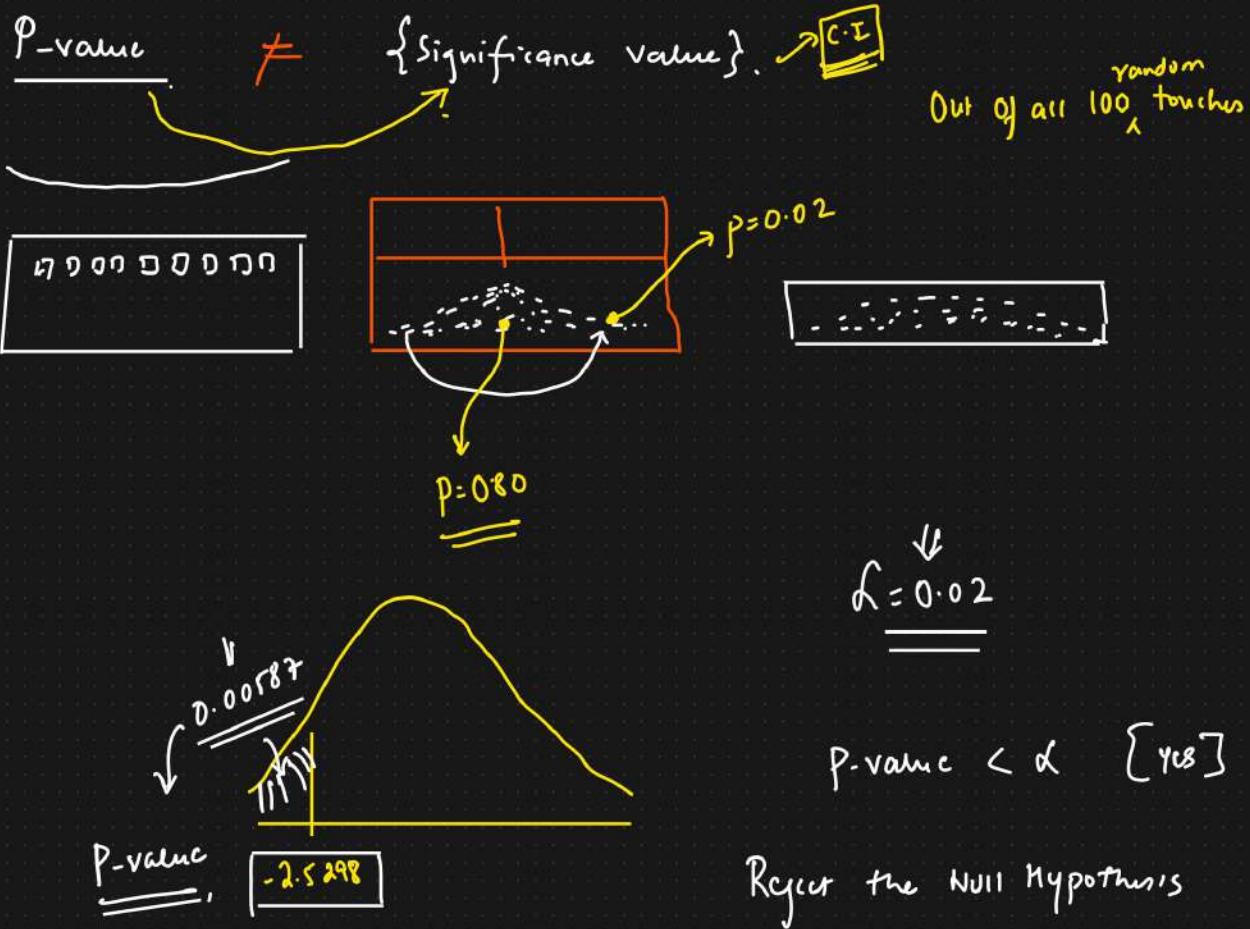
Reject the Null Hypothesis

Step 4:

$$Z\text{-score} = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{4.8 - 5}{0.50/\sqrt{40}} = -2.5298$$

Conclusion:  $-2.5298 < -2.05$  Reject the Null Hypothesis

Warranty needs to be revised.



- \* 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
- Null & Alternative hypotheses
  - 95%. Is there enough evidence to discard the null hypothesis?

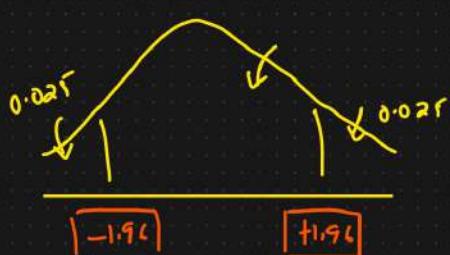
Ans)  $\bar{x} = 169.5$      $S = 3.9$      $n = 36$      $\mu = 168$      $(C.I = 0.95)$

Step 1  
 $H_0: \mu = 168$

$H_1: \mu \neq 168$

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

Step 3

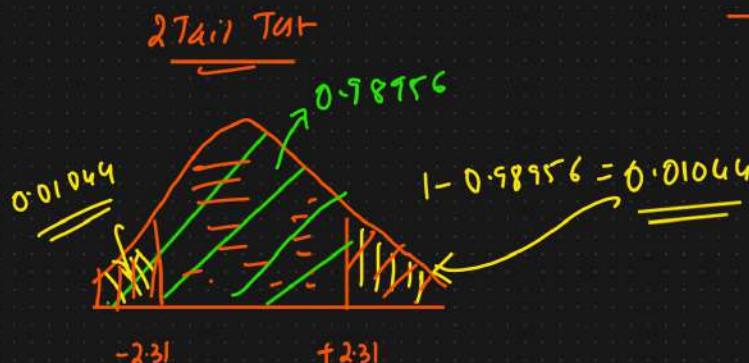


$$\text{Step 4 : } t\text{-score} = \frac{169.5 - 168}{\frac{3.4}{\sqrt{36}}} = \sqrt{2.31} \Rightarrow 2.31 \neq 0$$

$2.31 > 1.96$  {Reject the Null Hypothesis}

2 tail test

P-value



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

$0.02088 < 0.05$  {Reject the Null Hypothesis}.

\* 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 Alternative Hypothesis?

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

Ans) ①  $H_0 : \mu \geq 2$

② Step 2

$n < 30$

Sample Standard deviation

$$H_1 : \mu < 2$$

$$C.I = 0.99 \quad \alpha = 0.01$$

③ Step 3



$$\begin{aligned} \text{Degree of freedom} &= n - 1 \\ &= 10 - 1 = 9 \end{aligned}$$

④ Calculate t test statistics

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

⑤  $-4.216 < -2.82$  { Reject the Null Hypothesis }.

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

⑥ Z test with proportions

⑦ 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 finds that 130 responded Yes - owning a cell phone?

⑧ State Null And Alternative Hypothesis?

⑨ At a 95% CI, is there enough evidence to reject the Null Hypothesis?

Ans)

Step 1

Null Hypothesis :  $P_0 = 0.70$

Alternative Hypothesis :  $P_0 \neq 0.70$

$$q_0 = 1 - P_0 = 0.30$$

$$n = 200$$

$$x = 130$$

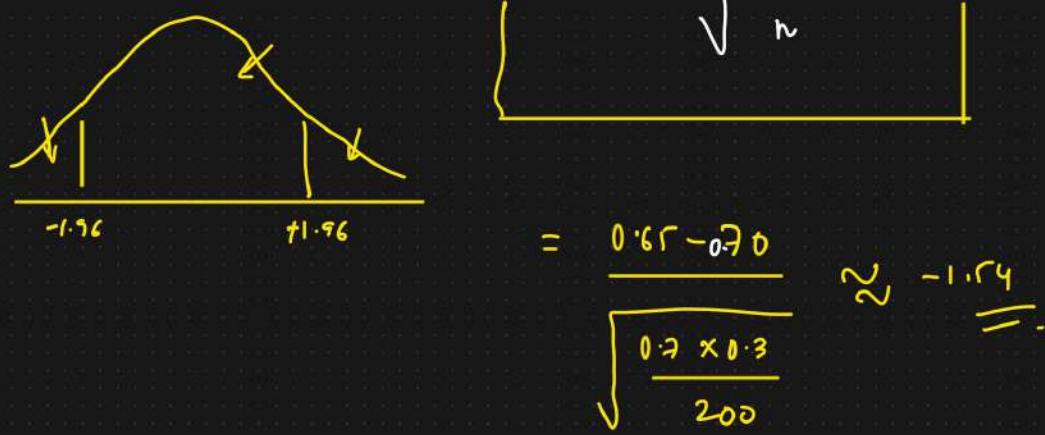
$$\hat{P} = \frac{x}{n} = \underline{\underline{0.65}}$$

Step 2 :  $\alpha = 0.05$

Step 4 : Z test with proportion

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

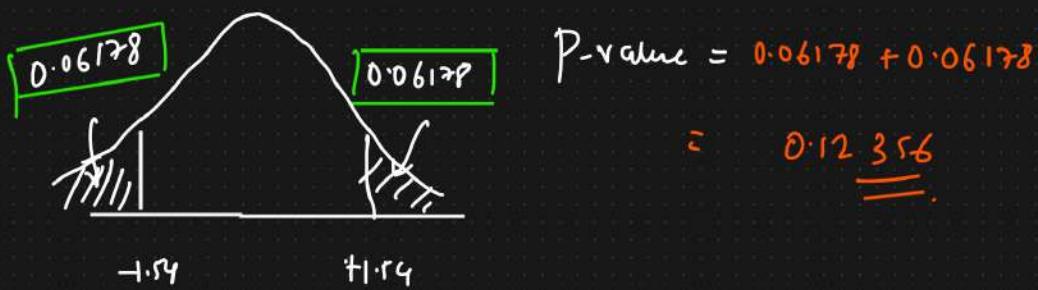
Step 3 :



Conclusion

$-1.54 > -1.96$  Fail to Reject the Null Hypothesis

Ratio



Pvalue > Significance value Fail To Reject Null Hypothesis

- ④ 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 & Alternative Hypothesis

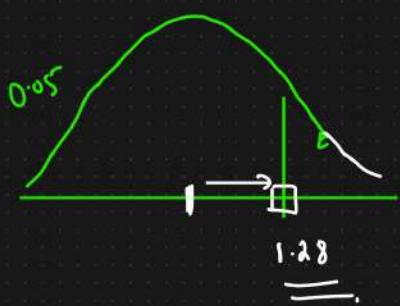
(b) At 10% significance level, is there enough evidence to support the idea that vehicle ownership in City ABC is 60% or less?

$$H_0 : p_0 \leq 0.60$$

$$H_1 : p_0 > 0.60$$

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

$$q_0 = 0.40$$



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

$$= \frac{0.08}{0.0309} = 2.599$$

Reject the Null Hypothesis

## ④ Chi Square test

ORDINAL DATA

① Chi Square test claims about population proportions.

NOMINAL DATA

It is a non parametric test that is performed on categorical data.

② In the 2000 US census the age of individuals in a small town found to be the following

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

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

<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)

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

H:  $H_0$

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

Step 1 : Null hypothesis  $H_0$ : The data meets the expected distribution  
 $H_1$ : The data does not meet the " "

Step 2 :  $\alpha = 0.05$  ( $C.I = 95\%$ )

Step 3 : Degree of freedom {category}

$$df = C - 1 = 3 - 1 = \underline{2} \quad \alpha = 0.05$$

$\hookrightarrow$  No. of categories.

Step 4 : Decision Boundary =  $\boxed{15.991}$  {chi square table}

Step 5 : Chi square Test Statistics

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

$$\boxed{\chi^2 = 232.494}$$

## Conclusion

$$\chi^2 > 5.99 \quad \text{Reject } H_0.$$

\* ↴

ANOVA TEST [F Test] → Assignment

↓

3 Distribution

① Binomial

② Bernoulli

③ Power Law of Pareto

④ Poisson