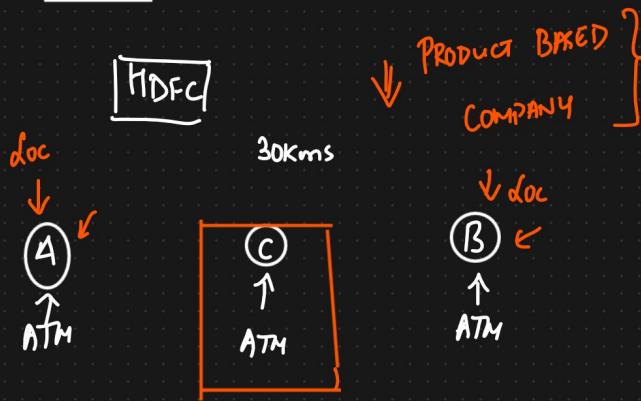


# Statistics

Use case

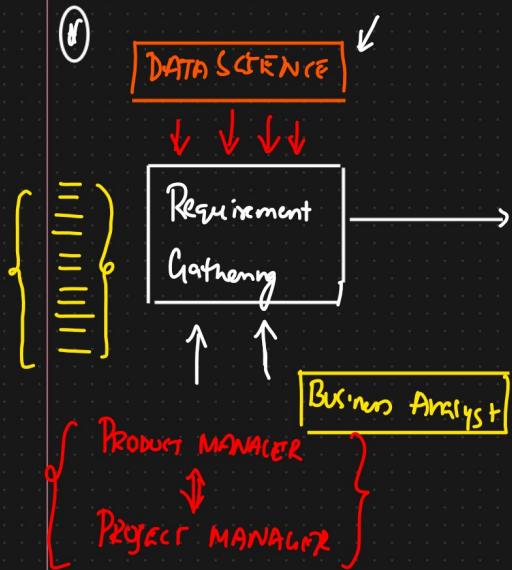


X Statistician → 5 years



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

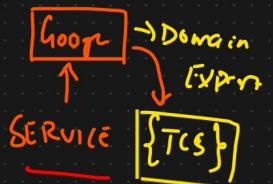
Statistics {Life cycle of DATA Scientific Project}



DATA ANALYST'S TEAM



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



PRODUCT BASED

{ Google }

Domain knowledge

PRODUCT MANAGER  
BA



YouTube, GPAY,  
Google Ads, GMAIL

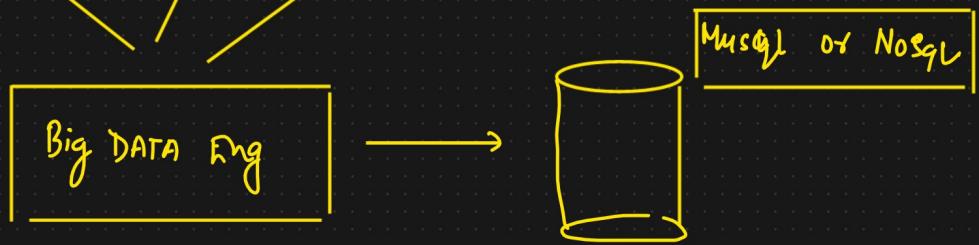
Domain Expertise

Product Manager

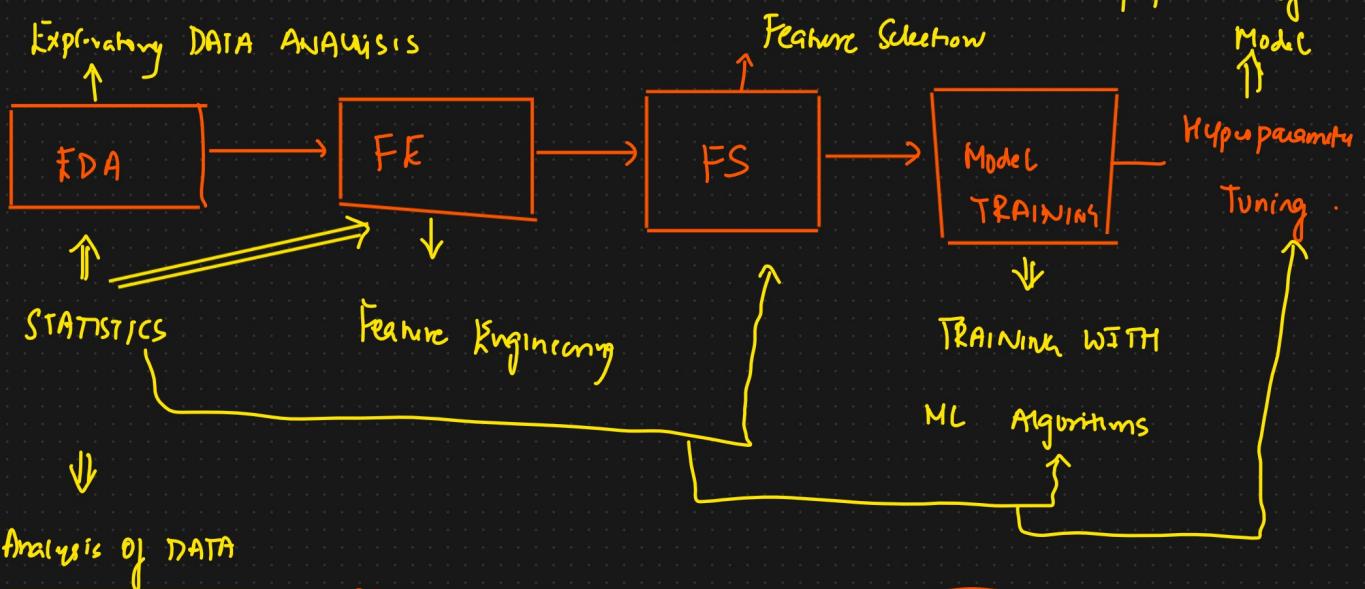
Internal  
DATABASE

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

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

Standard deviation

② Weights of students in classroom

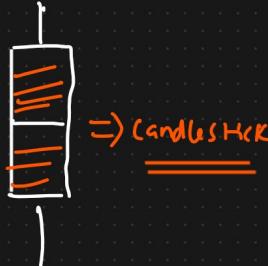
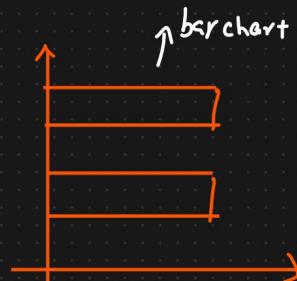
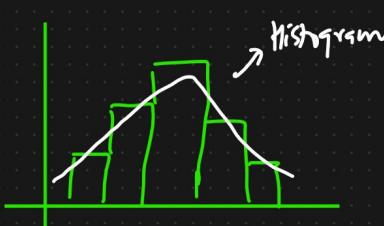
Descriptive Stats [EDA + FF]

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

Hypothesis Testing



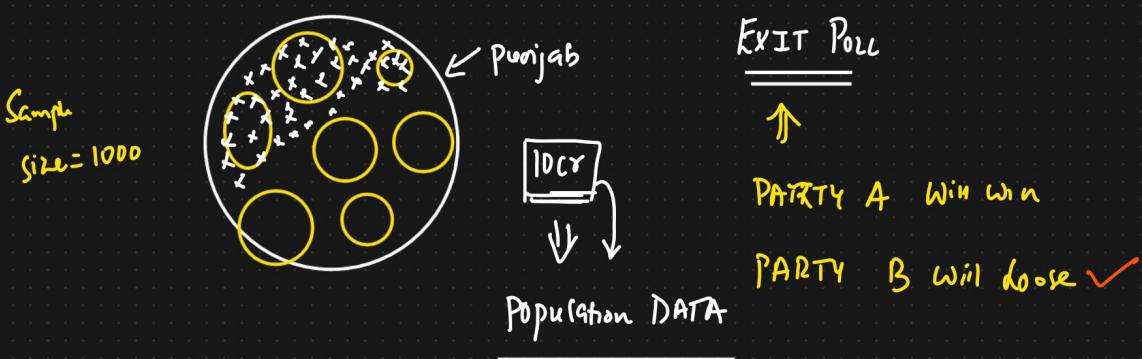
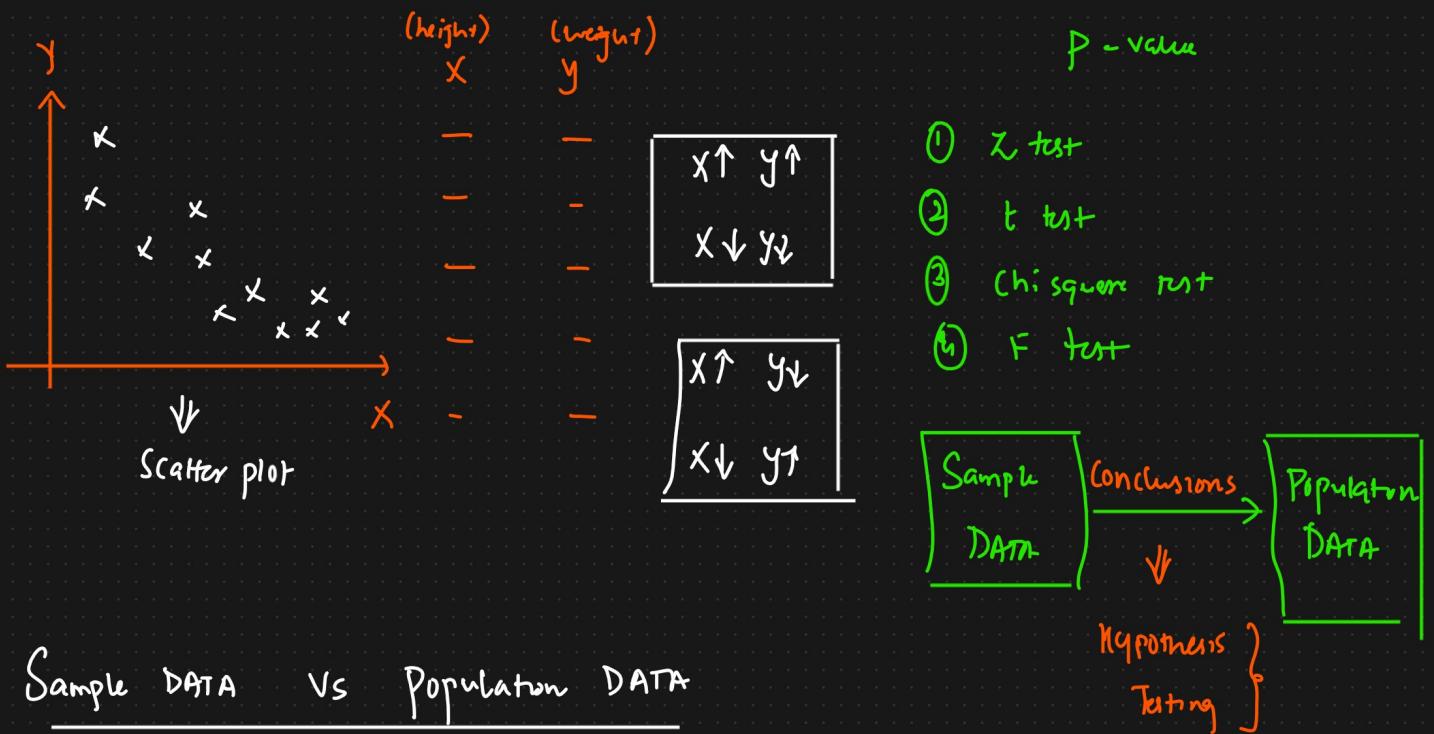
University → 500 people

CLASS A → 60 people

↓  
Sample data → Age → Average age of the entire university

C.I ⇒ Confidence Interval

Hypothesis Testing



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 { - - - - - }

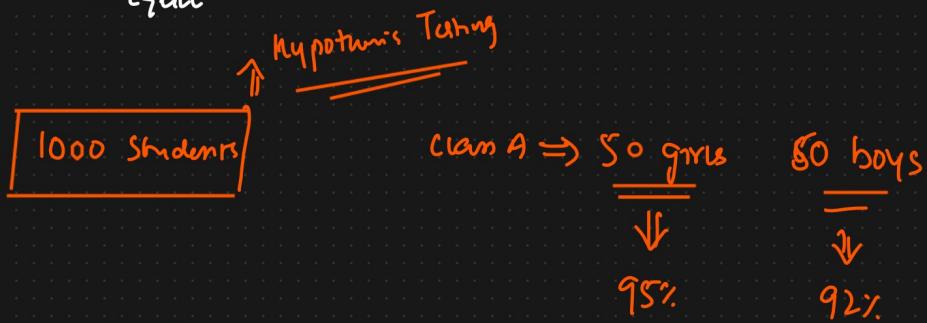
Descriptive Stats : What is the average age of students in the classroom?

Relationship between Age & Gender?

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

less than the average age of the students in the university?  
↓

$\left\{ \begin{array}{l} \text{Greater} \\ \Downarrow \\ \text{Equal} \end{array} \right.$



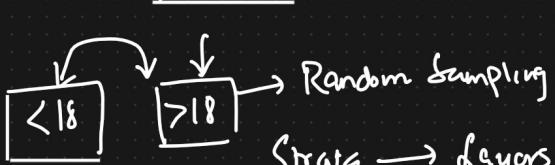
Choose a Sample  
Sampling Techniques

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

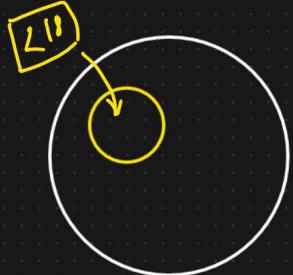
Education

- High School
- Master
- phd

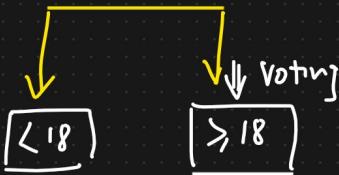
Degree

Blood groups

- 



Population {Exit poll}.



Not Vote

Random Sampling

③ Systamatic Sampling  $\rightarrow \{ \text{AIRPORT} \}$   $n^{\text{th}}$  person

Select every  $n^{\text{th}}$  individual out of  $N$  population (N)  $\left\{ \text{CREDIT CARD} \right\}$

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

inconvenient job for

$\left\{ \text{DATA SCIENCE SURVEY} \rightarrow \text{General AI Survey} \right\}$  a specific



$\Downarrow$   
 $\left\{ \text{Fill the Form} \right\}$

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

$age = 25$        $AgeS = [24, 25, 26, 27, 28, 29] \Rightarrow$  Collection

$age = 100$

Two different types of Variable

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

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 Marbles

Quantitative Variable



Continuous Variable.

Eg: Whole number  $\rightarrow$  fixed

Eg: No. of Bank Accounts

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

$25/X$

Eg: Continuous  $\rightarrow$  Decimal values

Eg: Height, weight, ages, Rainfall

Speed

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

Pincode = fixed

Categorical  
variable  
Marble  
 $\rightarrow$  Mixed  
 $\rightarrow$  Not Mixed }  
variable

## Assessment

Gender ? Categorical

① What kind of variable is Marital Status? Categorical variable

Length River length? Continuous

Movie duration? Continuous

Pincode ? Discrete

IQ ? Discrete

105.75, 90.5,

Pancard

Pincode

Fixed

Categorical

↓  
[FE]

Ans no. of Categories

[AMLPN - - -]

360099  
720058  
560092

}  $\Rightarrow$  It is many?

Categorical

Variables

Continuous



Discrete ←

Continuous

Whole number

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

Pincode = { }

Cities

5

Gender Pincode

M

F

[6]

Categorical

PAN ←

[ ]

[ ]

} Categorical

# 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

Agus = {<sup>0,</sup> 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100}

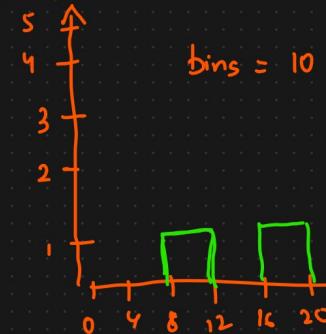
① Sort the Numbers

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

min = 10

max = 40

② Bins → No. of groups



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

③ Bin size → Size of Bins

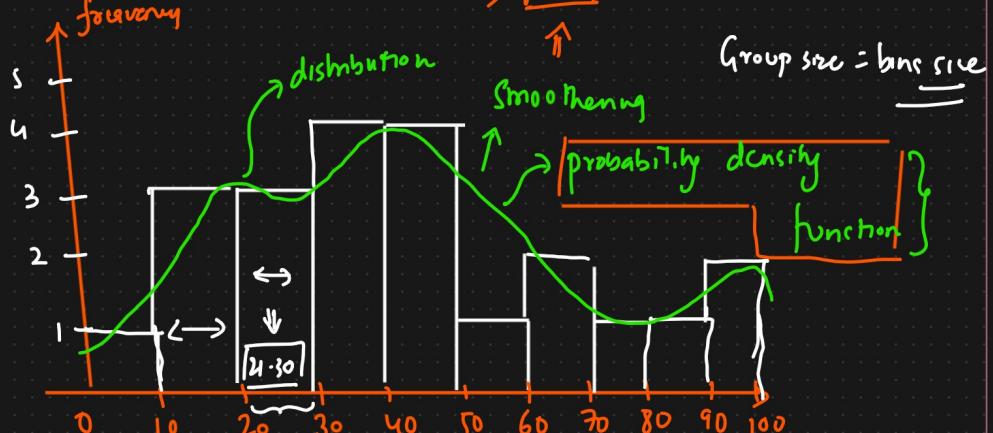
bins = 10

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

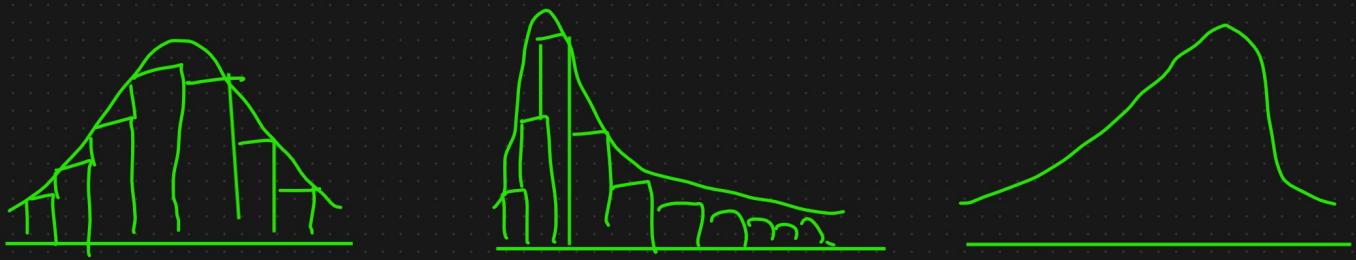
bin size =  $\frac{100}{10} = 10$  //

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

bins =



Agus = {<sup>0,</sup> 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100}



Assignment  
Weight = { $\boxed{30}, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, \boxed{80}, 90, \boxed{95}$ }.

bins = 10

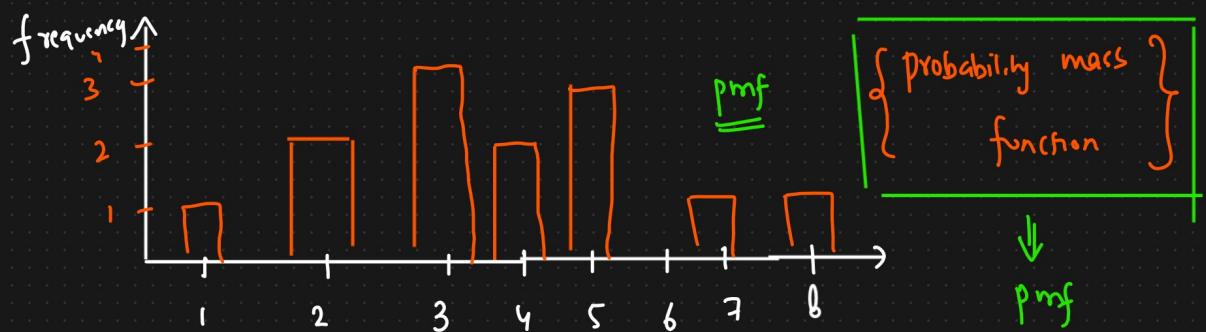
$$\text{bin size} = \frac{95 - 30}{10} = \frac{65}{10} = 6.5$$

{continuous value}

pdf

## ② Discrete

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



pdf : probability density function }  $\rightarrow$  continuous

pmf : probability mass function }  $\rightarrow$  discrete.

## ① 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 ( $n$ )

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

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

$$N = 6$$

$$N > n$$

$$n = 4$$

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

$$\text{Population Age} = \{24, 23, 2, 1, 28, 27\}$$

$$\text{Sample Age} = \{24, 2, 1, 27\}$$

$$\begin{array}{r} 13 \\ 54 \\ \hline 41 \end{array}$$

$$\text{Population mean } (\mu) = \frac{24+23+2+1+28+27}{6}$$

$$\mu = 17.5$$

$$\text{Sample mean } (\bar{x}) = \frac{24+2+1+27}{4}$$

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

$$\bar{x} = 13.5$$

$\boxed{\text{hp-null}}$   $\leftarrow$  Null values

### Practical Application (Feature Engineering)

Age	Salary	Family Size
-	-	-
-	-	-
-	-	-
$\rightarrow \text{NAN}$	-	-
-	-	-
-	$\text{NAN}$	-
-	-	$\text{NAN}$
-	$\text{NAN}$	-
$\rightarrow \text{NAN}$	-	-

$\leftarrow$  loss of Info

$$\boxed{\text{Age} = 29.6}$$

Mean

$$\begin{array}{l} \boxed{\text{NULL}} \downarrow \text{val} = 10/6 \\ 10/4 = \boxed{1, 2, 3, 4} \uparrow \text{NAN} \\ \boxed{\text{mean}} \uparrow \text{NAN} \end{array}$$

Age	Salary
24 ✓	45
28 ✓	50
29 ✓	$\boxed{\text{NAN}}$
$\boxed{\text{NAN}}$ ✗	$\boxed{\text{Salary} = 62}$
31 ✓	60
36 ✓	75
$\boxed{\text{NAN}}$ ✗	$\boxed{85}$
$\boxed{\text{NAN}}$ ✗	$\boxed{\text{NAN}}$

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

## ① Median

$$\{1, 2, 3, 4, 5\} = \{1, 2, 3, 4, 5, \boxed{100}\}$$

$\bar{x} = 3 \longrightarrow \bar{x} = 19.16$

Outlier  
 $\frac{1+2+3+4+5+100}{6} = \frac{115}{6} = 19.16$

### 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}

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

$$\boxed{2, 3}$$

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

## Dataset

Types of flower (Categorical variable).

Flower

Sunflower

Rose

NAN ← ROSE

Rose

Sunflower

Rose

NAN ← Rose

40%

FF  
Biased

Under 19

17, 18, 19, 16, 15, 32 → Outlier

## (F) Measure of Dispersion

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

② Standard deviation ( $\sigma$ ) ←

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

### Variance

Population Variance ( $\sigma^2$ ) { Degree of freedom }

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

Sample Variance ( $s^2$ )

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

$$\begin{aligned} \{1-100\} & \Downarrow \\ & = \Downarrow \text{First one.} = \{1-100\} \end{aligned}$$

Second one.  
↓

↓  
Assignment

$$\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\} \xrightarrow{\text{Variance}} \boxed{\text{Variance}}$$

$$\{1, 2, 3, 4, 50, 60, 70, 100\} \xrightarrow{\text{Variance}} \boxed{\text{Variance}}$$

Variance Given

$$\frac{21}{80} = \frac{101}{7}$$

$$\{1, 2, 3, 4, 5\} \xrightarrow{\text{Mean}} M = 3$$

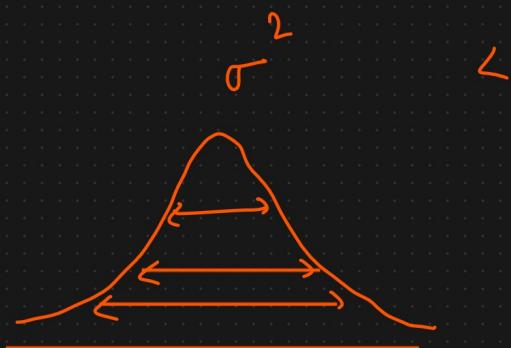
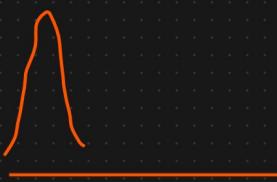
$$\{1, 2, 3, 4, 5, 6, 80\} \xrightarrow{\text{Mean}} M = 14.4$$

$$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5} = \frac{4+1+0+1+4}{5} = \frac{10}{5} = 2$$

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

**Variance ↑↑ Spread ↑↑**

$$\sigma^2 = 719.10$$



④ Standard deviation  $(\sqrt{\sigma^2}) \Rightarrow \boxed{4}$

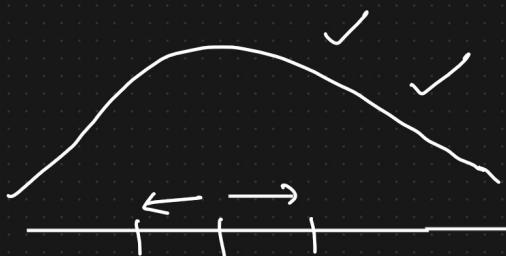
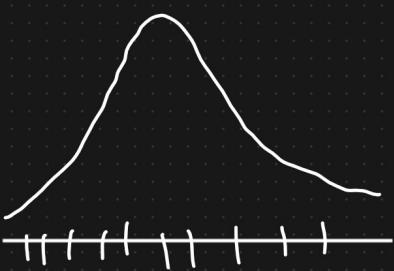
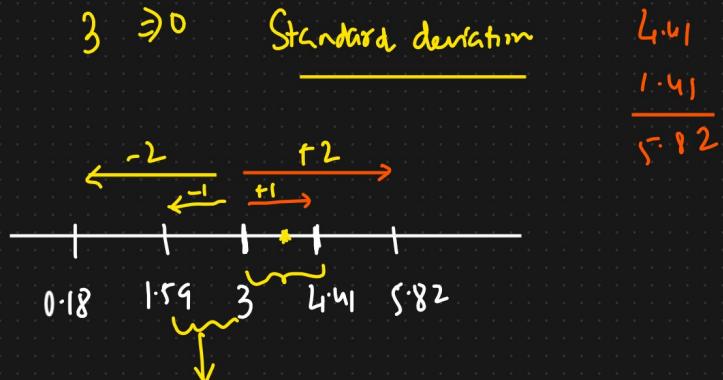
$$\begin{array}{r} 3.00 \\ 1.41 \\ \hline 1.59 \end{array}$$

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

$$M = 3$$

$$\sigma^2 = 2$$

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



## ④ Percentiles And Quartiles

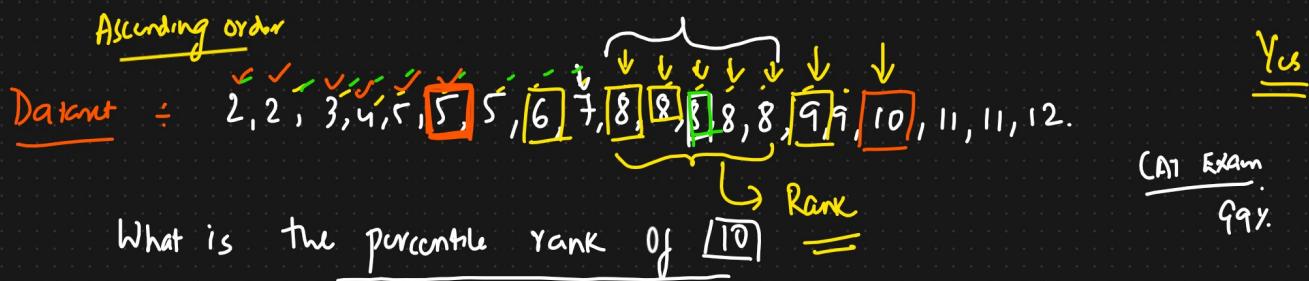
$$\text{Percentile} = \{ 1, 2, 3, 4, 5, 6, 7, 8 \}$$

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

Percentiles : CAT, IAT, JEEL, 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



Next item

= 0.8

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

45 percentile

$$= \frac{14}{20} = 70 \text{ percent.}$$

④ What is the value that exists at 25 percentile

75%

$$\text{Value} = \frac{\text{Percentile}}{100} \times \frac{n+1}{n}$$

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

$$\text{DfP} = 5$$

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

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

$$\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \underline{\underline{12, 14}}\} \quad \text{↓ outlier}$$

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



[Lower Fence  $\longleftrightarrow$  Higher Fence]

$\underline{\underline{\quad}}$

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

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

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

75 25

$$IQR = Q_3 - Q_1$$

$\downarrow = =$

Inter Quartile Range (IQR)

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

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

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

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

$$\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, \underline{\underline{12, 14}}\}.$$

-5

5 Number Summary.

① Minimum = 1 ✓

②  $Q_1 = 3$  ✓

③ Median = 5 ✓

④  $Q_3 = 7.5$  ✓

⑤ Maximum = 9 ✓

Box Plot

$\Downarrow$

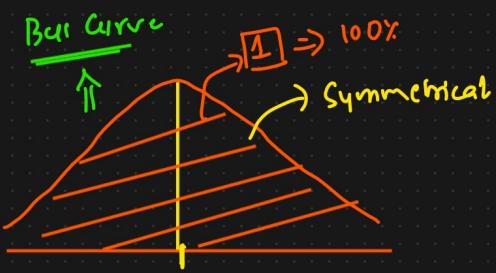


$\Downarrow$

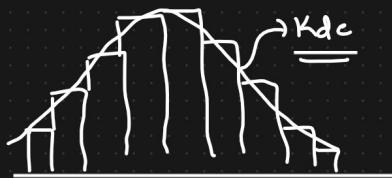
To Treat Outliers

## Day 3 - Stats

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



Kernel density estimator



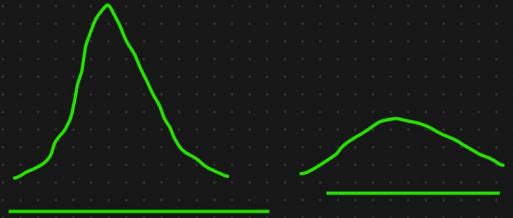
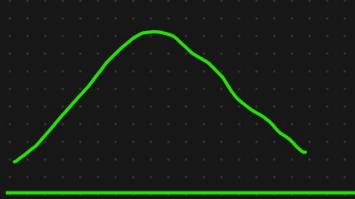
↓      ↓      ↓  
Age, weight, height       $\int$  Distribution  
↑  
 $\Rightarrow$  Doctors

Domain Expertise

[IRIS DATASET] ←  
↓

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

Gaussian Distri



① [Empirical Rule of Normal Distribution]

Empirical Rule

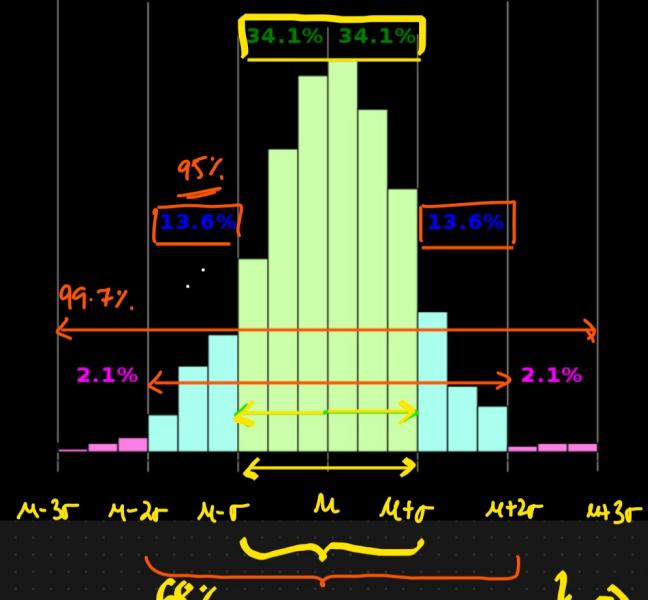


$$68 - 95 - 99.7\%$$



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

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



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

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



$$\mu = 3$$

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

$$= =$$

$$\sigma = 1.41$$

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

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

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

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

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

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

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

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

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



Why?

[Standardization]  $\Rightarrow [\mu=0 \text{ & } \sigma=1]$

<u>(years)</u> <u>Age</u>	<u>(kg)</u> <u>Weight</u>	<u>(cm)</u> <u>Height</u>
$\mu=0$	72	170
$\sigma=1$	38	160
24	84	165
26	92	170
32	87	150
33	83	180
34	80	175
28		
29		

[0-1]

Same Scale



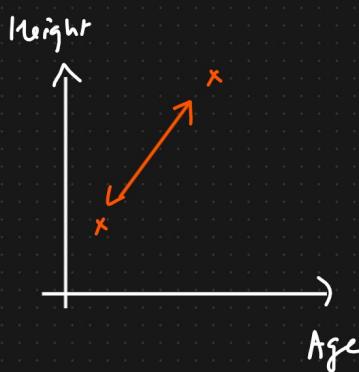
Machine Learning

Maths Equations

Algorithm  $\Rightarrow$  Mathematical Model

Mathematical

Calculation Time ↑↑↑



Feature Scaling

Normalization  $\Rightarrow$

Standardization  $\left\{ Z\text{-Score} \right\}$

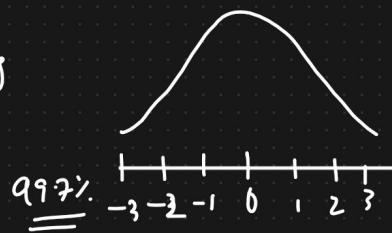
$$\frac{x_i - \mu}{\sigma}$$

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

$$[\mu=0, \sigma=1]$$

$$\begin{bmatrix} 0 & 5 \end{bmatrix} \quad \begin{bmatrix} -3 & 3 \end{bmatrix}$$

$$[0-4]$$



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

① Min Max Scaler

[0-1]

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

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

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

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

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

1
2
3
4
5

$y \downarrow$

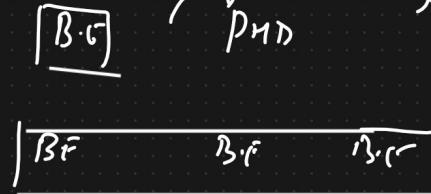
0
0.25
0.5
0.75
1

$y'$

-1.414
-0.302
0
0.702
1.414

Apply ??

Deep learning



B.T  
B.F  
B.C

B.T

B.F

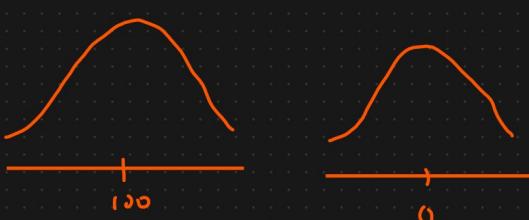
B.C

① Standardization

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

$x \rightarrow$  Normal Distribution ( $\mu, \sigma$ )  
 $\downarrow Z\text{-score}$

$y \rightarrow$  SND ( $\mu=0, \sigma=1$ )



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

Normalization  $[0 - 1]$

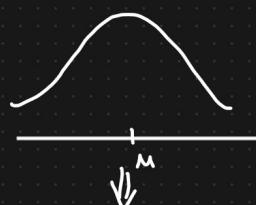


① Min Max Scaler  $\curvearrowright \Leftrightarrow$

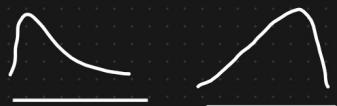
$\Downarrow$   
ML

Standardization

$\Downarrow$   
ML



Min Max Scaler



Min Max Scaler

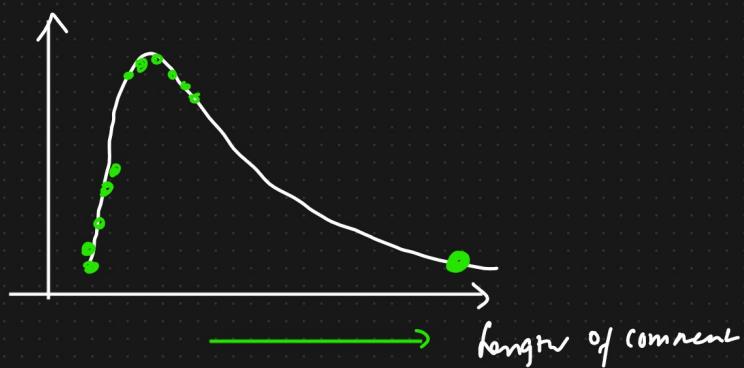
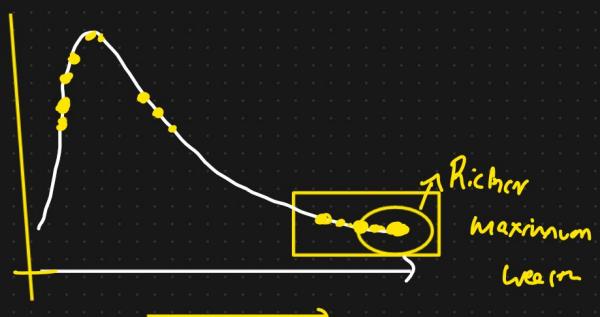
① log Normal Distribution

$\Rightarrow$  Normal/Gaussian  
Distribution.



log Normal Distribution

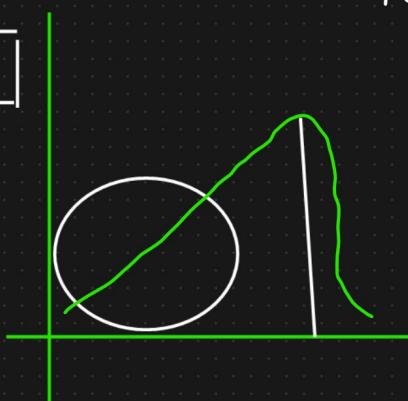
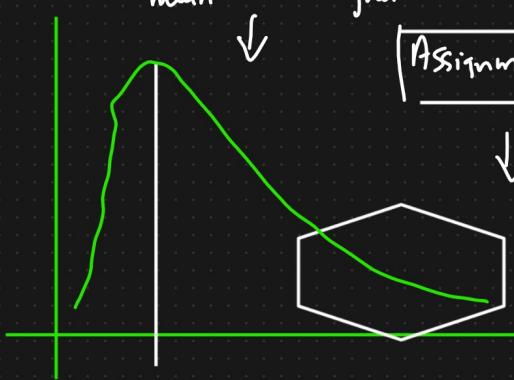
Richer maximum  
mean  
WCAT  
distribution



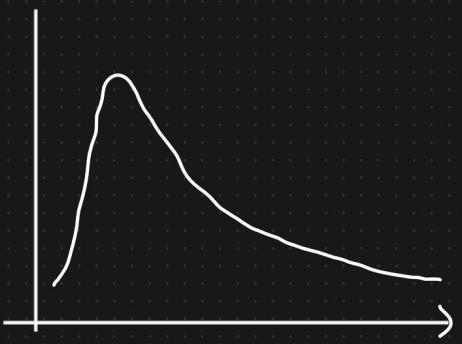
mean will be higher

$\boxed{\text{Assignment}}$

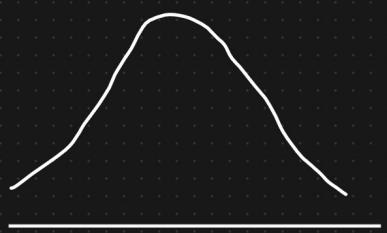
Relation of mean,  
median, mode



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

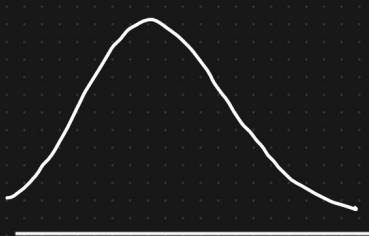


$X \sim \text{log Normal Distribution}$



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

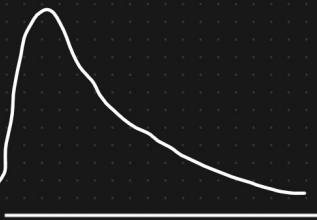
$$Y = \ln(X)$$



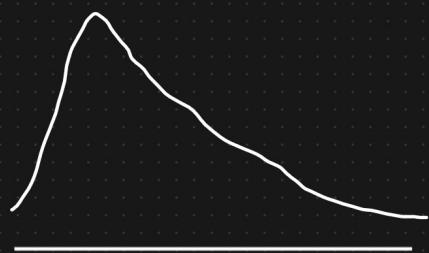
Antilog

$$\Rightarrow \exp(Y)$$

$\Downarrow$



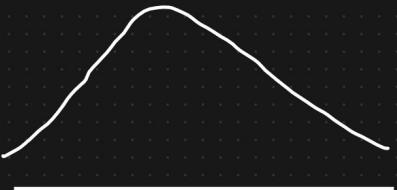
(\*)



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

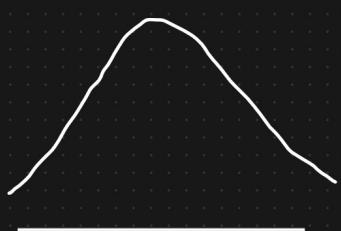
Natural log  
 $\uparrow$   
 $\log_e$

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

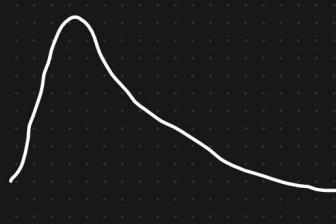


$\Downarrow$

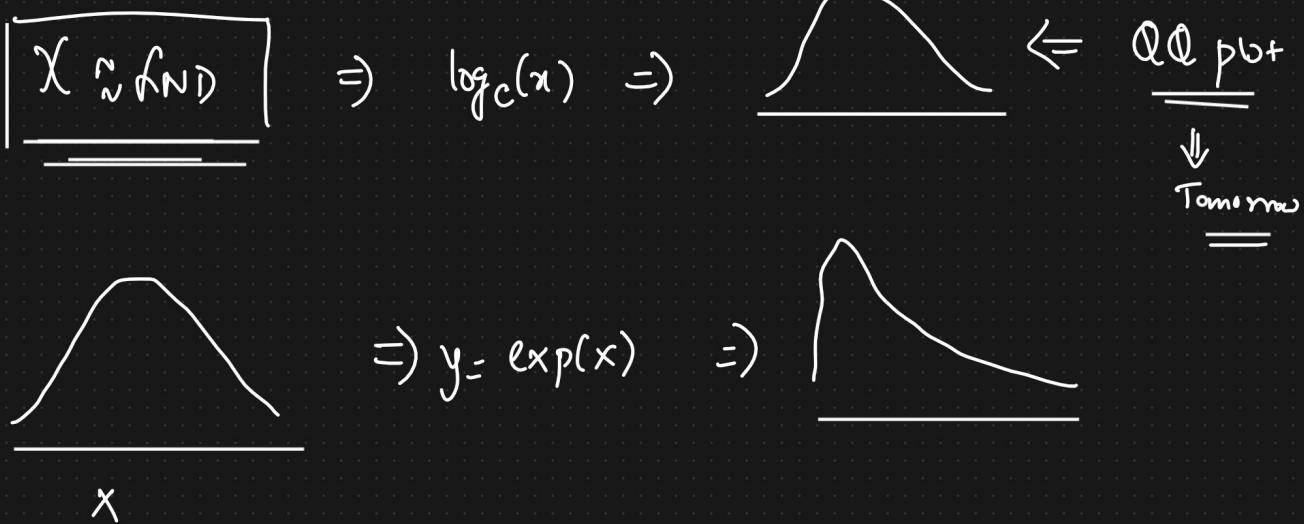
Inverse



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



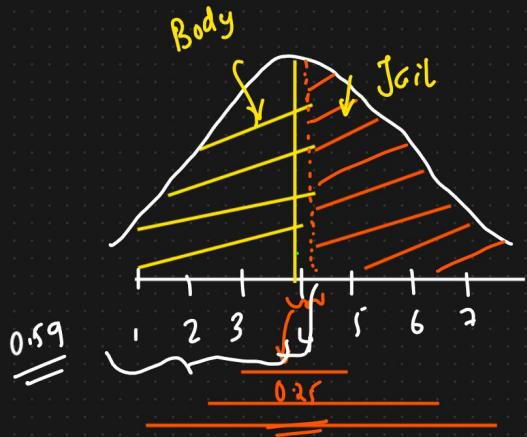
$Y$



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?

$0.59 \Rightarrow 59\%$

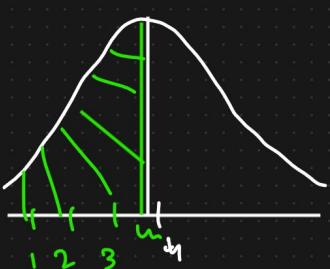
$1 - 0.59 = 0.41 \Rightarrow 41\%$

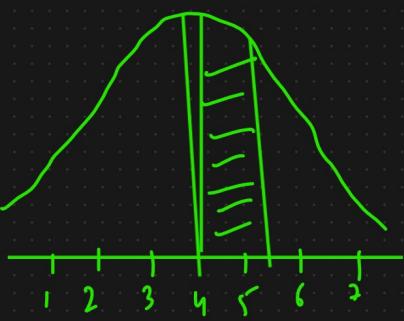
⑦ Z-score =  $\frac{\pi_i - M}{\sigma} = \frac{4.25 - 4}{1} = \boxed{0.25}$

⑧ Z-table (area under the curve)



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





$$4.25 \quad 4.5 \quad 4.75$$

- Q 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).

⑧ Central Limit Theorem }

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

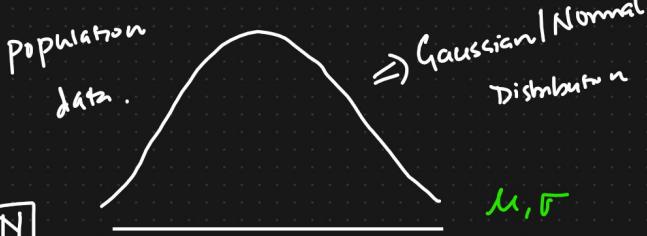
Size of sample

The larger the value the better

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

$$\rightarrow m$$

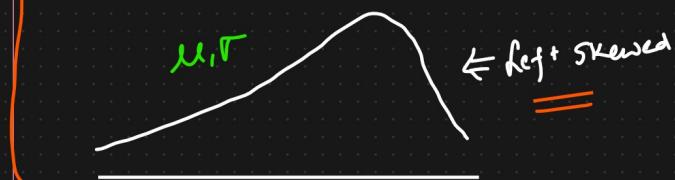
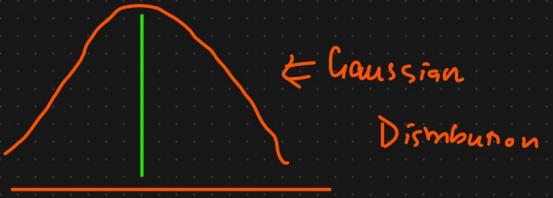
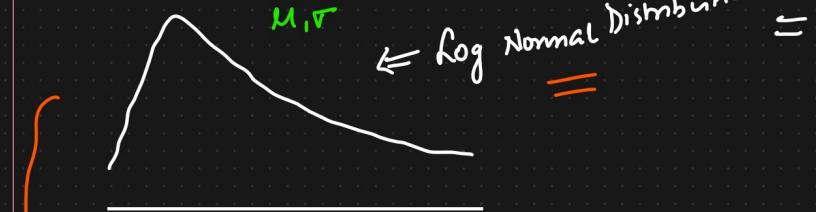
No. of samples



$$M, \sigma$$

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

Sampling with replacement



10 different region

$$n > 30$$

Size of shark through out the world?  $\rightarrow$  Assumptions

$N < 30$

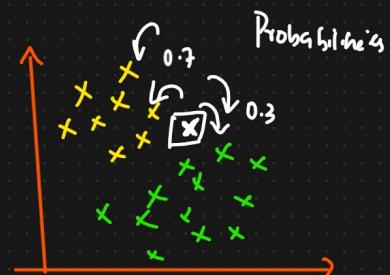
$m \uparrow \uparrow \uparrow$

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

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

Strongly → coin  $P(H) = 1$   
unfair coin

Strong  
Basic  
↑



Rolling a Dice

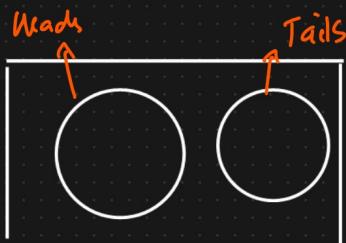
$$P(1) = \frac{1}{6} \quad P(2) = \frac{1}{6} \quad P(3) = \frac{1}{6}$$

① Mutually Exclusive Events

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

① Tossing a coin

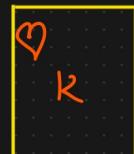
② Rolling a dice



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?

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

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

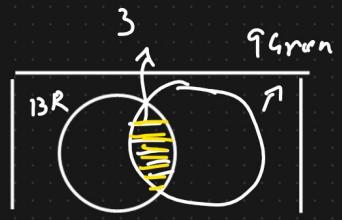
## Non Mutual Exclusive Event

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

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

$$= \frac{13}{52} + \frac{4}{52} - \frac{1}{52} = \boxed{\frac{16}{52}}$$

### # Multiplication Rule

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

Bag of marbles ○ ○ ○ X  
○ ○ ○

$$\Rightarrow P(W) = \frac{4}{7} \longrightarrow P(Y) = \frac{3}{6}$$

$\uparrow$   
white  
1 marble

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

Ans)  $P(1) = \frac{1}{6}$     $P(2) = \frac{1}{6}$     $P(3) = \frac{1}{6}$     $P(4) = \frac{1}{6}$

Independent   Dependent ↗

Multiplication Rule for Independent Events

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

$$= \frac{1}{6} * \frac{1}{6} = \boxed{\frac{1}{36}}$$

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

- Mutual Exclusive
- Non Mutual Exclusive

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

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

Dependent and Independent Events

Event A      ↓      Event B

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

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

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

②



$\Rightarrow$  Dependent Events

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

marble from the bag?

Ans)

$$P(O) = \frac{4}{7} \quad \boxed{P(Y/O)} \quad \text{conditional probability}$$

Orange Marble

Naive Bayes

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

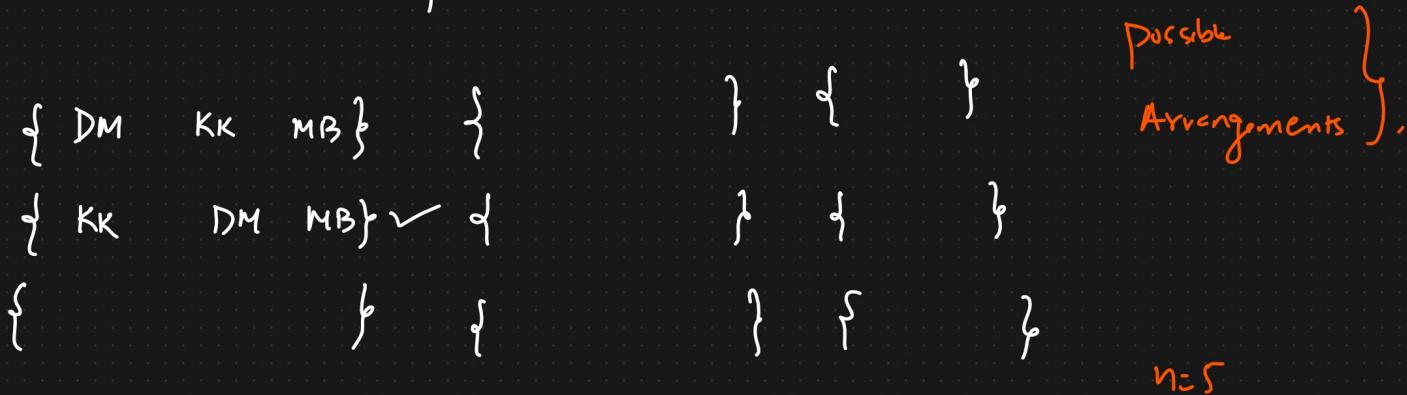
$$= \frac{4}{7} * \frac{3}{6} = \frac{4}{7} * \frac{1}{2} = \frac{2}{7} \approx 0.286$$

④ Permutation

$$\begin{array}{c} 5 \times 4 \times 3 \\ \hline = \boxed{\underline{60 \text{ ways}}} \end{array} \Rightarrow \text{Permutation}$$

{ Dairy Milk, Kit Kat, Milky Bar, }  
Sneakers, 5 star

With permutation, order matters



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

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

 $n$ : Total No. of Objects $r$ : # of selection⑤ Combination

Repetition will not occur

$$\{ DM \quad KK \quad MB \} \quad \text{Unique Combination}$$

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

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

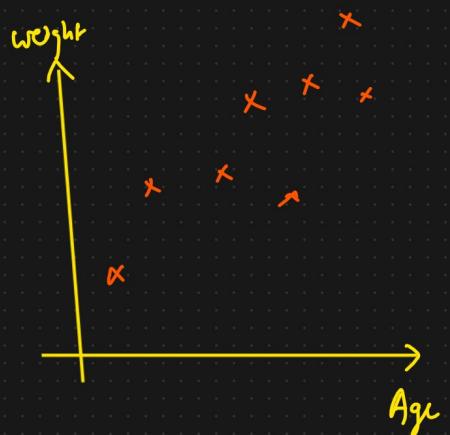
# DREAM 11

(\*) Covariance ✓

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

{Feature Selection}

Age ↑	Weight ↑
Age ↓	Weight ↓



Quantify the relationship

x & y using mathematical  
question

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

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

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

$\Downarrow$

$\boxed{2y} \Rightarrow +ve$        $\text{Cov}(x, x) \quad \Leftarrow$        $\boxed{ }$

$$\boxed{\text{Cov}(x, x) = \text{Var}(x)} \Leftarrow$$

+ve Covariance

$x \uparrow$	$y \uparrow$
$x \downarrow$	$y \downarrow$

-ve Covariance

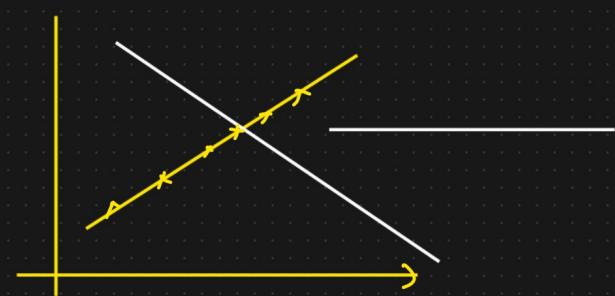
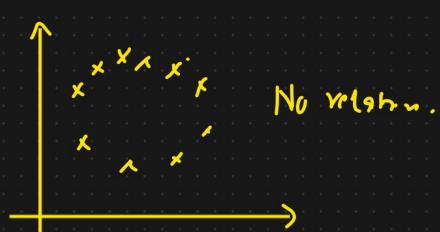
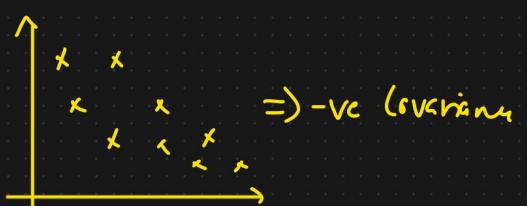
$x \uparrow$	$y \downarrow$
$x \downarrow$	$y \uparrow$

Covariance = 0

No relation  
with x & y



$\Rightarrow +ve$ .



X	Y
10	4
8	6
7	8
6	10
7.75	7

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

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

$$= -3.25$$

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

$$\underline{\underline{}}$$

① Pearson Correlation Coefficient (-1 to 1)

$$\rho(x, y) = \frac{\text{Cov}(x, y)}{\sqrt{x} \cdot \sqrt{y}}$$

More the value towards +1

More +ve correlated it is

-1  
negative correlated

Scale

-ve Covariance = +ve

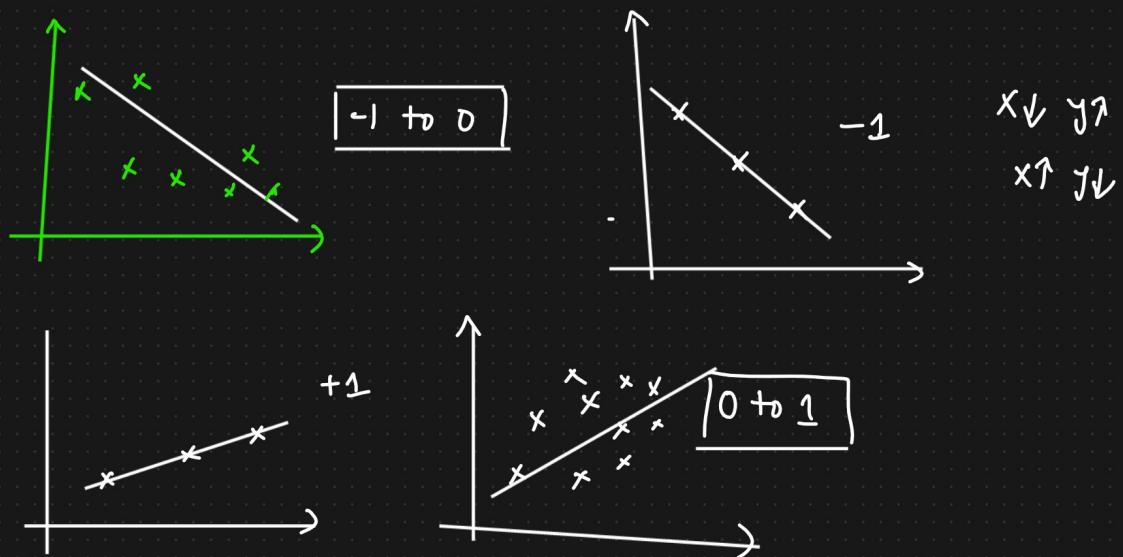
+ve

-ve

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

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

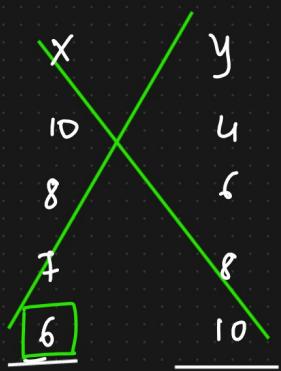
$$\boxed{0.5}$$



## (4) 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

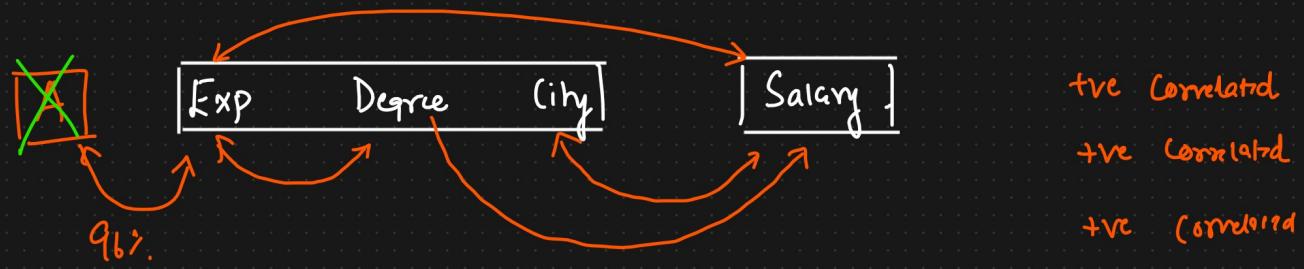
Ascending

Why this Correlation will be used?

$$0.95 \\ 95\%$$



$$\begin{array}{c} \downarrow \\ \text{+ve} \\ \text{-ve} \\ \uparrow \end{array} \} \text{ Good} \\ 0.2 \quad 0.01$$



## Inferential Statistics

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

Z test  
t test  
Chi square test  
Anova test (F-test)

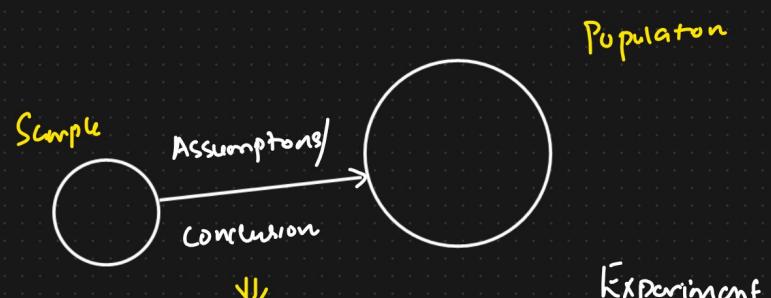
### 3 Distributions

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

Transformation

## Inferential Stat

### Steps of hypothesis Testing



### Hypothesis Testing

- ① Null Hypothesis: Coin is fair  $\Rightarrow$  Accepted  $\rightarrow$  [Coin is fair or not]
- ② Alternative hypothesis: Coin is not fair

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

### ③ Perform Experiments

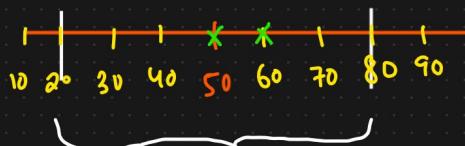
10  $\rightarrow$  Null Hypothesis is Rejected  
 $\rightarrow$  Alternative Hypothesis is Accepted

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

$$CI = [20 - 80]$$



Coin is fair



C-I  $\Rightarrow$  Confidence Interval

70 times  $\Rightarrow$  Domain Export



Confidence Interval

- ↳ We fail to Reject the Null Hypothesis [within C-I]  $\Rightarrow$  Conclusions
- ↳ We Reject the Null Hypothesis [outside C-I]  $\Rightarrow$  Conclusions

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  $\Rightarrow$

Vaccines  $\Rightarrow$  Medical  $\Rightarrow$  critical

Conclusions

Confidence Interval : (CI)

$\Rightarrow$  Domain Experiment

=

Significance Value

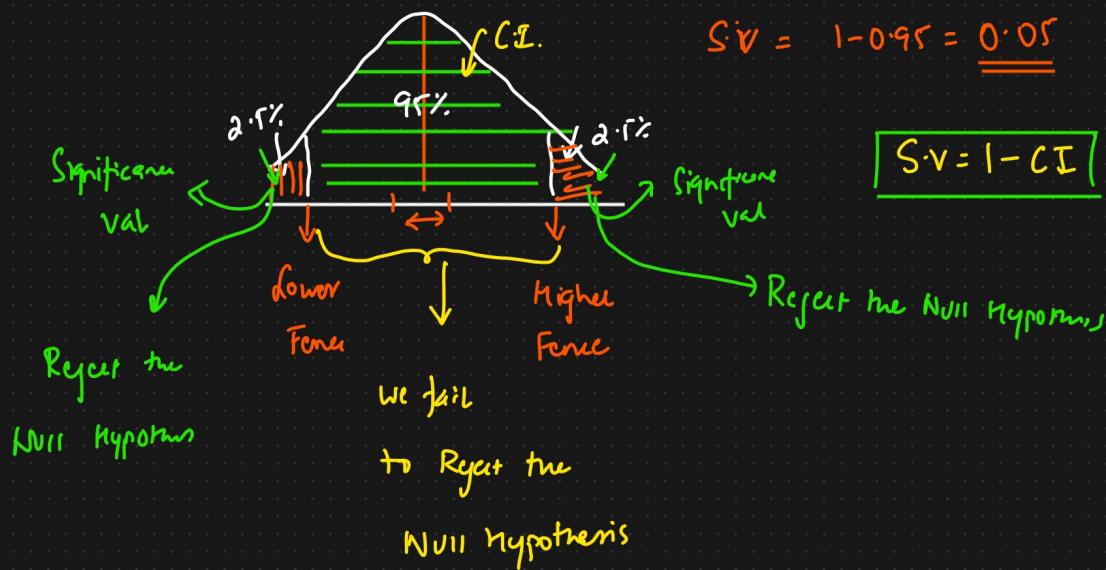
$\boxed{95\%}$

$C.I = 95\%$

$S.V = 1 - C.I$

$S.V = 1 - 0.95 = \underline{\underline{0.05}}$

$S.V = 1 - C.I$

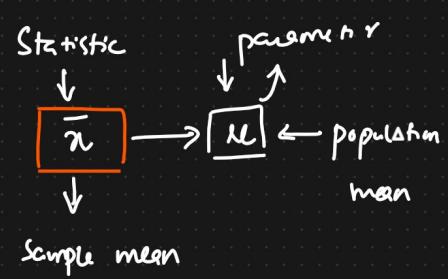


Point Estimate : The value of any statistic that estimates the value of a parameter is called Point Estimate

Point Estimate

$$\bar{x} \xrightarrow{\uparrow} \mu$$

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



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

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



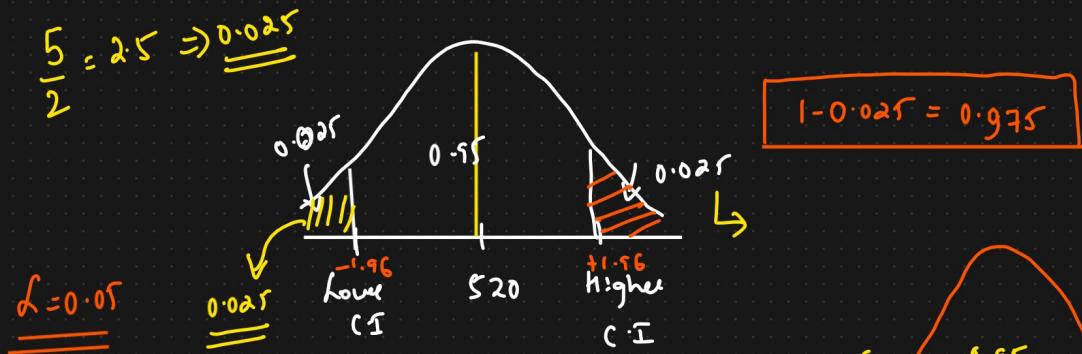
Higher C.I :- Point Estimate + Margin of Error

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

Population Std  
 $\alpha$  = 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?

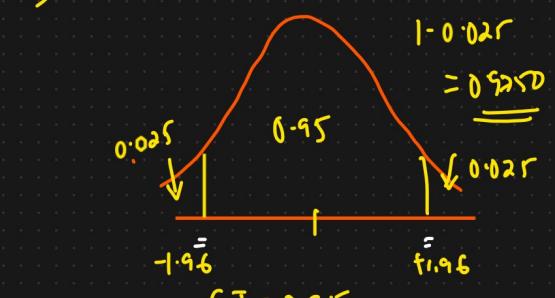
Ans)  $n=25 \quad \bar{x}=520 \quad \sigma=100 \quad C.I = 95\% \quad S.V = 1-C.I = 0.05$



Lower C.I = Point Estimate - Margin of Error

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

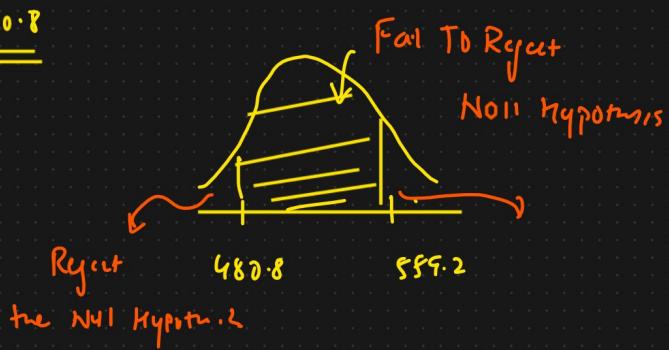
$$= 520 - Z_{0.025} \frac{100}{\sqrt{25}}$$



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

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

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

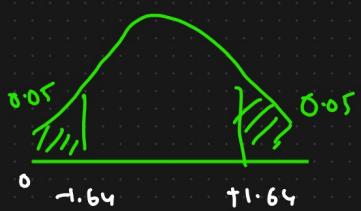


$$\text{Higher CI} = 520 + 1.96 \times 20 = \underline{\underline{559.2}}$$

Reject the Null Hypothesis

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

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



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

$$\text{Higher CI} = 480 + Z_{0.10/2} \left[ 17 \right]$$

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

$$= 480 + 27.8$$

$$= 480 - 1.64 \left[ 17 \right]$$

$$= 507.8$$

$$= 480 - 27.8 = 452.12$$

$$\left[ 452.12 \leftrightarrow 507.8 \right].$$

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?

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

$$\bar{x} \pm t_{f/2} \left( \frac{s}{\sqrt{n}} \right).$$

t test



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

19 20 21

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

$$= 520 - 2.064 \times 16$$

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

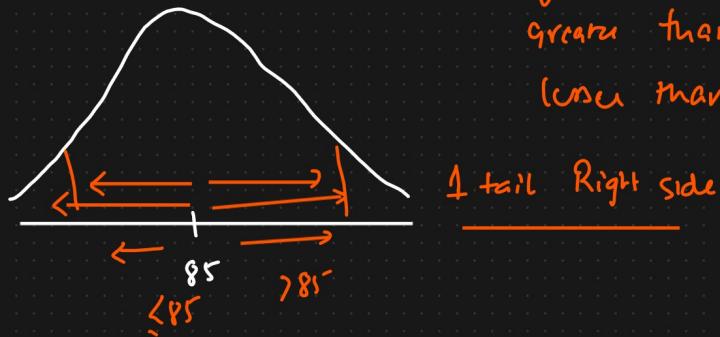
$$\text{Higher C.I} = 553.024$$



### 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?

Two Tail:



↓  
greater than 85%.

less than 85%.

1 tail Right side

- ① Z-test }  
② t-test }.

① 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% C.I, is there enough evidence to support Machine is working properly or not

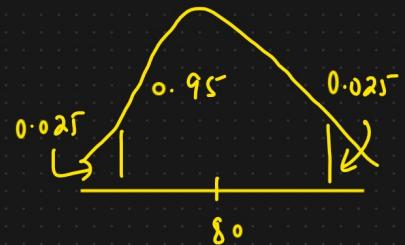
Step 1

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

$$\mu = 80 \text{ ml} \quad n = 40 \quad \bar{x} = 78 \quad s = 2.5$$

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

Step 2  $\therefore C.I = 0.95 \quad S.V(\alpha) = 1 - 0.95 = 0.05$



Step 3  $\therefore$

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

$$n = 40$$

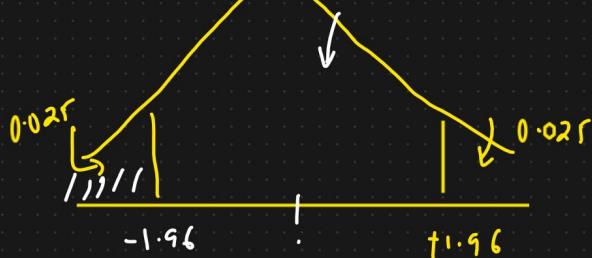
$$s = 2.5$$

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)

$$Z = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} \rightarrow \text{Standard Error}$$
$$= \frac{78 - 80}{\frac{2.5}{\sqrt{40}}} = -5.05$$

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 } There is some fault in the  
machine.

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.

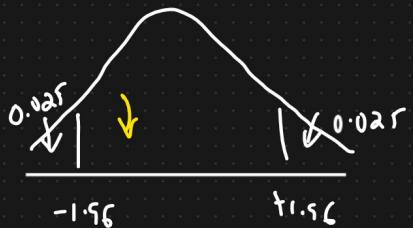
### Step 1

①  $H_0 : \mu = 32$

$H_1 : \mu \neq 32$

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

### Z test



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

$\sigma / \sqrt{n}$

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 hypothesis

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

= -

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?  $\{ \frac{95\%}{=} \}$

## Inferential Statistics

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

Z test  
t test  
Chi square test  
Anova test (F-test)

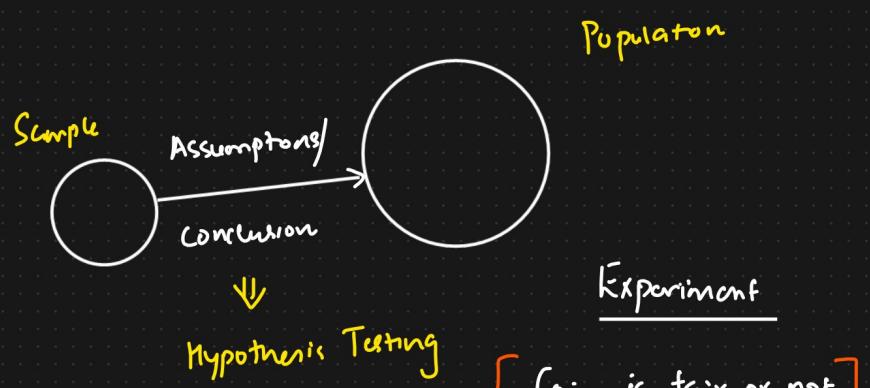
### 3 Distributions

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

Transformation

## Inferential Stat

### Steps of hypothesis Testing



- ① Null Hypothesis: Coin is fair  $\Rightarrow$  Accepted  $\rightarrow$  [Coin is fair or not]  $P(H)=0.5 \quad P(T)=0.5$
- ② Alternate hypothesis: Coin is not fair

### ③ Perform Experiments

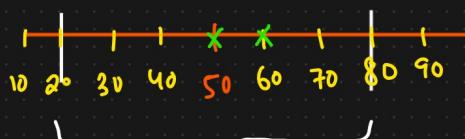
10  $\rightarrow$  Null Hypothesis is Rejected  
 $\rightarrow$  Alternate Hypothesis is Accepted

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

$$CI = [20-80]$$



Coin is fair



$C-I \Rightarrow$  Confidence Interval

70 times  $\Rightarrow$  Domain Export



Confidence Interval

- ↳ We fail to Reject the Null Hypothesis [within C-I]  $\Rightarrow$  Conclusions
- ↳ We Reject the Null Hypothesis [outside C-I]  $\Rightarrow$  Conclusions

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  $\Rightarrow$

Vaccines  $\Rightarrow$  Medical  $\Rightarrow$  critical

Conclusions

Confidence Interval : (CI)

$\Rightarrow$  Domain Experiment

=

Significance Value

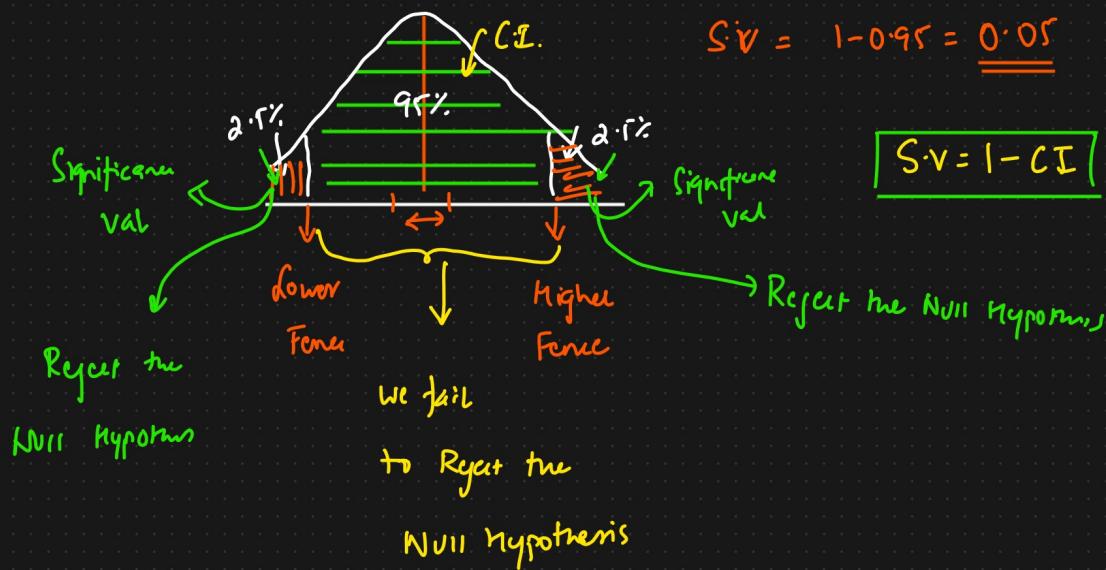
$\boxed{95\%}$

$C.I = 95\%$

$S.V = 1 - C.I$

$S.V = 1 - 0.95 = \underline{\underline{0.05}}$

$S.V = 1 - C.I$

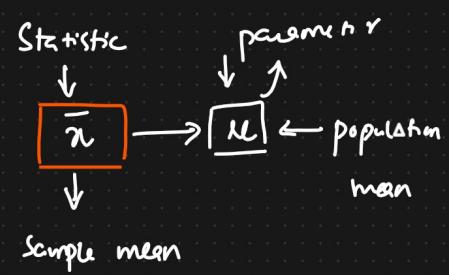


Point Estimate : The value of any statistic that estimates the value of a parameter is called Point Estimate

Point Estimate

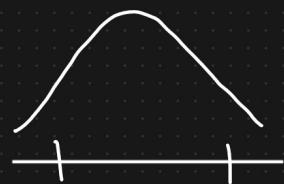
$$\bar{x} \xrightarrow{\uparrow} \mu$$

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



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

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



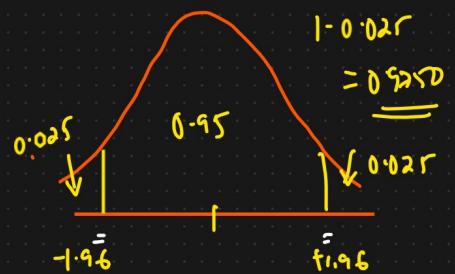
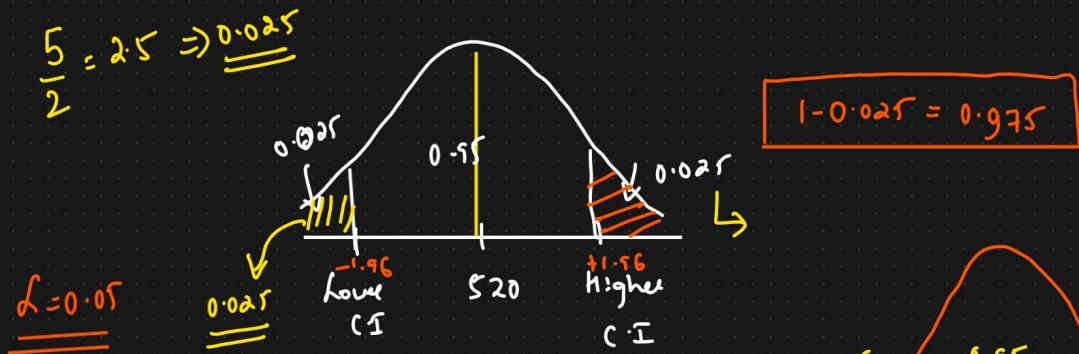
Higher C.I :- Point Estimate + Margin of Error

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

Population Std  
 $\alpha$  = 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?

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



Lower C.I = Point Estimate - Margin of Error

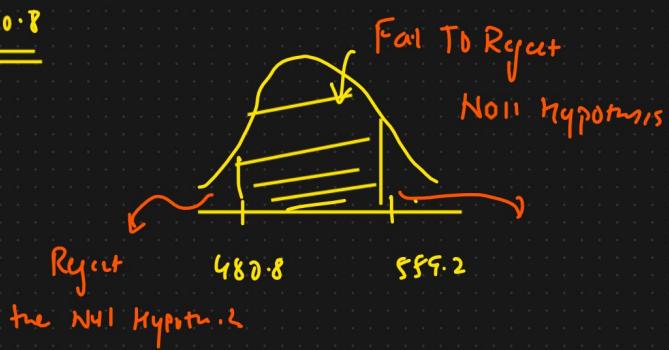
$$= 520 - Z_{0.05/2} \frac{\sigma}{\sqrt{n}}$$

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

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

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

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

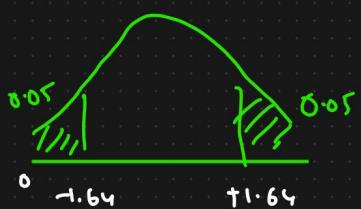


$$\text{Higher CI} = 520 + 1.96 \times 20 = \underline{\underline{559.2}}$$

Reject the Null Hypothesis

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

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



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

$$\text{Higher CI} = 480 + Z_{0.10/2} \left[ 17 \right]$$

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

$$= 480 + 27.8$$

$$= 480 - 1.64 \left[ 17 \right]$$

$$= 507.8$$

$$= 480 - 27.8 = \underline{\underline{452.12}}$$

$$\left[ 452.12 \leftrightarrow 507.8 \right].$$

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?

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

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

t test



0.05

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

19 20 21

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

$$= 520 - 2.064 \times 16$$

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

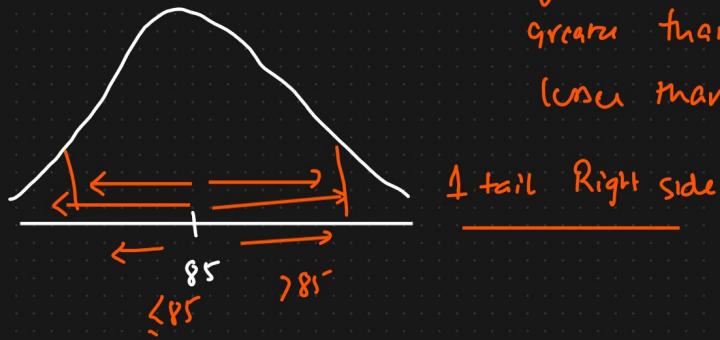
$$\text{Higher C.I} = 553.024$$



### 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?

Two Tail.



↓  
greater than 85%.

less than 85%.

1 tail Right side

- ① Z-test }  
② t-test }.

① 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% C.I, is there enough evidence to support Machine is working properly or not

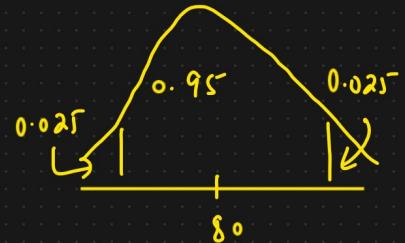
Step 1

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

$$\mu = 80 \text{ ml} \quad n = 40 \quad \bar{x} = 78 \quad s = 2.5$$

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

Step 2  $\therefore C.I = 0.95 \quad S.V(\alpha) = 1 - 0.95 = 0.05$



Step 3  $\therefore$

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

$$n = 40$$

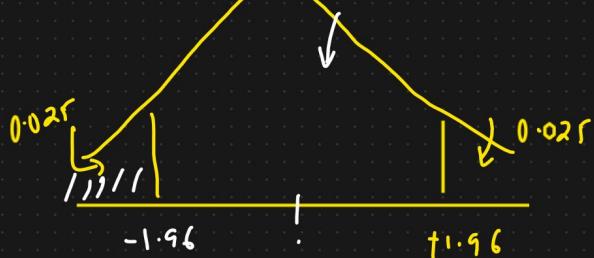
$$s = 2.5$$

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)

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

$\rightarrow$  Standard Error

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 } There is some fault in the  
machine.

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.

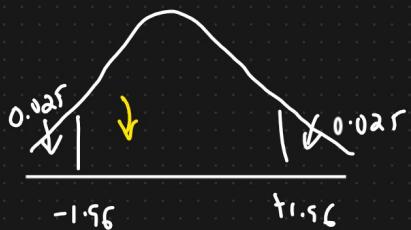
### Step 1

①  $H_0 : \mu = 32$

$H_1 : \mu \neq 32$

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

### Z test



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

$\sigma / \sqrt{n}$

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?

$$\text{Z-score} = \frac{-2.5 - 2.98}{\sqrt{0.25}} = -2.5$$

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?  $\left[ \begin{array}{l} 95\% \\ \hline \end{array} \right]$

$$\text{Z-score} = \frac{140 - 100}{\sqrt{15^2 / 30}} = 14.96$$

① A factory manufactures cars with a warranty of 5 years <sup>or more</sup> <sub>1</sub> 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?

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

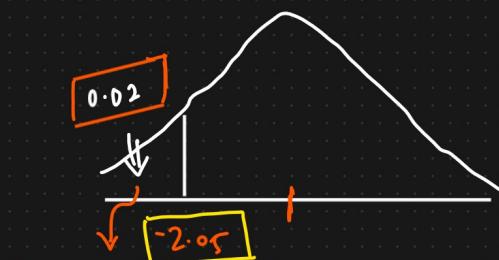
Step 1

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

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

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.

P-value  $\neq$  {Significance value}.  $\rightarrow$  C.I  
Out of all 100 random touches



$$P=0.02$$

$$P=0.02$$

$$\alpha = 0.02$$

P-value  $< \alpha$  [yes]

P-value,  $\boxed{-2.5298}$

Reject the Null Hypothesis

\* 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 & Alternative hypotheses }

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

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

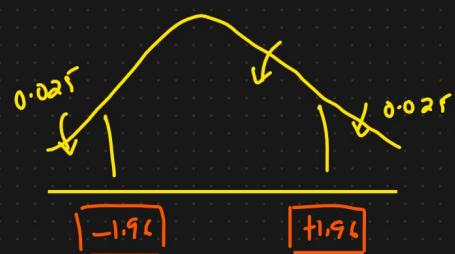
Step 1

$$H_0: \mu = 168$$

$$H_1: \mu \neq 168$$

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

Step 3

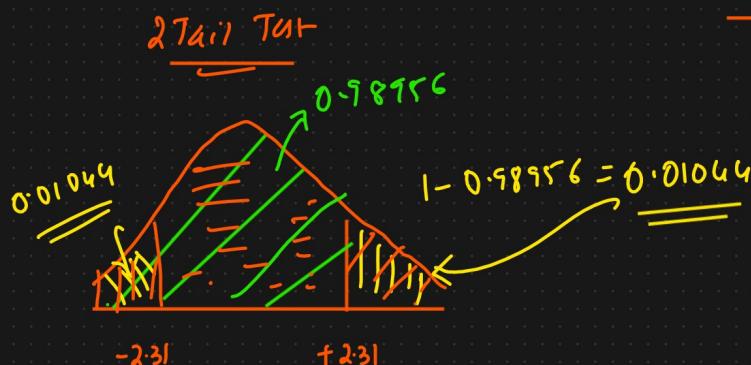


$$\text{Step 4 : } \bar{x}\text{-score} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{36}}} = \sqrt{2.31} \Rightarrow 2.3076$$

$2.31 > 1.96$  {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}.

④ 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 H<sub>0</sub>?

Ans) ①  $H_0 : \mu \geq 2$

② Step 2

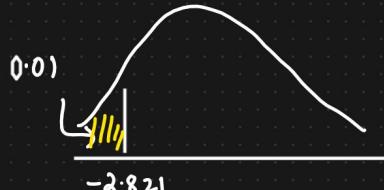
[n < 30]

Sample Standard deviation

$H_1 : \mu < 2$

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

③ Step 3



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

④ Calculate 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

⑥ 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?

⑦ State Null And Alternative Hypothesis?

⑧ At a 95% CI, is there enough evidence to reject the Null Hypothesis?

Ans) Step 1

Null Hypothesis:  $P_0 = 0.70$  ✓

Alternative Hypothesis:  $P_0 \neq 0.70$  ✓

$$q_0 = 1 - P_0 = 0.30$$

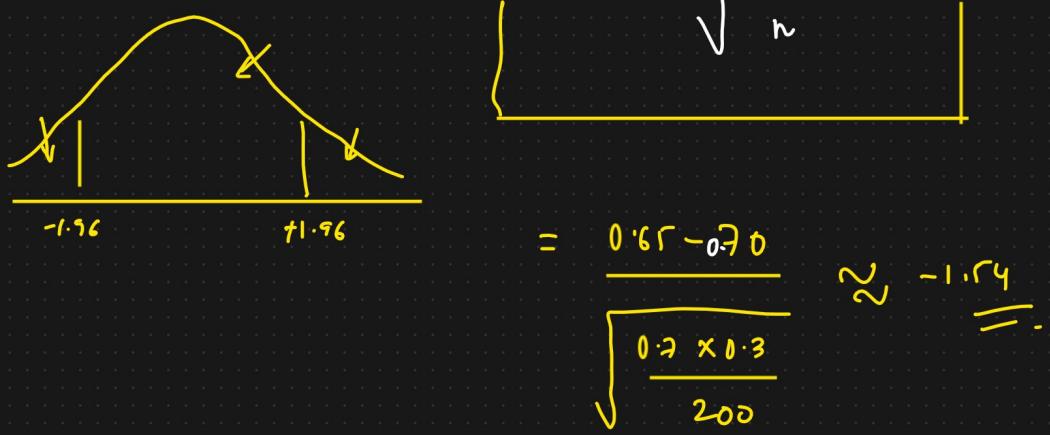
Step 2 :  $\alpha = 0.95$      $\delta = 0.05$

Step 3 :

$$\begin{aligned} n &= 200 \\ \hat{P} &= \frac{130}{200} = 0.65 \end{aligned}$$

Step 4 : Z test with proportion

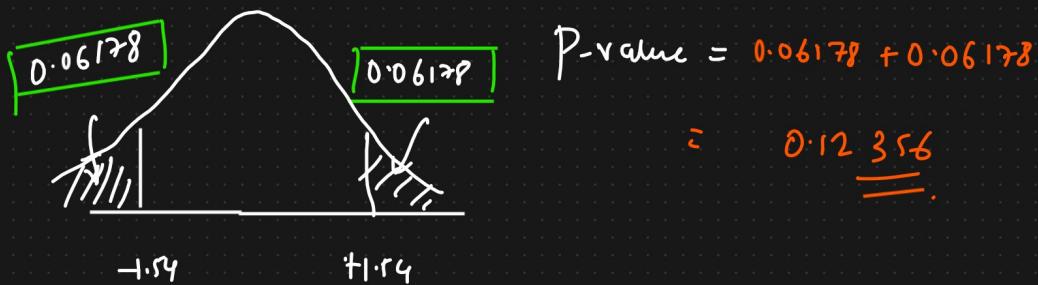
$$Z_{\text{test}} = \frac{\hat{P} - P_0}{\sqrt{\frac{P_0 q_0}{n}}}$$



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.

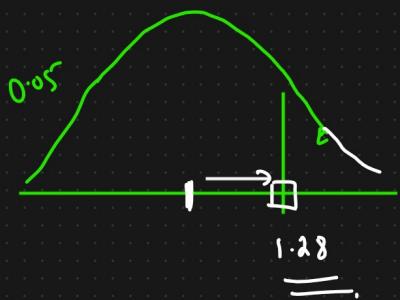
- State the Null & Alternative Hypothesis
- At 10% significance level, is there enough evidence to support the idea that vehicle ownership in City ABC is 60% or less?

$$H_0 : P_0 \leq 0.60$$

$$H_1 : P_0 > 0.60$$

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

$$q_0 = 0.40$$



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

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

Reject the Null Hypothesis

## ④ Chi Square test

① Chi Square test claims about population proportions.

ORDINAL DATA

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: NO

	$<18$	$18-35$	$>35$
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 :  $\alpha = 0.05$   $\rightarrow C.I = 95\%$ .

Step 3 : Degrees of freedom {Categories}.

$$df = C - 1 = 3 - 1 = \boxed{2}$$

$\hookrightarrow$  No. of categories.

Step 4 : Decision Boundary =  $\boxed{5.991}$  {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

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

