

## 2) Inferential Stats:

- Z-test
- t-test
- ANOVA → f-test
- CHISQUARE
- Hypothesis testing {p-values}
- Confidence interval.
- 2-table, t-table.

### # What is Statistics?

→ Statistics is the science of collecting, organising + analysing the data.

→ { Better Decision Making }

- Data? facts or pieces of information that can be measured.

### ① Descriptive Stats:

→ It consists of Organising + summarizing data.

### ② Inferential Stats:

→ Techniques where in we use the data that we have measured to form conclusions.

## # Sampling Techniques :

- Population
- population ( $N$ )
- Sample ( $n$ )

### 1) Simple random Sampling :

↳ Every member of the population ( $N$ ) has an equal chance of being selected for sample ( $n$ ).

### 2) Stratified Sampling :

↳ where the population ( $N$ ) is split into non overlapping groups (strata).

e.g.: Gender : 1) Male      2) Female      Survey

### 3) Systematic Sampling :

↳ We pick every  $n^{th}$  individual from the population ( $N$ ).

e.g.: Survey in a mall

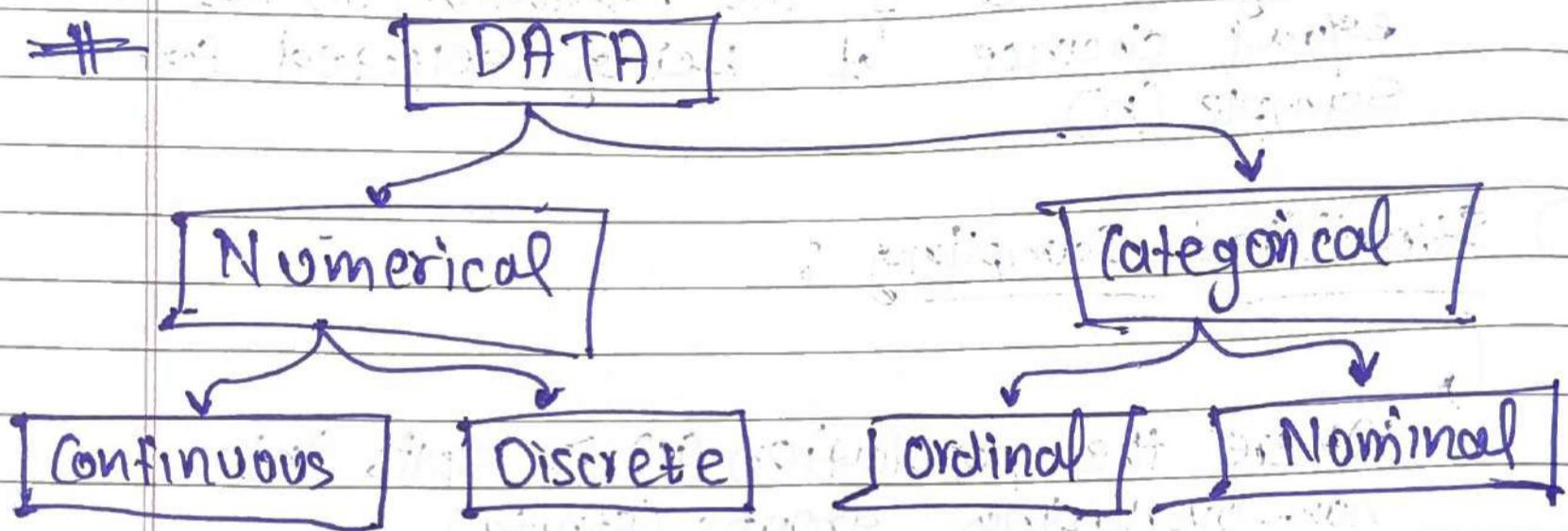
↳ pick every  $5^{th}$  person from the row.

## 4) Convenience Sampling :

↳

For ex : Doing a survey :

↳ pick people who have knowledge of data science.



## # Frequency Distributions :

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

Flower	Frequency	Cumulative frequency
Rose	3	3
Lilly	4	$3+4$
Sunflower	2	$9+2$

## ① Bar Graph :



(Categorical axis)

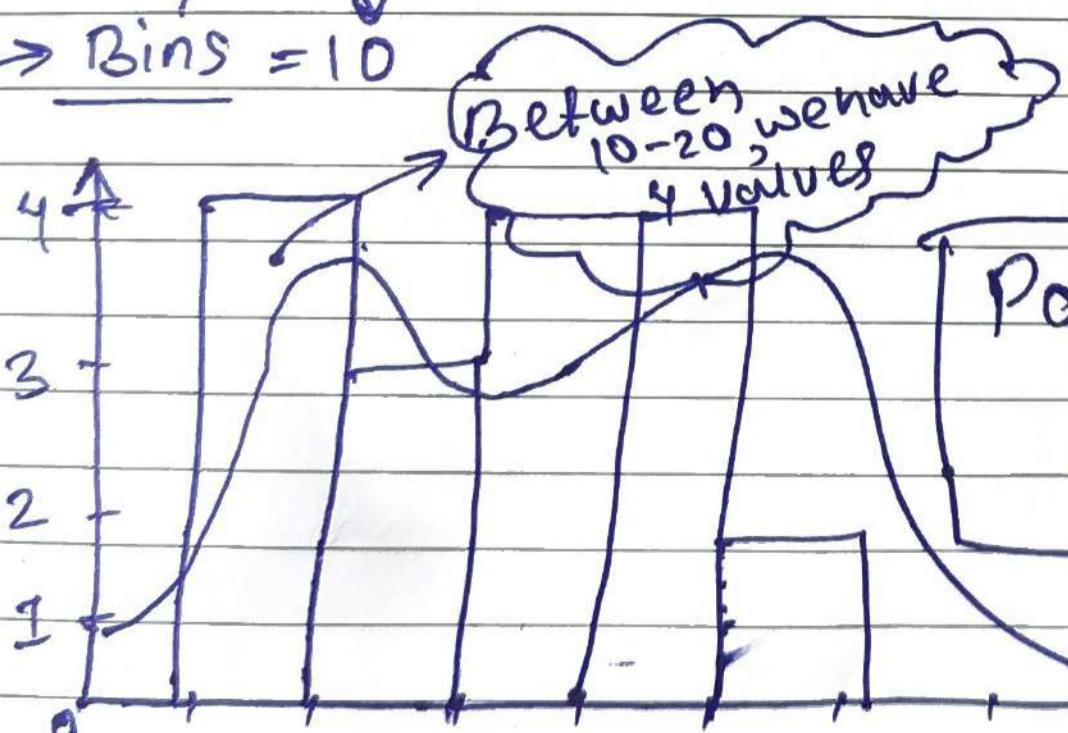
## ② Histograms :

(continuous data)

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

By default

Bins = 10



Between 10-20, we have 4 values

Pdf  
smoothening  
of histogram

- Pdfl - Probability density function.
- Kde - kernel density estimator.

Bar vs histogram

- $\rightarrow$  discrete data
- $\rightarrow$  continuous data

discrete data      continuous data

(class boundaries)

• 2000 observations

• Evaluation  
plots

0.0 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1.0 1.1 1.2 1.3 1.4 1.5 1.6 1.7 1.8 1.9 2.0 2.1 2.2 2.3 2.4 2.5 2.6 2.7 2.8 2.9 3.0 3.1 3.2 3.3 3.4 3.5 3.6 3.7 3.8 3.9 4.0 4.1 4.2 4.3 4.4 4.5 4.6 4.7 4.8 4.9 5.0 5.1 5.2 5.3 5.4 5.5 5.6 5.7 5.8 5.9 6.0 6.1 6.2 6.3 6.4 6.5 6.6 6.7 6.8 6.9 7.0 7.1 7.2 7.3 7.4 7.5 7.6 7.7 7.8 7.9 8.0 8.1 8.2 8.3 8.4 8.5 8.6 8.7 8.8 8.9 9.0 9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 9.9 10.0

spike plot

0.1 = 2000 s.

0.05 0.1 0.15 0.2 0.25 0.3 0.35 0.4 0.45 0.5 0.55 0.6 0.65 0.7 0.75 0.8 0.85 0.9 0.95 1.0

• histograms

• probability plots

0.0 0.2 0.4 0.6 0.8 1.0 1.2 1.4 1.6 1.8 2.0 2.2 2.4 2.6 2.8 3.0 3.2 3.4 3.6 3.8 4.0 4.2 4.4 4.6 4.8 5.0 5.2 5.4 5.6 5.8 6.0 6.2 6.4 6.6 6.8 7.0 7.2 7.4 7.6 7.8 8.0 8.2 8.4 8.6 8.8 9.0 9.2 9.4 9.6 9.8 10.0

## # Measure of central tendency :

i) Arithmetic Mean for population & sample

Mean  $\rightarrow$  Average

Population ( $N$ )

Sample ( $n$ )

$$X = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\} \quad \bar{x} = \frac{\sum_{i=1}^n x_i}{n}$$

$$\mu = \frac{\sum_{i=1}^N x_i}{N} = 3.2$$

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

= 3.2 (mean of random data taken)

$\Rightarrow$  central Tendency

① Mean ② Median ③ Mode.

Refers to the measure used to determine the centre of the distribution of data.

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

Mean : 3.2

new Data :  $\{1, 100\}$

Mean :  $\frac{32 + 100}{11} = 12.1$

Old Mean : 3.2

~~old~~

New Mean : 12

- By just adding a new data point
- which is far away from the other datap.
- we called it outliers.
- worst impact on the distribution.

Now, let's take median.

⇒ Median :

↳ Sort the numbers in ascending

\* ↳

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

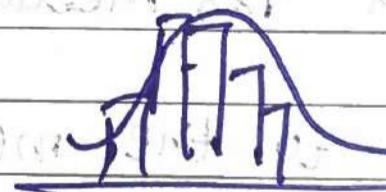
Total number

of elements or

data points : 11

middle element : 3

Distribution :



→ if, even then take the mean of  
the two middle datapoints.

So, our median is : 3

It means if we have outliers in our data, then mean will fluctuate very much so, we consider median, it handles it well.

$\Rightarrow$  Mode: Most frequent element.

Data:  $\{1, 2, 3, 3, 4, \underbrace{6, 6, 6}\}$

Mode: 3

Outlier data:  $\{1, 2, 3, 3, 6, 6, 6, 100, 100, 100, 100\}$

Mode: 100

So, everytime when outliers are there we will use median as a central tendency.

# mode is used when we have missing values we can replace it with mode.

$\hookrightarrow$  Categorical data.

## # Measure of Dispersion $\rightarrow$ Spread

① Variance

② Standard Deviation

# What is Spread?

$$\text{data-1} : \{1, 1, 2, 2, 4\}, \bar{x}_1 = 2$$

$$\text{data-2} : \{2, 2, 2, 2, 2\}, \bar{x}_2 = 2$$

Both the means are same, but data points are different.

Question: If the means are same

then what makes the dist. different.

i) Variance

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

For population

For Sample :

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

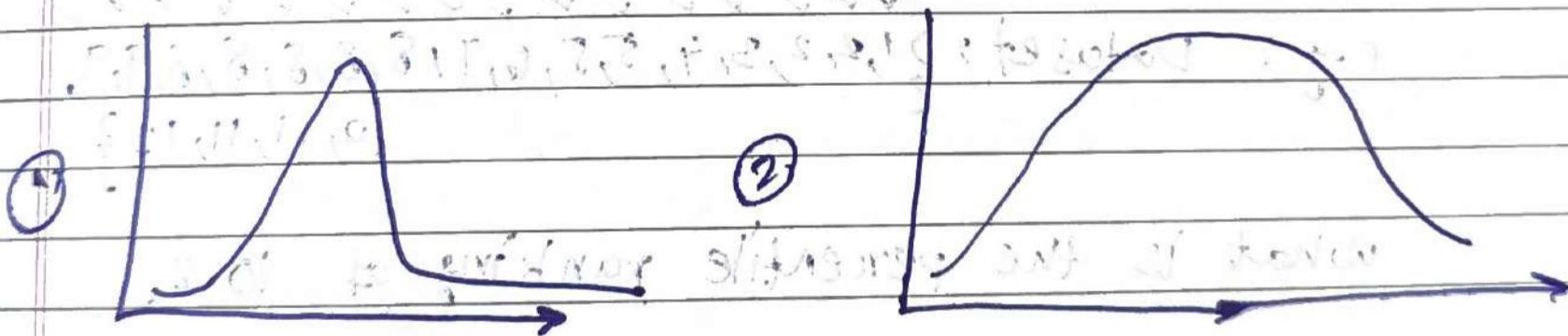
- The difference is in denominator,

Population  $\circ N$

Sample  $\circ$

$$n-1$$

Why?



Less variance

more variance  
high variance

③ Standard Deviation

$$\sigma = \sqrt{Var}$$

# # Percentiles + Quartiles {Find Outliers}

Data: 1, 2, 3, 4, 5

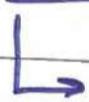
% of the odd numbers?

% of odd = Amount of odd no.

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

Total no. of odd

- Percentiles:



A percentile is a value below which a certain percentage of observation lie.

e.g.: Dataset: {1, 2, 2, 3, 4, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12}

What is the percentile ranking of 10?

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

$$= \frac{16}{20} \times 100 = [80^{\circ}\text{ile}]$$

\* it means 80% of the entire population lies below 10.

• What is value exists at percentile ranking of 25%?

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

$$= \frac{25}{100} \times 21 = 5.25$$

↳ Index position

↳ So, 5.25 lies b/w 5, 5 in the dataset.  
Taking avg  $\frac{5+5}{2} = 5$

↳ 75%?

$$= \frac{75}{100} \times 21 \leftarrow 18.75^{\text{th}} \text{ position.}$$

↳ 9  $\rightarrow 75\%$

$$= (14\%) \times 21 = 29.4 \approx 30$$

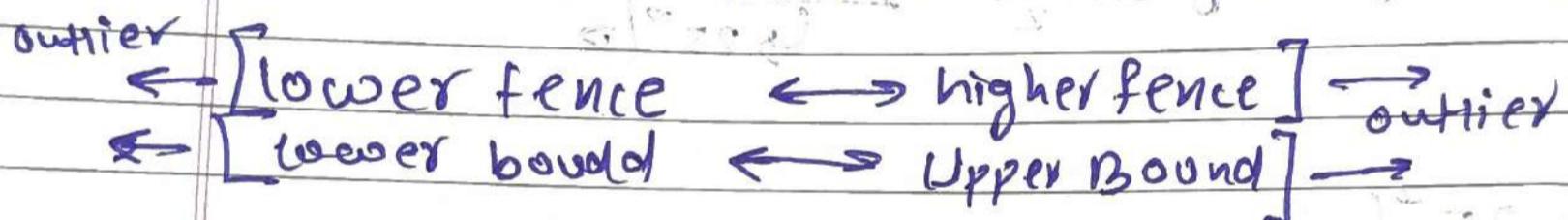
$$= 18.75 - 18 = 0.75 \times 30 = 22.5$$

## # Five Number Summary

- ① Minimum
- ② First quartile ( $Q_1$ )
- ③ Median
- ④ Third quartile ( $Q_3$ )
- ⑤ Maximum

$\Rightarrow$  Removing the outlier

$$\{1, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 8, 9, 27\}$$



$$\text{lower fence} = Q_1 - 1.5 \times \text{IQR}$$

$$= Q_1 - 1.5 \times (Q_3 - Q_1)$$

$$\text{upper fence} = Q_3 + 1.5 \times (\text{IQR})$$

$$\Rightarrow \text{IQR} = \text{Inter Quartile Range.} = Q_3 - Q_1$$

$$Q_1 = 25\% = \frac{25}{100} \times (19+1) = \frac{5}{5}^{\text{th}} \text{ elem}$$

$\boxed{3} \rightarrow 25\%$

$$Q_3 = 75\% = \frac{75}{100} \times (19+1) = \frac{15}{5}^{\text{th}} \text{ elem}$$

$\boxed{17} \rightarrow 75\%$

$$\text{IQR} = 7 - 3 = 4$$

$$\text{lower fence} = \cancel{15}(0)$$

$$\begin{aligned} &= Q_1 - 1.5(\text{IQR}) \\ &= 3 - 1.5(4) \\ &= -3 \end{aligned}$$

$$\text{upper fence} = Q_3 + 1.5(\text{IQR})$$

$$\begin{aligned} &= 7 + 1.5(4) \\ &= 13 \end{aligned}$$

any value  $< -3$   
(outlier)

any value  $> 13$   
(outlier)

Outlier we got  $\therefore 27$

remove  $\therefore 27$

Remaining

Data:

1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9

Minimum = 1

$$Q_1 = 3$$

$$\text{Median} = 5$$

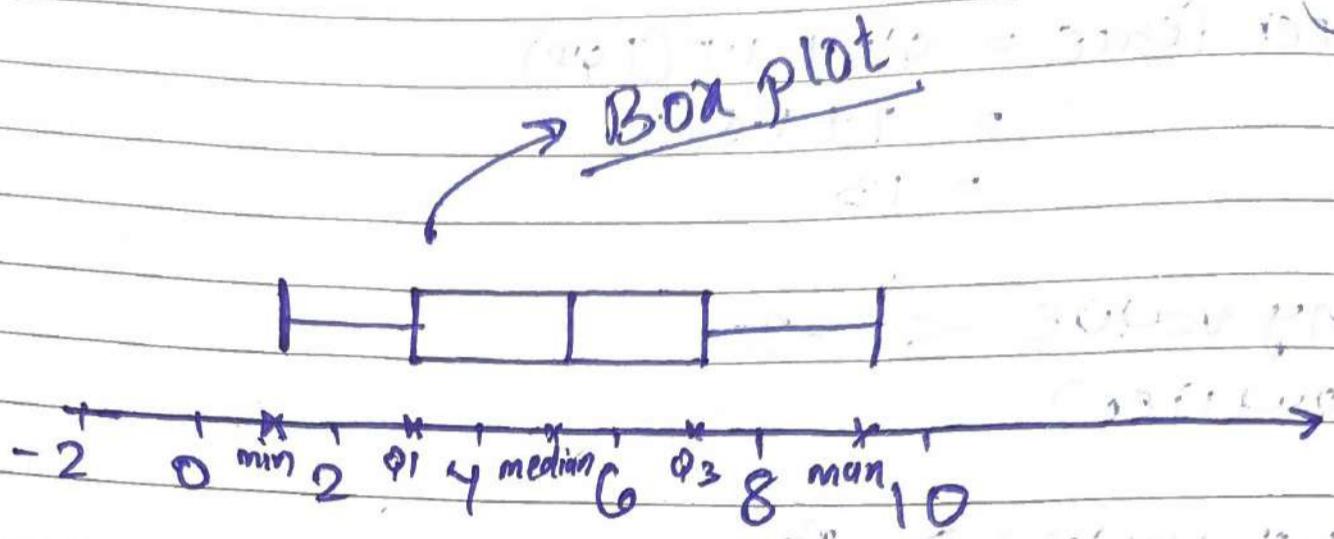
$$Q_3 = 7$$

$$\text{Max} = 9$$

Five number Summary

Now,

∴ Drawing Boxplot



⇒ Boxplot

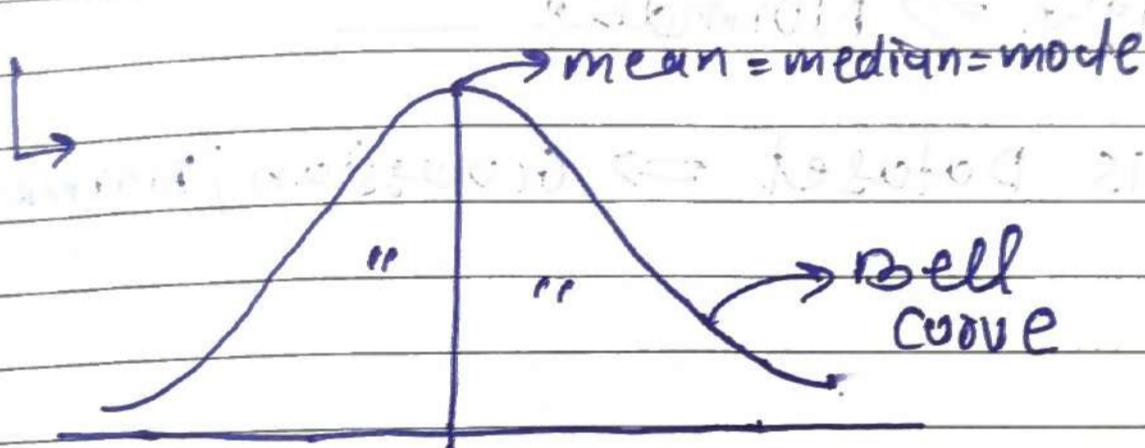
↳ can be used to determine outlier.

#

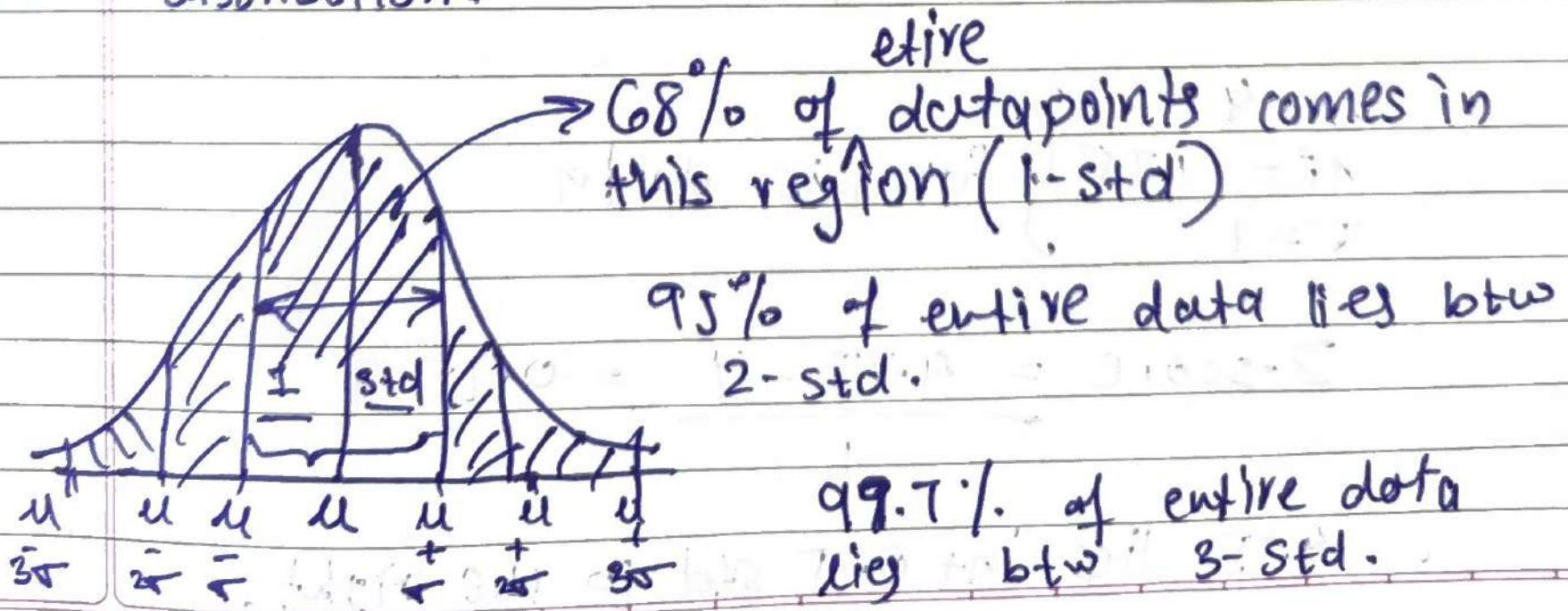
## Distributions:

- ↳ normal
- ↳ Standard normal
- ↳ z-score
- ↳ log normal
- ↳ Bernoulli
- ↳ Binomial.

### (1) Gaussian / Normal Distribution



- Symmetrical
- Mean = median = mode
- Best or perfect for ML models
- we can derive very important things from this distribution.



## → Empirical formula %

68 - 95 - 99.7 %. Rule

For e.g.: Height  $\Rightarrow$  Normally distributed  
 ↳ Domain Expertise

2) Weight  $\Rightarrow$  Normally —

3) IRIS Dataset  $\Rightarrow$  Gaussian / Normal

# 2-score

↳ How much std away it is from the mean.

$$\text{2-score} = \frac{x_i - \mu}{\sigma}$$

$$\mu = 4$$

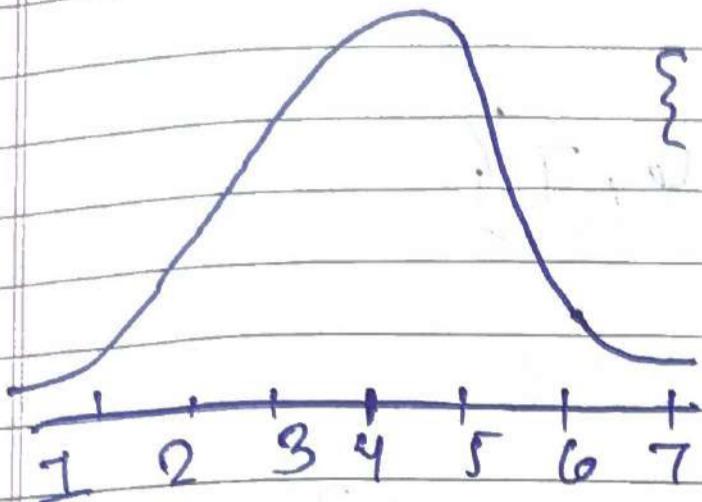
$$x_i = 4.75 \text{ for this entry}$$

$$\sigma = 1$$

$$\text{2-score} = \frac{4.75 - 4}{1} = 0.75$$

4.75 lies int 0.75 std to the right.

NOCO,



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

$$\mu = 4$$

$$\sigma = 1$$

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

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

$$z(1) = \frac{1-4}{1} = -3 \quad \{ -3, -2, -1, 0, 1, 2, 3 \}$$

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

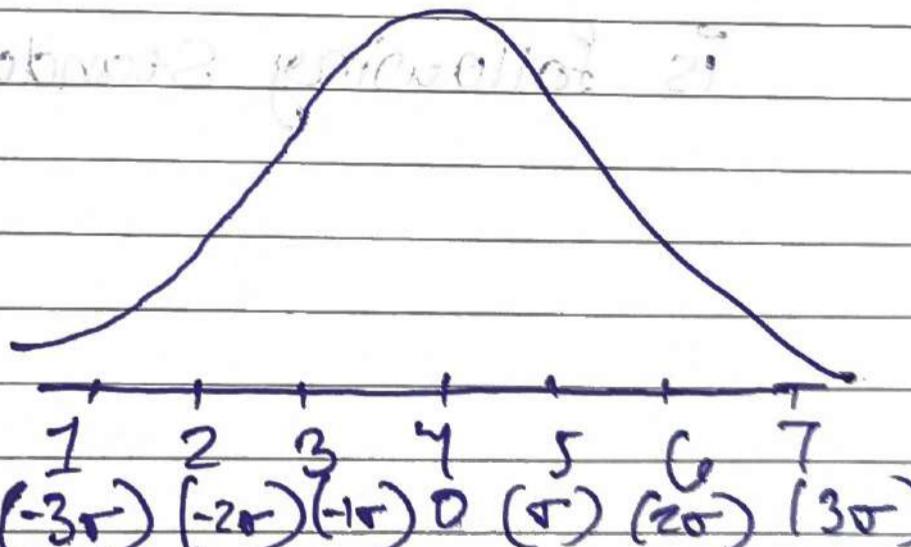
$$z(3) = \frac{3-4}{1} = -1$$

$$z(4) = 0$$

$$z(5) = 1$$

$$z(6) = 2$$

$$z(7) = 3$$



Initial data

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

$\downarrow$  Z-score

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

This distribution is called

[STANDARD NORMAL DISTRIBUTION]

any data satisfying ( $\mu=0$  &  $\sigma=1$ )

is following Standard Normal Distribution

## Now, Practical Application :-

Dataset :-

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

↳ Z-score

↳ Standard Normal distribution.

This process is called

\* Standardization.

↳ converts distribution into standard normal distribution

and the parameter is  $\{ \mu = 0 \text{ & } \sigma = 1 \}$

# Normalisation:

$$\hookrightarrow (0 \text{ to } 1)$$

we range our data points between 0 to 1.

$\Rightarrow$  Min Max Scaler

we can range in any lower &

upper value like from (-1 to +1) :

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

$$(-3 \text{ to } +3)$$

normalization helps in breaking

values of different size

breaking of redundant features

non-redundant features

the order of columns will be

## # Probability :

↳ is the measure of the likelihood of an event.

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

$$P(E) = \frac{1}{6} = \frac{\text{number of ways an even can occur}}{\text{number possible outcomes.}}$$

## # Addition Rule :

↳ (probability, "or")

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

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

we can get only one output after rolling only single event is occurring at the same time.

~~Non-mutual exclusive~~

Eg: Deck of cards  $\{\text{Queen}, \heartsuit\}$

Q] Probability of A or B (head or tail)  
in tossing a coin.

$\Rightarrow$  mutually exclusive

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

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

$\Rightarrow$  Roll a dice

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

Since each face has equal probability  
So  $P(1) = P(3) = P(6) = \frac{1}{6}$

$$\therefore P(1 \text{ or } 3 \text{ or } 6) = 3 \times \frac{1}{6} = \frac{1}{2}$$

# For Non-mutually Exclusive

Q] You are picking a deck from decks  
card.

Prob of queen on a heart.

$\Rightarrow$  Non-mutually exclusive

$$P(Q) = \frac{4}{52}$$

$$P(B) = \frac{13}{52}$$

$$P(Q \text{ and } B) = \frac{1}{52}$$

Addition Rule for Non-mutually exclusive events

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

$$\begin{aligned} P(Q \text{ or } B) &= P(Q) + P(B) - P(Q \text{ and } B) \\ &= \frac{4}{52} + \frac{13}{52} - \frac{1}{52} \end{aligned}$$

$$P(Q \text{ or } B) = \frac{16}{52}$$

### ③ Multiplication Rule :

#### 1) Independent events

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

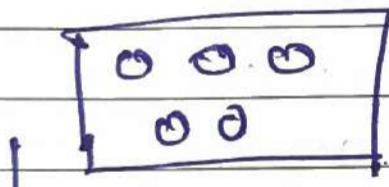
1 is not depend on 2  
any event is not depend on any other event.

#### 2) Dependant Event

Eg :

We have a bag with 3 red balls & 2 green.

$$P(\text{red}) = \frac{3}{5}$$

 picking a red ball.

Eg : but after picking red ball the total number of red balls become 2 and total balls become 4.

Now,  $P(\text{green}) = \frac{2}{4} = \frac{1}{2}$ ,

• <sup>2nd</sup> green event is depend on the <sup>1st</sup> event

g] Rolling a dice, 1<sup>st</sup> event got 4 & 2<sup>nd</sup> event got 5  
 probability of 5 + 4 ? (Independent event)  
 $P(5 \text{ and } 4) = P(5) \times P(4)$

Multiplication Rule.

$$P(5 \text{ and } 4) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

g] What is the probability of drawing a queen and then a Ace from a deck of cards?  
 (Dependent event)

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

Naive Bayes

conditional probability

$P(B|A) \rightarrow$  Probability of B after A event occurred.

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

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

### ③ Multiplication Rule :

#### 1) Independent events

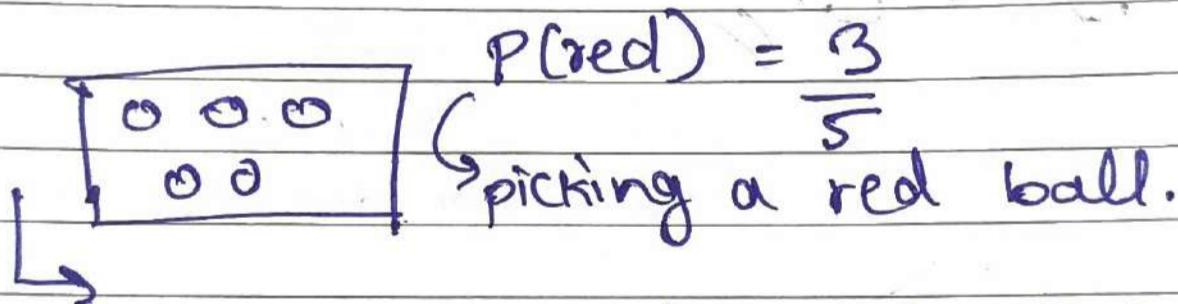
Eg : Rolling a dice :  $\{1, 2, 3, 4, 5, 6\}$ .

1 is not depend on 2  
any event is not depend on any other event.

#### 2) Dependant Event

Eg :

We have a bag with 3 red balls & 2 green.



Eg : but after picking red ball the total number of red balls become 2 and total balls become 4.

Now,  $P(\text{green}) = \frac{2}{4} = \frac{1}{2}$ ,

2nd <sup>good</sup> event is depend on the 1<sup>st</sup> event

Q] Rolling a dice, 1st event got 5 & 2nd event got 4

probability of 5 & 4? (Independent event)

$$P(5 \text{ and } 4) = P(5) \times P(4).$$

Multiplication Rule.

$$P(5 \text{ and } 4) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

Q] What is the probability of drawing a queen and then a Ace from a deck of cards? (Dependent event)

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

Naive Bayes  
Conditional probability

$P(B|A)$  → Probability of B after A event occurred

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

$$= \frac{4}{52} \times \frac{3}{51}$$

## # Permutation & Combination

### i) Permutation:

School trip  $\Rightarrow$  (chocolate factory)

$\rightarrow$  Dairy milk, 5-star, Milky bar,  
Eclairs, Gems, Silks.

Student  $\Rightarrow$  Assignment.

ask Student to write 3 names of chocolates they like.

S-I

(1) (2) (3)  
How many options this student have?

(1) at first place student has 6 choices

6

(2) Now 5 choices left and so on.

6    5    4

So, Total permutation:

$$6 \times 5 \times 4 = 120,$$

Permutation :  ${}^n P_r = \frac{n!}{(n-r)!}$   
 formula :  $\uparrow$  Total no. of choc.

how many  
no. asked.

$$\text{Given } n=6, r=3 \Rightarrow {}^6 P_3 = \frac{6!}{(6-3)!} = \frac{6 \times 5 \times 4 \times 3!}{3!} = 6 \times 5 \times 4 = 120.$$

$$= 120$$

Combination:

$$\frac{n!}{r!(n-r)!} = \frac{6!}{3!(6-3)!} = \frac{6!}{3!3!} = \frac{6 \times 5 \times 4 \times 3!}{3!3!} = \frac{6 \times 5 \times 4}{3 \times 2 \times 1} = 120$$

$$= 120$$

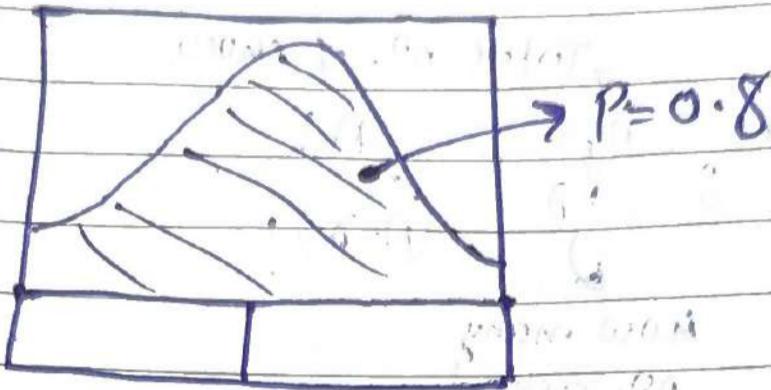
$$120 = (5!)^2 = (5 \times 4)^2$$

(5!)  $\neq$  120

#

## P-value

Mouse Pad:



$$P = 0.8$$

The shaded region is the region where we use the mousepad most

$$P = 0.8$$

↳ means every 100 times I touch the mouse pad I will touch here 80 time.

Coin  $\Rightarrow$  whether it is fair or not by performing 100 tosses.

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

↳ for a fair coin.

# Inferential STATS

## # Hypothesis Testing

- (1) Null hypothesis: coin is fair.  
(default)
- (2) Alternative hypothesis: coin is unfair  
(opposite of null hypothesis).
- (3) Experiment.
- (4) Reject or accept Null Hypothesis.

So, after deciding Null & Alternative hypo -  
we will perform experiment

we get result : 30 times head  
70 times tail

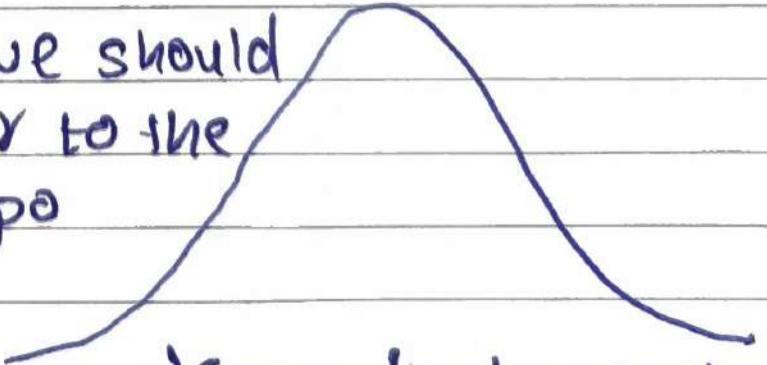
Conclusion?

#  $\rightarrow$  Significance value ( $\alpha$ )

$\alpha = 0.05$   $\rightarrow$  considered maximum.

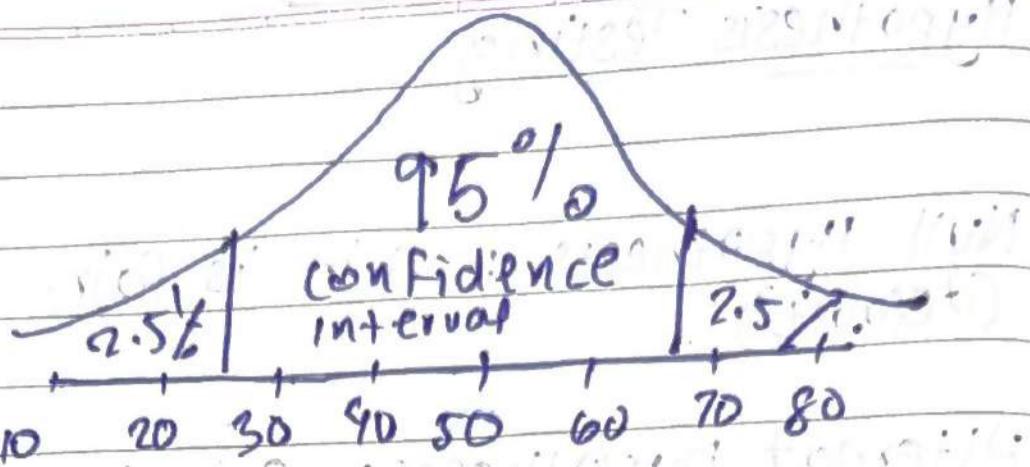
To make conclusion we should have the result nearer to the mean, to proof null hypo correct.

How much = significance near to mean? value ( $\alpha$ )



$$\alpha = 0.05$$

s.l.



$$C.I \text{ (confidence interval)} = 100 - \alpha \\ = 95\%$$

Now, if the experiment b/w H & T inside the C.I then its fair.

$\begin{cases} 30-H \\ 70-T \end{cases}$  if the experiment b/w H & T inside the C.I then its fair.

→ it's inside the C.I

[coin is fair]

NULL-hypothesis

- we accept null hypothesis

if it's 10-Head

- we reject null hypothesis

## # Type-1 & Type-2 Errors

- Null hypothesis ( $H_0$ )  $\Rightarrow$  coin is fair
- Alternate hypothesis ( $H_1$ )  $\Rightarrow$  not fair.

Reality check:

- $\hookrightarrow$  Null hypothesis is true
- $\hookrightarrow$  Null hypothesis is false

Outcome :

1]

We reject the null hypothesis, when in reality it is false :

- $\hookrightarrow$  So, it's a good thing isn't it?

[This is good one].

Outcome :

2] we reject null ( $H_0$ )  
in reality it's true:

- $\hookrightarrow$  So, it's bad thing, because:  
we said ( $H_0$ ) coin is unfair, but it's fair

[this is Type-1 error]

Outcome-3 : We accept ( $H_0$ )  
 (coin is fair)

but in reality its false

$\hookrightarrow$  error!

[Type-II error]

e.g: email spam

e.g: Death sentence

Outcome-4 : We accept  $H_0$

(coin is fair)

but in reality true

$\hookrightarrow$  good case

# Confusion

Matrix :

P N

TP TN

FP FN

$\Rightarrow$  we said unfair  
 Fair, actual its unfair  
 (Type-II error)

T = TRUE

F = FALSE

P = positive

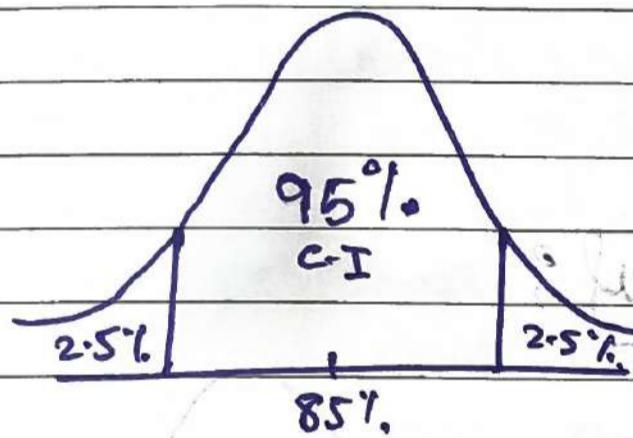
N = negative

$\Rightarrow$  we said unfair,  
 actual its fair

## # 1-tail & 2-tail test

E.g.: Colleges in Karnataka have an 85% placement rate. A new college was recently opened and it was found that a sample of 100 students had a placement rate of 88% with a standard deviation 4%. Does this college has a different placement rate? ( $\alpha=0.05$ )

### 2-tailed Test :



We have to check whether 88% is lying on the graph.

it may lie below and above 85%

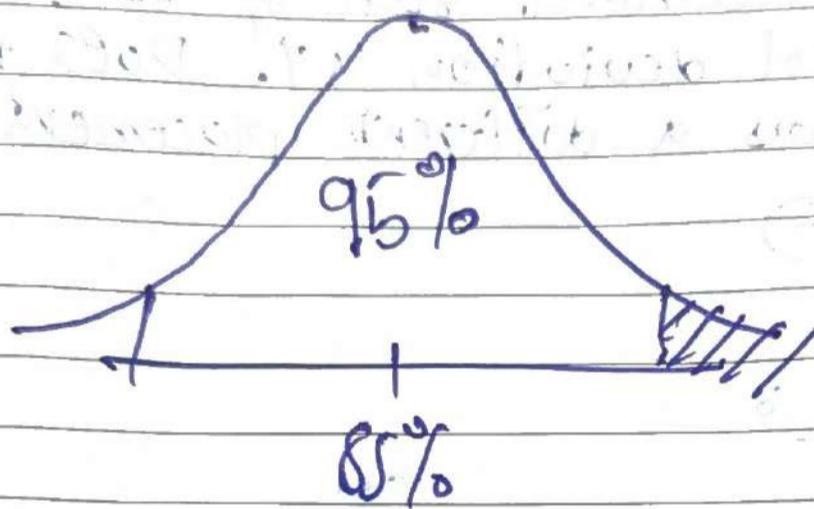
That's why it is two-tailed-test

So we will have to calculate the Z-score to calculate the probability.

let's change the question:

Does this college have a placement rate greater than 85%?

We have to check only above (greater)



1-tailed-test

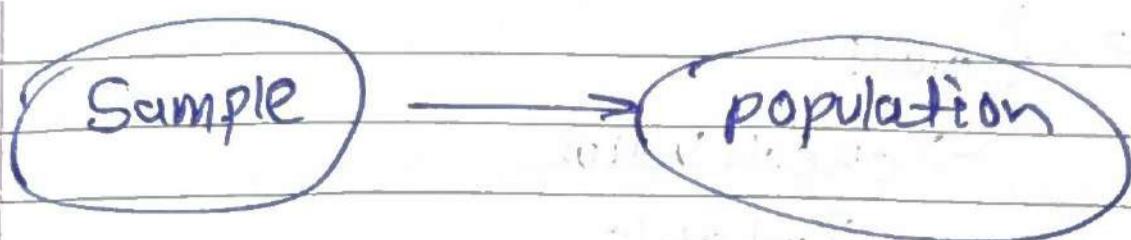
# Confidence Interval:

⇒ point estimate

The value of any statistics that estimates the value of a parameter.

Inferential stats

↳ based on sample data we take decision or estimate about the population.



Mean

$$\bar{x} \rightarrow 2.9$$

Mean

$$\mu \rightarrow 3$$

$\Rightarrow$  C.I.

Point estimate  $\pm$  margin of error

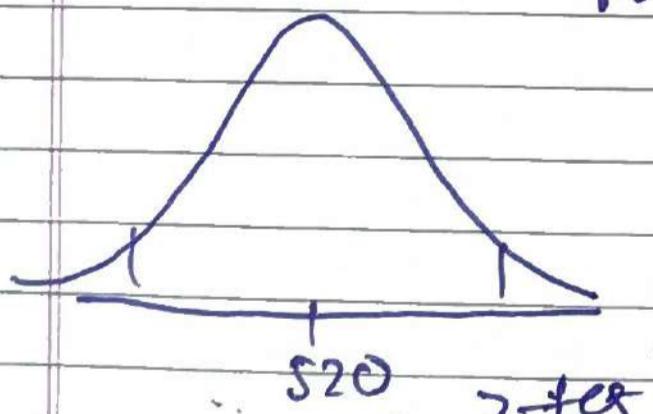
$$T = 100$$

$$n = 25$$

$$\alpha = 0.05$$

$$\bar{x} = 520$$

e.g.: On the quant of CAT exam the std is known to be 100. A sample of 25 test takers has a mean of 520 score. Construct a 95% CI about the mean.



③ Whenever population std is given, we apply a test:

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

$\hookrightarrow$  Standard error.

lower bound

$$\bar{x} - Z_{0.05} \frac{100}{\sqrt{25}}$$

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

↗ check value  
on Z-table.

$$1 - 0.025$$

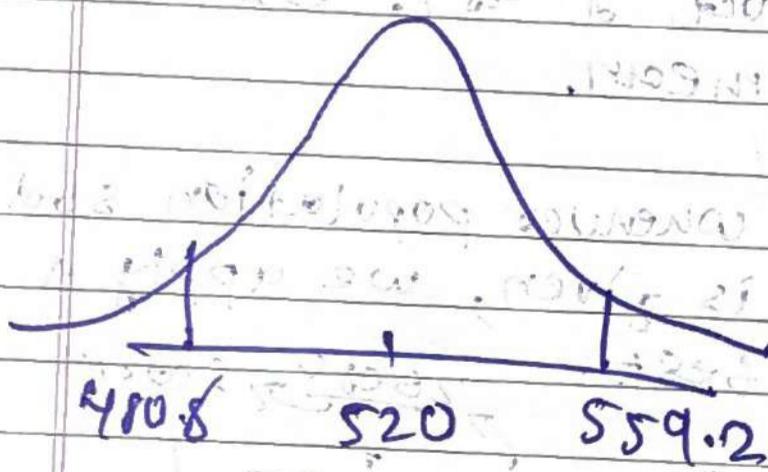
$$= 0.975$$

↗ 2-table  
1.96

$$\text{Upper} = 520 + 1.96(20)$$

$$= 559.2$$

$$\text{Lower} = 520 - 1.96(20) = 480.8$$



g] On the quant ~~that~~ test of CAT exam, a sample of 25 test taken has a mean of 520 with a std of 80. Construct 95% confidence interval about the mean.

$\Rightarrow$  No population standard dev given

$\hookrightarrow$  Here, we use t-test.

Point estimate = Margin of error

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

$$\text{Upper: } \bar{x} + t_{\alpha/2} \left( \frac{s}{\sqrt{n}} \right)$$

\* Degree of freedom  $D.F. = n - 1 = 25 - 1 = 24$

$$t\text{-table} \Rightarrow t_{\frac{\alpha/2}{2}} = 2.064$$

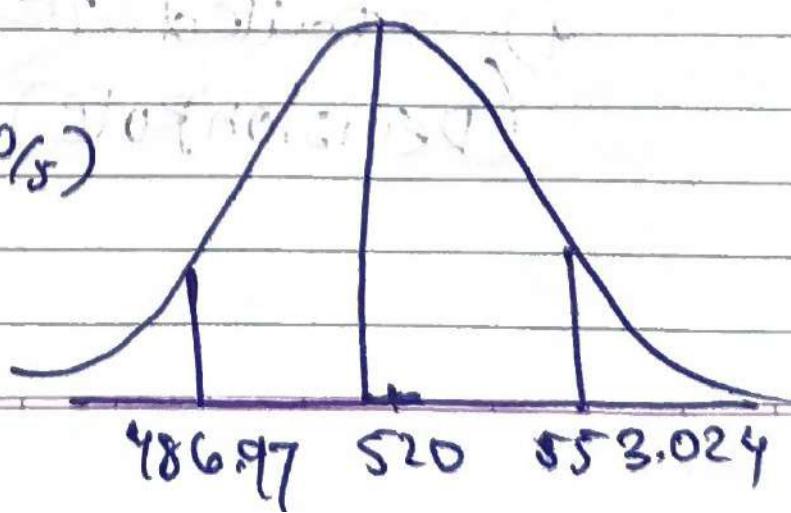
$\hookrightarrow$  from

$$\text{Upper} = 520 + 2.064 (80/\sqrt{5})$$

$$= 553.024$$

$$\text{Lower} = 520 - 2.064 (80/\sqrt{5})$$

$$= 486.97$$



1)

One Sample

( $Z$ -test)

① Population std is given

②

Sample size  $n \geq 30$

\* In the population, the average IQ is 100 with a std of 15.

Researchers wants to test a new medication to see if there is

positive or negative effect on intelligence, or no effect at all.

A sample of 30 participants who have taken the medication has a mean of 140. Did the medication affects the intelligence. ( $\alpha = 0.05$ ).  $T-I = 95\%$



1) Define Null ( $H_0$ )

$$\rightarrow \mu = 100 = H_0$$

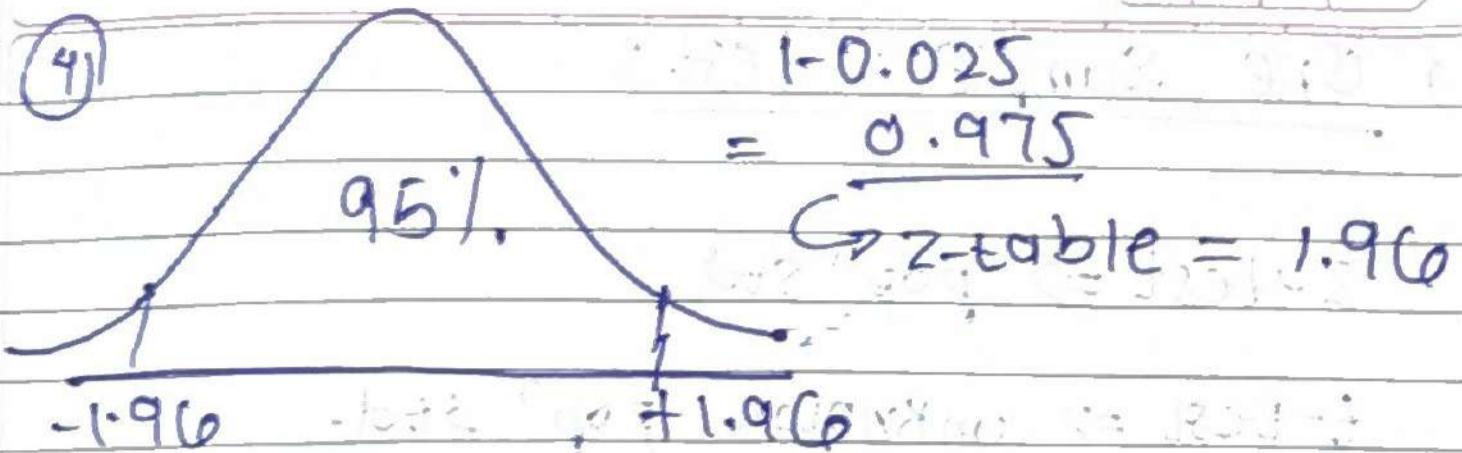
2) ( $H_1$ )

$$\rightarrow [\mu \neq 100 = H_1]$$

3)  $\alpha = 0.05$

4) 2 tailed test

(Decision Rule)



(5) Calculate Test Statistics (Statistics of the test)

$$Z = \frac{\bar{x} - \mu}{\text{standard error}} \rightarrow \text{Real formula}$$

$$\text{standard error} = \sqrt{\frac{\sigma^2}{n}}$$

before,  $n=15$   
so,  $\sqrt{\frac{\sigma^2}{15}} = \sigma$

$$= \frac{140 - 100}{\sqrt{\frac{15}{30}}} = \frac{40}{\sqrt{15}} = \frac{40}{3.87} = 10.34$$

State our Decision,

$$10.34 > 1.96$$

if  $Z$  is less than  $-1.96$  or greater than  $1.96$ , we reject null hypothesis.

Did the medication improved the intelligence?  
Yes!

## ② One Sample t-test:

Z-test  $\Rightarrow$  pop<sup>n</sup> std.

t-test  $\Rightarrow$  unknown pop? std.

g) Population avg IQ = 100

$n = 300$   $s = 20$

$\bar{x} = 140$

Did the medical affect the intellig?

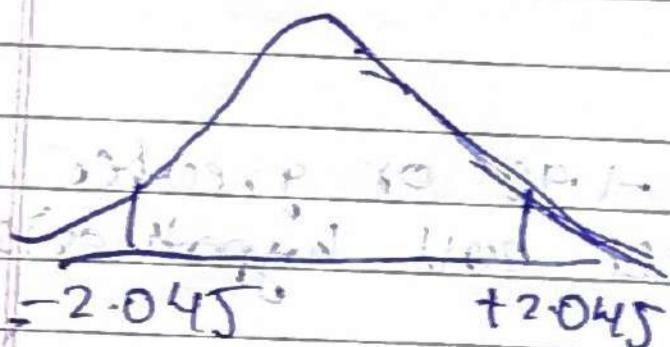
$\Rightarrow ① H_0: \mu = 100$

②  $H_1: \mu \neq 100$

③ degree of freedom

$n - 1 = 30 - 1 = 29$

④ state decision



(5)

### T-test

$$t = \frac{\bar{x} - \mu}{\frac{s}{\sqrt{n}}} = 10.96,$$

$t = 10.96 > 2.045$   
 ↳ we reject  $H_0$

$P \leq$  significance value

$H_0 X$

$H_1 \checkmark$

### # CHI SQUARE TEST :

↳ Chi square test claims that about population proportions it is a non parametric test that is performed on categorical (nominal or ordinal) data.

g] In the 2000 Indian census, the age of the individual in a small town were found to be the following!

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

?

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

⇒ What is non-parametric test?

→ A statistical that does not rely on any assumptions about the population distribution,

use's ;

when data is not normally distributed even small sample size.  
ordinal & nominal data!

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

$<18$	$18-35$	$>35$	$n=500$
121	288	91	Observed
$500 \times 0.2$	$500 \times 0.3$	$500 \times 0.5$	Expected

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

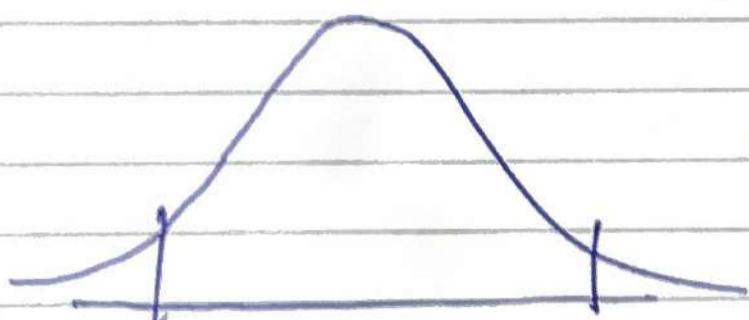
①  $H_0$  = The data meets the distribution  
2000 census

②  $H_1$  = does not meet .....

③  $\alpha = 0.05$  (95% CI)

④ Degree of freedom =  $n-1 = 3-1 = 2$

⑤ Decision Boundary



$$\text{Chi square} = \chi^2$$

2tailed test

if  $\chi^2$  is greater  
than 5.99 reject  $H_0$ .

5) Calculate Test Statistics.

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

$f_o$  = observed.

$f_e$  = expected.

$$= \frac{(121 - 100)^2}{100} + \frac{(288 - 150)^2}{150} + \frac{(914 - 250)^2}{250} = 232.94$$

$$\chi^2 = 232.94 > 5.99$$

{ Reject the null ( $H_0$ )

(2)

Covariance :

P-value &lt; significance value

 $H_0 X$  $X$  $Y$ Weight

50

60

70

75

Height

100

170

180

180

 $x \propto y$ 

According to dealer.

No. of hours  
Study

playing

2

3

4

0

4

3

 $x \propto \frac{1}{7}$ 

Covariance

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

Sample

$\text{Cov}(x, y) \Rightarrow +\text{ve}$

$$\downarrow \\ +\text{ve} \rightarrow [x \propto y]$$

directly proportional.

$\text{Cov}(x, y)$

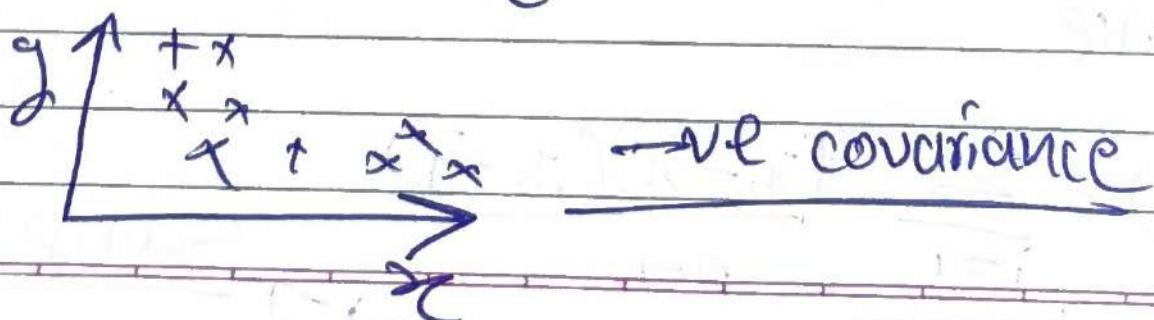
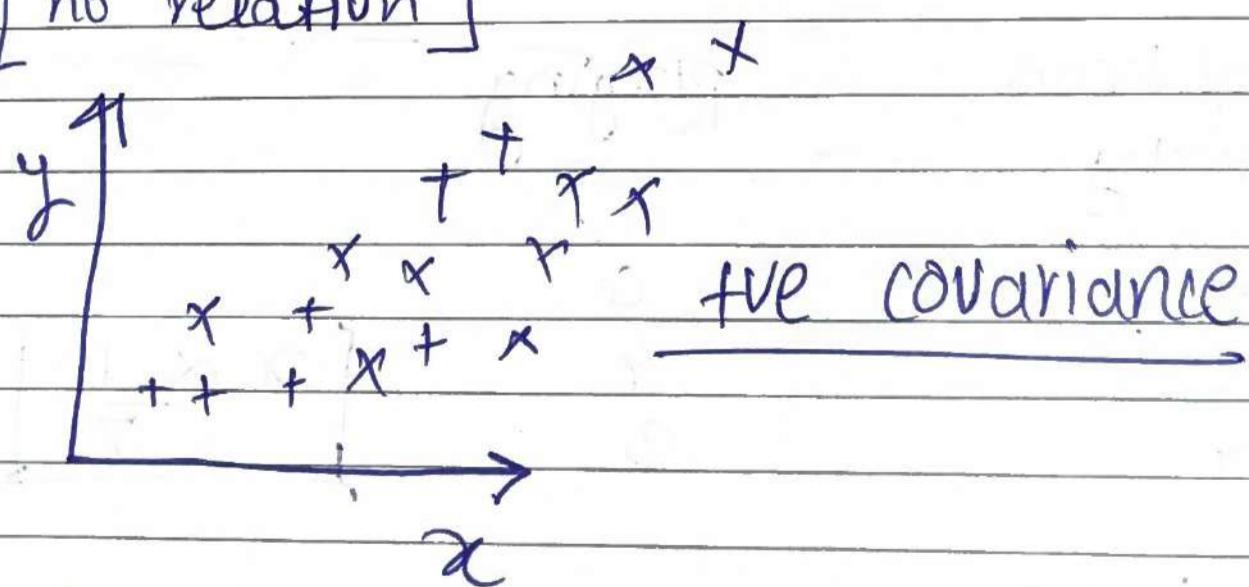
$$\downarrow \\ -\text{ve} \rightarrow [x \propto \frac{1}{y}]$$

$\text{Cov}(x, y)$

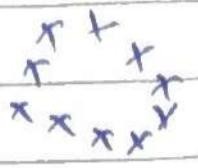
$$\downarrow \\ 0$$

$\rightarrow$  if  $x$  is increasing  
no effect on  $y$ .

[no relation]



$y$



$x$

Covariance = 0

No relation

### Disadvantage of Covariance

- ① There is no fixed range  
no limit to magnitude  
such as:  $(-200, +10, +1000)$

We have to restrict this value within some range.

②

### Pearson Correlation :

$\hookrightarrow (-1 \text{ to } +1)$

more negative side (more -ve correlation)

more on +1 side (more +ve correlation)

- is good at satisfying linear properties
- for linear Spearman correlation.

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

### ③ Spearman Correlation:

$$\text{Spear}(x, y) = \frac{\text{Cov}(R(x), R(y))}{R_{xy} \times R_{yj}}$$

↑ Rank

- For Non-linear Data.

What is  $R(x)$ ?

(G)

Weight (x)	Age Height (y)	<del>Age</del> Height (x)	$R(x)$	$R(y)$
75	170	2	2	2
62	160	3	3	3
60	150	4	4	4
55	145	5	5	5
<del>50</del>	8			
85	180	1	1	1

$P \leq 0.05 \rightarrow$  Reject the

$H_0$

## # P-value & Significance

g) avg weight of all residence in a city is 168 pounds with a std 3.9. we take a sample of 30 individuals and the mean is 169.5 pounds.

$$C-I = 95\%$$

→ Z-test

$$\mu = 168$$

$$\sigma = 3.9$$

$$\bar{x} = 169.5$$

$$n = 30$$

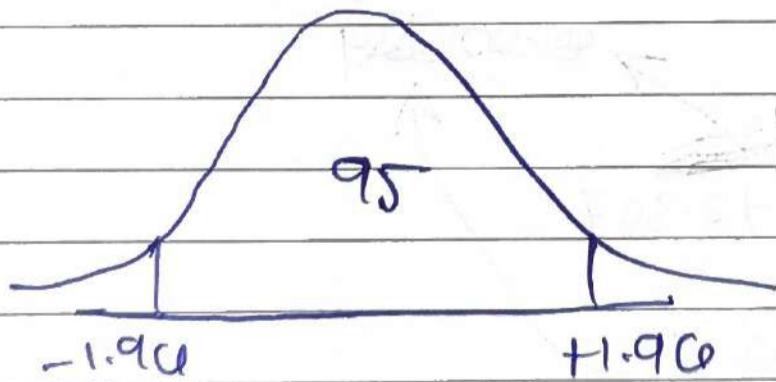
$$\alpha = 0.05$$

$$① H_0 = \mu = 168$$

$$② H_1 = \mu \neq 168$$

$$③ \alpha = 0.05$$

### ④ Decision



### ⑤ Z-test:

$$Z = \frac{x - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{169.5 - 168}{\frac{3.9}{\sqrt{30}}}$$

$$= \frac{1.5}{\frac{3.9}{\sqrt{30}}} \times \sqrt{30} = 2.307$$

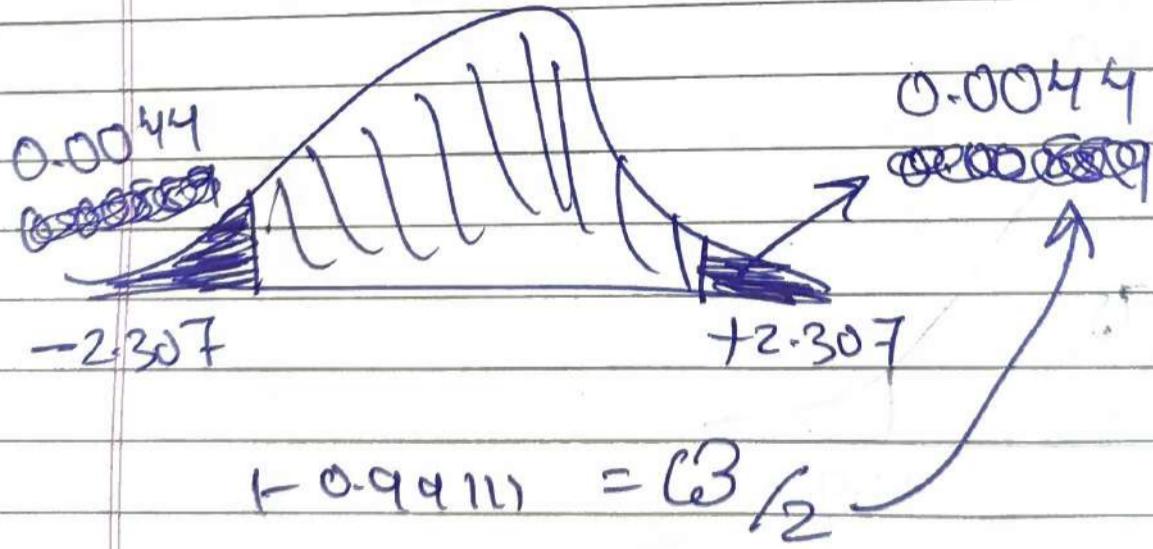
$$Z = 2.307 > 1.96$$

↳ True

↳ reject  $H_0$ .

$$2.307 \Rightarrow z\text{-table} \Rightarrow \underline{0.99111}$$

Area under  
the curve



$$1 - 0.99111 = 0.0088$$

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

$$0.0088 < 0.05$$

↳  $XH_0$  //

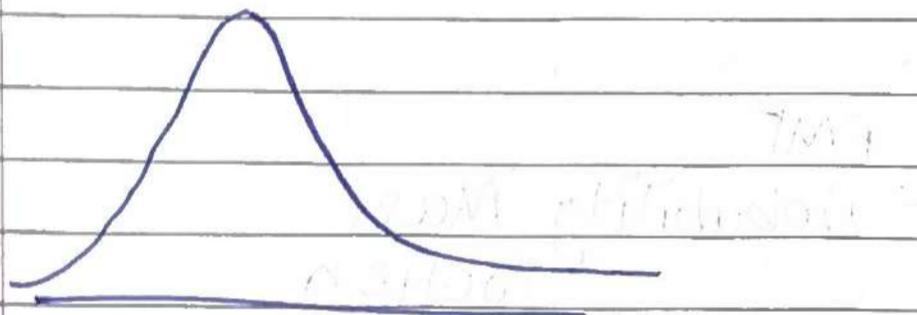
# P-value  $\geq$  significance

↳  $\times H_0$

P-value  $<$  significance

↳  $\checkmark H_0$

# log Normal distribution.



{  
y  $\sim$  lognormal  
dist}

$\downarrow$   
 $\log(y) \Rightarrow$  Normal  
dist

if y belongs to lognormal dist

it will satisfy  $\log(y) \Rightarrow$  will converted  
into normal distribution.

## # Bernoulli's Dist:

↳ 2 outcomes

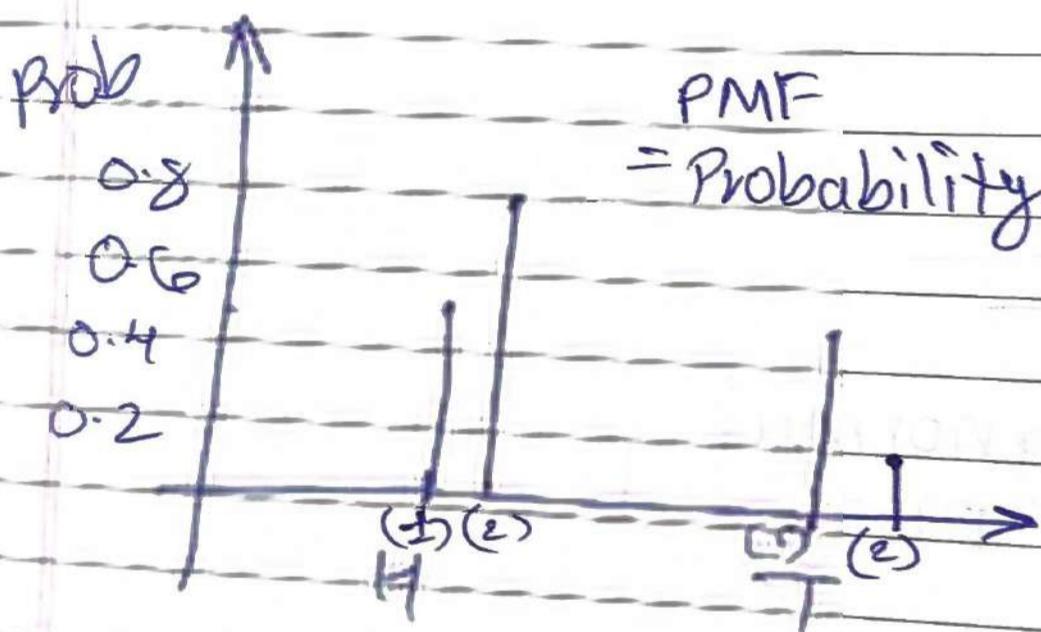
↳ (0 or 1)

e.g. Tossing a coin

$$P(H) = 0.5 = P$$

$$q = 1 - p$$

$$\boxed{q = 0.5}$$



PMF  
= Probability Mass function

Pdf  $\Rightarrow$  for continuous variables,  
PMF  $\Rightarrow$  for categorical variables

## # Binomial Distribution:

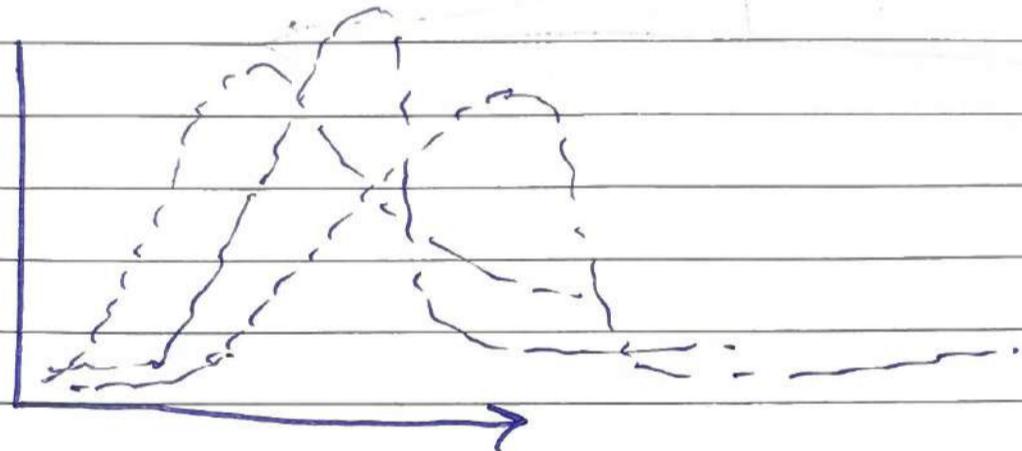
Every trial  $\rightarrow$  Bernoulli dist

multiple trials combined

combination Bernoullis distribution.

$$B(n, p)$$

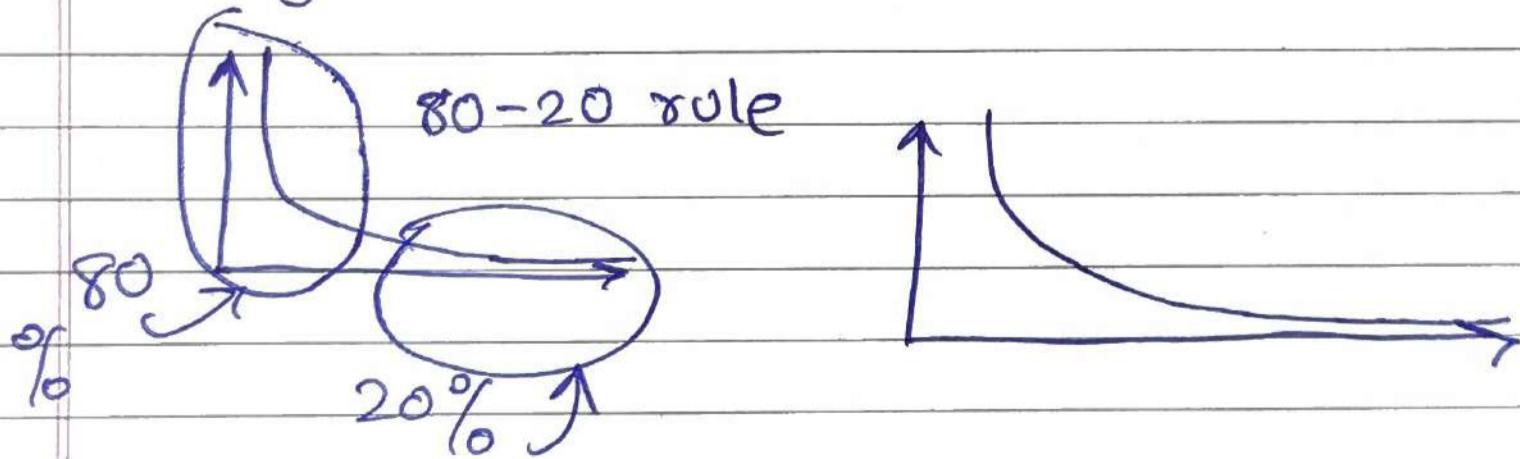
↓      → success prob  
no. of      for each trial  
trials



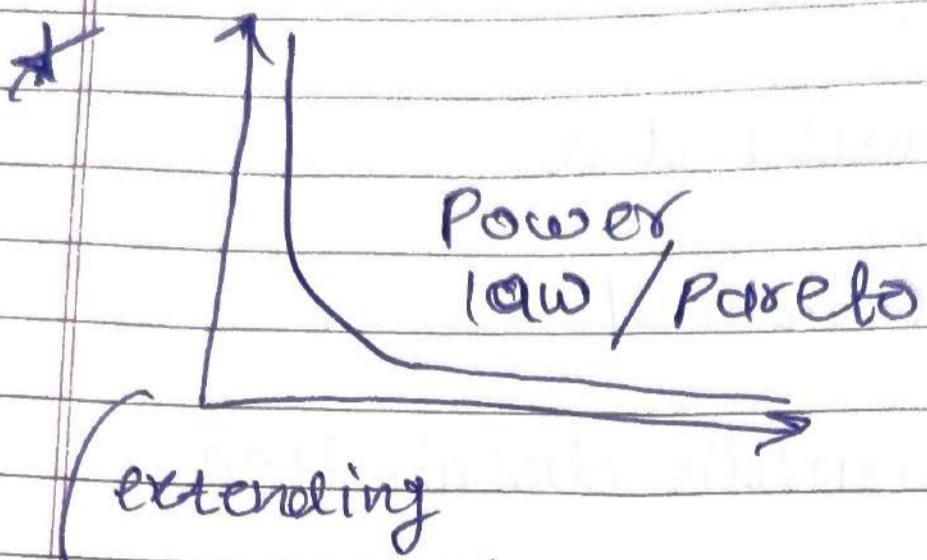
## # Pareto dist:

→ Power law

Non-gaussian dist.

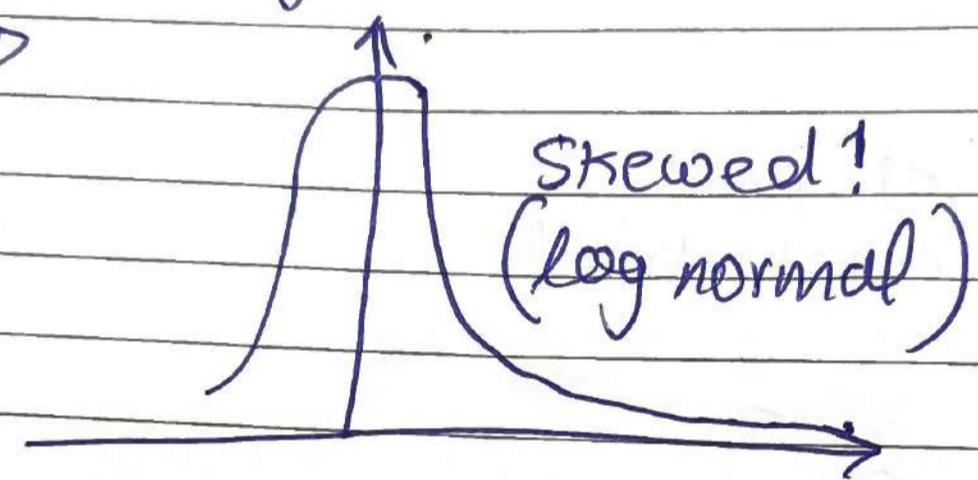


e.g.: 80% of wealth is distributed in 20% of people.



Power  
law / Pareto

extending



Skewed!  
(log normal)