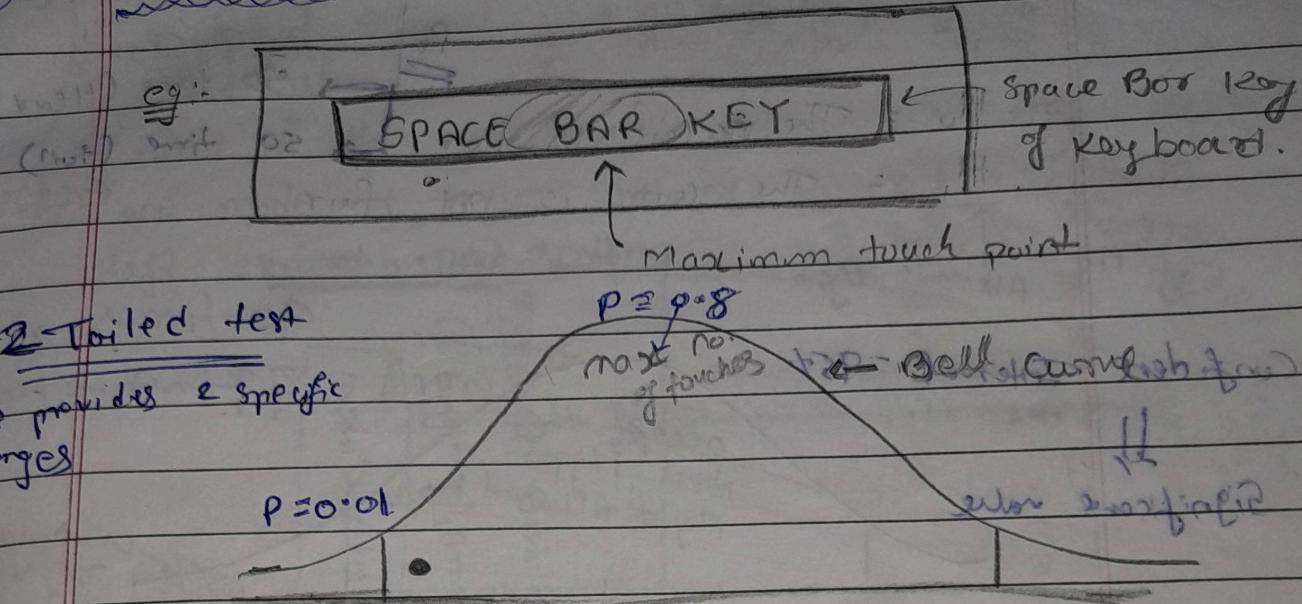


(P) \* P-value (Probability value)



$P = 0.01$

↳ This  $P = 0.01$  represent that,

If we repeat this experiment 100 times, then the number of time, ~~the number of times~~ <sup>we will be</sup> touching this area will be  $\Rightarrow$  1 times ~~below - 9~~.

$P = 0.05$   $\Rightarrow$  less than 5 times below - 9

Some explanation  $\Rightarrow$  80 times hardly in

80 times hardly in

Define  $\Rightarrow$  P Value  $\Rightarrow$  It is the probability of the ~~null hypothesis~~

"Null Hypothesis" to be true. ~~is significantly~~

whether it is right or wrong in no of observations

Null Hypothesis  $\Rightarrow$  Treats all everything

Same or Equal



P value  $\xrightarrow{\text{we also called}}$  Significant value.

Eg: Coin  $\Rightarrow$  Fair coin  $\rightarrow$  (Test) Experiment 100 times

Null Hypo:  $H_0$   $\leftarrow$   $P(\text{Head}) = \frac{1}{2}$  almost sure  
always  $\therefore H_0$  = The coin is fair  
opposite of  $H_0$   $\uparrow$  Alternative Hypo

Pr of Head  $\approx$  almost sure

Pr

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

so time (Head)

so time (tail)

$\Pr(T) = \frac{1}{2}$

perform Experiment =

40 times  $\rightarrow$  Head

Confidence Interval = 95%

$\alpha$  is determined by

Significance value

accept the

Null Hypo

Pr of Alternative Hypo

Symmetric graph

97 times  $\rightarrow$  tail

we have to reject

10

50

50

60

60

70

70

80

80

we should always get this value nearer to

### \* P-value

$\Rightarrow$  P-value is a measure of the probability that an observed difference could have occurred just by random chance. The lower the p-value, the greater the statistical significance of the observed difference. P-value can be used as an alternative to or in addition to pre-selected confidence levels for hypothesis testing.

(14)

## \* Confidence Interval

point estimate; to

basically calculate the parameter ( $\mu$ )

This is v. important  
in order to find out

① Point Estimate:

$$\bar{x}$$

$\mu$

With the help of prop. we will try to obtain  
estimates the parameters of population mean

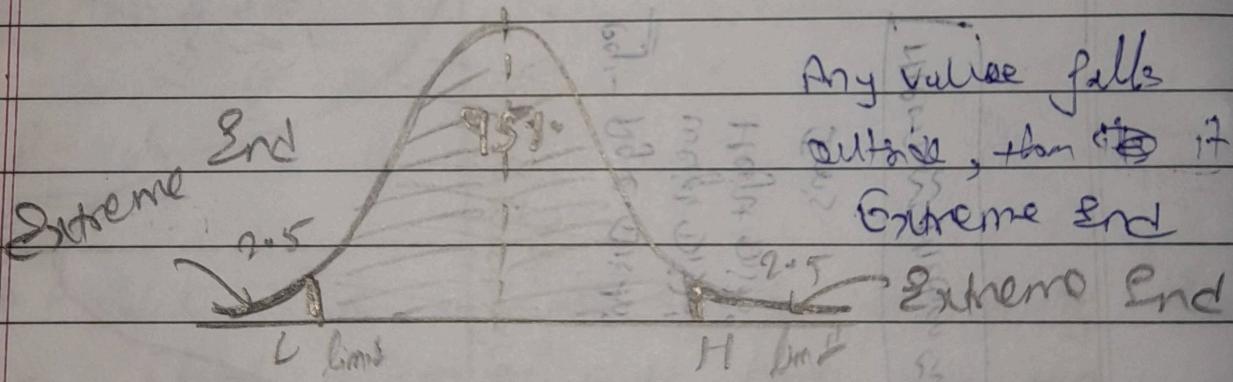
(ii)

with some  $CI = 95\%$ , i will able to find out the range of the population means from [L to H].

95% {  
CI range  
Lower limit

Higher  
limit

This is Basically my range that i can get  
w.r.t the population mean ( $\mu$ )



Any value falls between Lower limit and  
Higher limit, they are basically CI

(15)

\* Bernoulli distribution - Mean, Variance and Standard Deviation

~~In BD will have only~~

The outcome ~~would be~~ ~~BD~~ will always 2  
outcome = 2

\* Eg: Tossing a coin

$$\begin{array}{c} \text{outcomes} \\ \hline \text{H, T} \end{array} \quad \begin{array}{l} P(H) = 0.5 \\ P(T) = 0.5 \end{array}$$

\* Suppose if I say that whenever my random variable gets a tail,

$$\underline{X = 1} \rightarrow (\text{success})$$

If my random variable has zero this will basically be my fail,

$$\underline{X = 0} \rightarrow (\text{fail})$$

for this the

$$P(X=x) = p^x (1-p)^{1-x} \quad \left. \begin{array}{l} \text{Probability Mass} \\ \text{function} \end{array} \right\}$$

$\uparrow$   
random variable X will have some value (x)

Note:

In Probability Density function  $\xrightarrow{\text{outcomes in}} \text{Continuous Value}$

In Probability Mass function  $\xrightarrow{\text{outcomes in}} \text{Discrete Value}$

Only

two outcomes  $\rightarrow P$

$$1-P = q$$

a will usually have,  $\underline{x=1}$  & 1

P of getting success

$$\boxed{P(\text{success}) = p}$$

P of getting failure

$$\boxed{P(\text{failure}) = 1-p = q}$$

based on this two probabilities are defined Bernoulli Distribution.

+ Derive  $P(X=0) = p^0 (1-p)^{1-0} = (1-p) \Rightarrow q$

$$P(X=1) = p^1 (1-p)^{1-1} = p (0=1)$$

$$\text{PMF} = \begin{cases} q = 1-p & \text{if } x=0 \\ p & \text{if } x=1 \end{cases}$$

$$\therefore \text{PMF} = p^x (1-p)^{1-x}$$

e.g. fair coin

$$\left\{ P(H) = p = \frac{1}{2} \right.$$

$$\left. P(T) = 1-p = 1-\frac{1}{2} = \frac{1}{2} \right.$$

$$\left\{ P(H) = 0.4 = p \right.$$

$$\left. P(T) = 1 - 0.4 = 0.6 \right.$$

"if have one Experiment make"

$\Rightarrow$  In BD we have values like  $p$  and  $q = 1-p$

Always Values between  $[0 \text{ } \& \text{ } 1]$

Mean, Variance & SD

out with measured  
bright end  
Expected value graph

$$\text{Exp. value} = \sum_{i=1}^q x_i \cdot p(x)$$

Q1.  $x = 0 \text{ or } 1$

$$P(X=0) = 0.4 \Rightarrow 1 - P$$

$$P(X=1) = 0.6 \Rightarrow p$$

$$Ex = \sum_{i=1}^q x_i \cdot p(x)$$

$$= 0(0.4) + 1 \times (0.6)$$

$$= \underline{\underline{0.6}}$$

Variance ( $\sigma^2$ )

$$P(1-p) = Pq$$

$$\sigma(\text{SD}) = \sqrt{Pq} = \sqrt{0.4 \times 0.6} = \sqrt{0.24}$$

Mean

P

Variance

Pq

SD

$\sqrt{Pq}$

(16)

## 5 Number Summary And How to handle outliers using IQR [removing the outlier]

- ① 1) Minimum → 1
- ② 2) 25% → 3 first Quartile ( $Q_1$ )
- ③ 3) Median → 6 (Central Element)
- ④ 4) 75% → 9 Third Quartile ( $Q_3$ )
- ⑤ 5) Maximum → 10

### Percentiles

value which is my  $50\%$  element  
 50% of this entire distribution  
 is less than Grade  
 1, 2, 3, 4, 5, 6, 7, 8, 9, 9, 10  
 ↪ 25% value      ↪ 75% value

⇒ Calculate 25%

$$\text{Percentile} = \frac{25}{100} \times (n+1) \quad | \quad n = \text{no. of Elements}$$

$$= \frac{25}{100} (13) = 3.25 \approx 3$$

Calculate 75%

$$\text{Percentile} = \frac{75}{100} (n+1)$$

$$= \frac{75}{100} (13) = 9.75 = 10 \approx 10$$

$$IQR = Q_3 - Q_1$$

$$IQR = 9 - 3 \Rightarrow 6$$

Dataset  $\Rightarrow$  1, 2, 3, 4, 5, 5, 6, 7, 8, 9, 9, 10, 12

- ↳ [lower bracket]  $\rightarrow$  Higher bracket
- ↳ whichever values within this particular distribution follows between this lower bracket and higher bracket only those elements will be kept
- ↳ Remaining all the elements that are outside of this particular bracket will be treated as outlier.

$$\text{Lower Bracket} = Q_1 - 1.5(IQR)$$

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

$$LB = 3 - 1.5(6)$$

$$= -6$$

$$HB = 9 + 1.5(6)$$

$$= 18$$

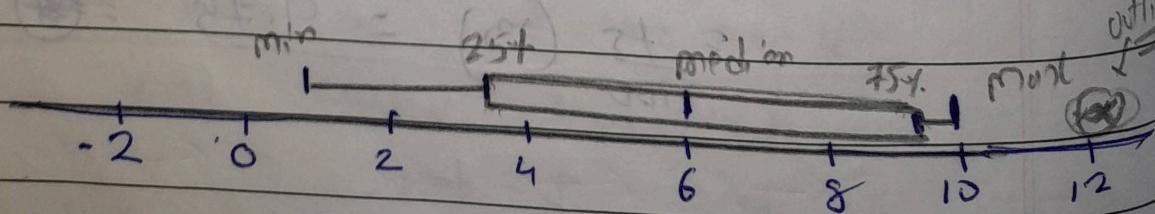
LB

(-6)

HB

18

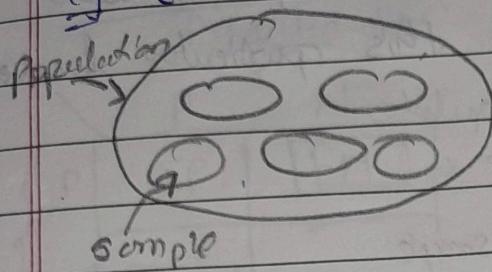
Elements not falling in between can be outliers



## (17) \* Types of Sampling Techniques

### (1) Random Sample Techniques

e.g. Ballot Paper Selection)



from Population we take sample  
Party A →  
B →  
C →

Here we are randomly select some kind of samples and we some kind of conclusions w.r.t the population.

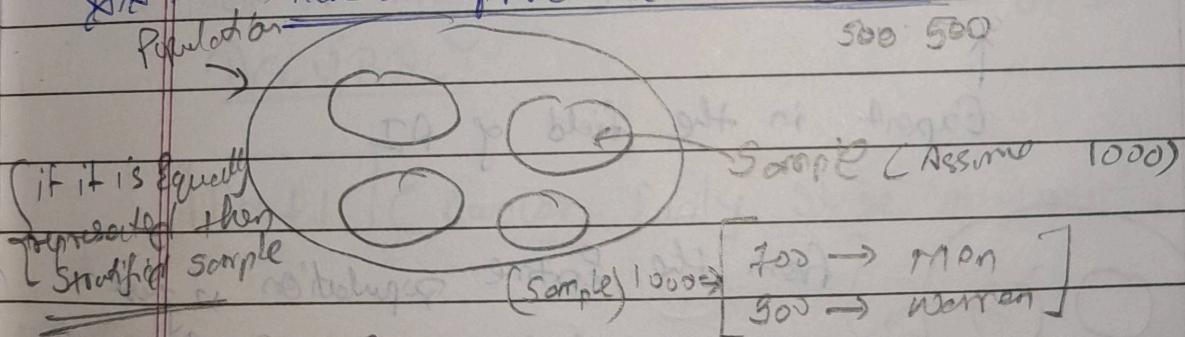
~~Design~~ Layers  $\rightarrow$  Groups

~~Design~~ They are non overlapping groups

### (2) Stratified Sampling Technique

⇒ Stratified sampling technique, we know in a technique where we basically train our model we do the train-test-split in stratified manner

~~Design~~ → Here we provide ratio in [1:1]



Ratio of Gender

in [7:3] [There should be equal proportion of Gender]

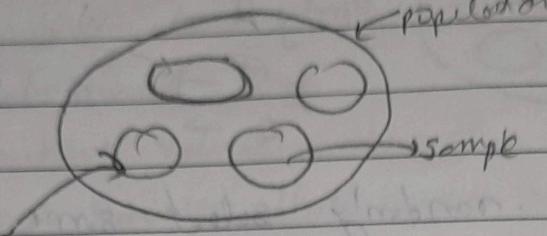
If we try to pick up the sample with this particular ratio, this may become completely Biased

Because, if this we may not get proper output or proper inferences based on this entire population

### ③ Systematic Sampling Techniques

18

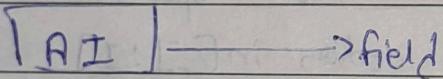
- Here we are just selecting the  $n^{\text{th}}$  person  
and every  $n^{\text{th}}$  person will try to select  
 $10^{\text{th}}$  person i.e. in this particular location



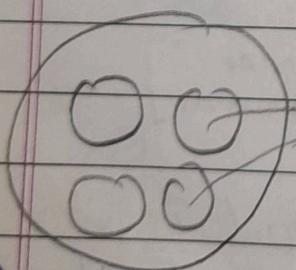
Selecting this particular location and every 10<sup>th</sup> person that i see i will try to record the statement, like which party this particular guy has actually voted.

### ④ Cluster Sampling Technique → Domain

Taking some in the AI domain



↑  
Expert in the field of AI



{ from the entire population, i just pick up people who are familiar with this domain knowledge

Live class

④

### Convenience Sampling (or) { Voluntary Response Sampling }

Survey → 1000 people → No → Interested → All  
generated at Survey and mailed to 1000 people  
only interested people will fill the form

e.g.: A survey on Privacy

18

Reduce false Positive (or) Negative confusion matrix

Date : \_\_\_\_\_  
Page : \_\_\_\_\_

Confusion Matrix

Predicted data

		P	N	
		TP	FN	Type 2 error
Actual data	P	FP	TN	Type 1 error
	N			

9 use cases

- we should always try to reduce FN & FP
- ① whether the person has a disease or not
- ② whether the market will crash or not
- ③ Vaccination side effect or not

\* Most of time the Matrix we are going to use is

## ACCURACY

| FN ↓ | (Cancer Early Stage treatment, covid)

use case

- ① Person has a disease or not

→ we should always try to reduce FN ↓

In this use case (cancer)

TP ↑ TN ↑ FN ↓ FP ↓

- ② Market will crash or not

→ Aim: Person should be able to save his money. FN ↓ we should try to reduce.

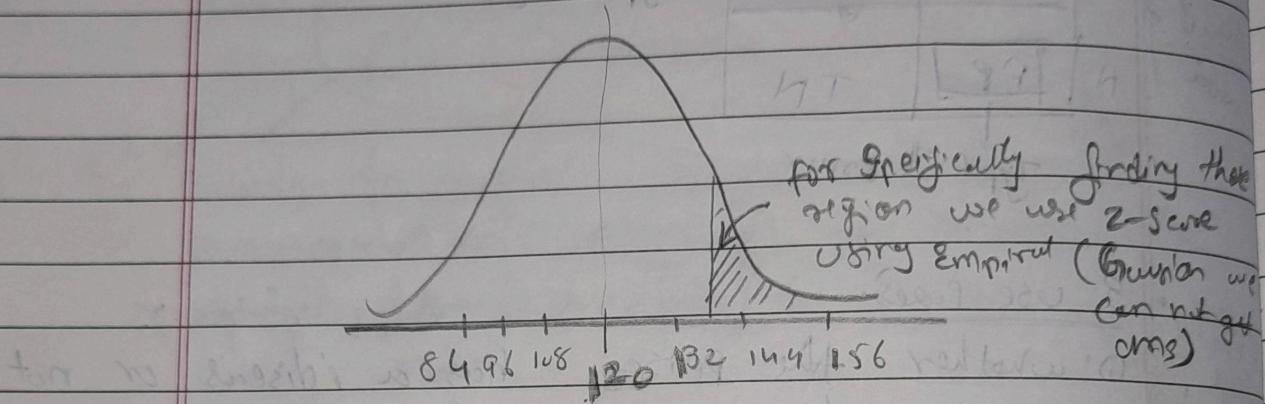
19

## Vaccination site effect

### Z-Score And Its Applications.

$$\mu (\text{population mean}) = 120$$

$$\sigma (\text{standard Deviation}) = 12$$



→ Z-score will help us to find out whether a value like how far this value is from the mean in the terms of standard deviation.

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

$$= 144 - 120$$

$$12$$

How much SD away from the mean

### Application ① Standardization

② Helps us to compare scores between different distribution

~~India perform good in 2021~~

Eg:- India Cricket Team

2020 (Cricket)

Assume

$$\text{Avg} = 180$$

$$\sigma = 12$$

$$\text{final score} = 187$$

2021 (Cricket)

$$\text{Avg} = 182$$

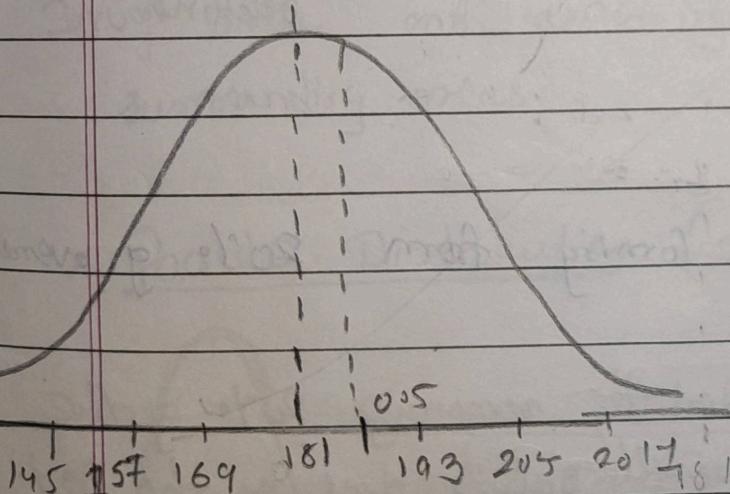
$$\sigma = 5$$

$$\text{final score} = 185$$

$$Z_{2020} = \frac{187 - 180}{12}$$

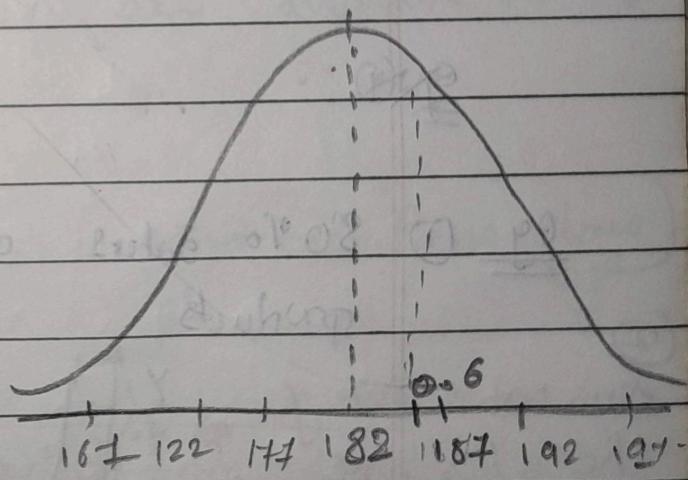
$$= 0.5$$

=



$$Z_{2021} = \frac{185 - 182}{5}$$

$$= 0.6$$

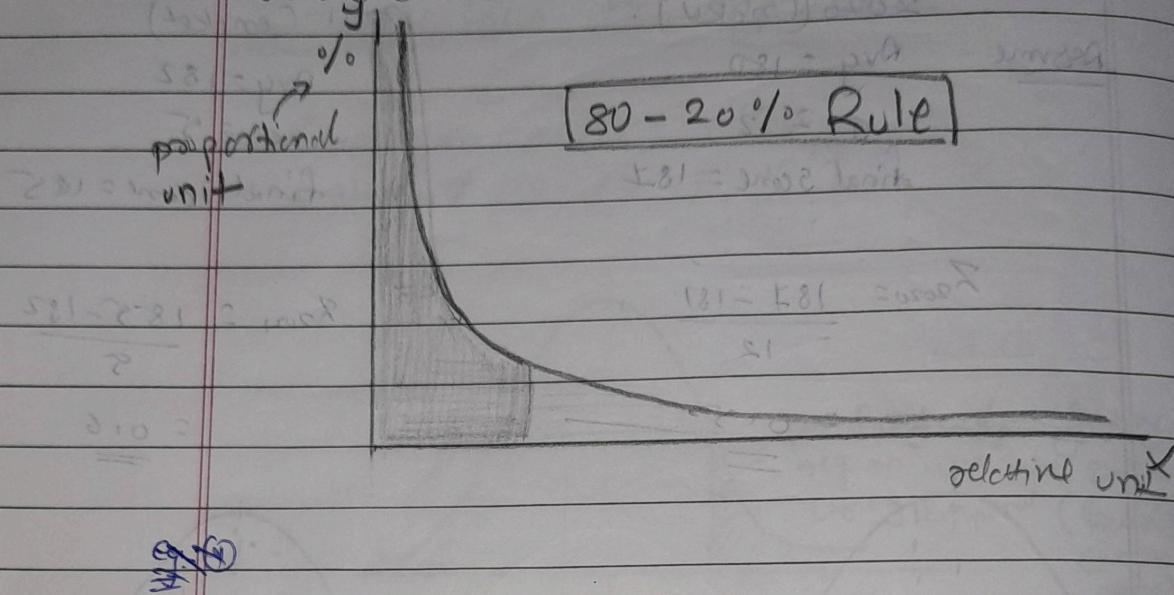


~~A power law~~ is the relationship between two quantities where in a relative change in one quantity results in a proportional change in the other quantity.

Date : \_\_\_\_\_  
Page : \_\_\_\_\_

(20)

## Power law Distribution And Its Examples and Applications



e.g. ① 80% sales are coming from 20% of overall products

② This entire 80% of the sales are done by the 20% of the overall products.

③ And the remaining 20% of the overall sales done by the 80% of the overall products.

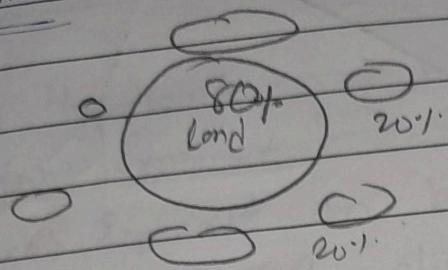
④ This 20% of the overall products is actually doing the 80% of the sales.

⑤ 80% of windows crash are because of 20% of all the overall bugs.

⑥ 80% of Data Scientists uses 20% of overall software products.

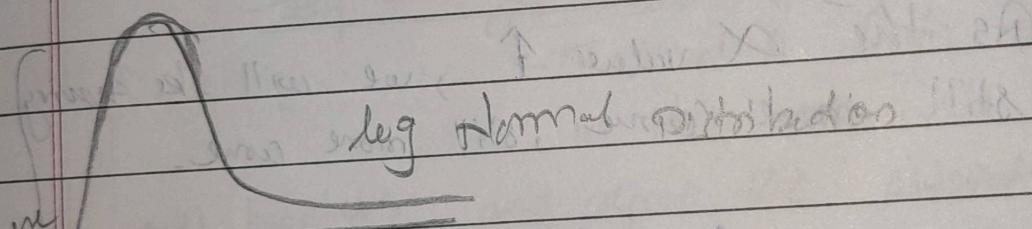
④ 80% of the IPL matches are won by 20% of the team

⑤ Oil fields



80% of the oil fields are available in this 20% of area and remaining 20% are available on the surrounding areas

\* Pareto Distribution (It is a kind of Power law distribution)



Difference between a normal distribution and a Pareto distribution:  
Normal distribution is symmetric and bell-shaped.  
Pareto distribution is skewed to the right, with a long tail extending to the right.  
In a Pareto distribution, the probability of an event occurring is proportional to  $x^{-\alpha}$ , where  $\alpha > 0$ .  
The mean of a Pareto distribution is finite if and only if  $\alpha > 1$ .

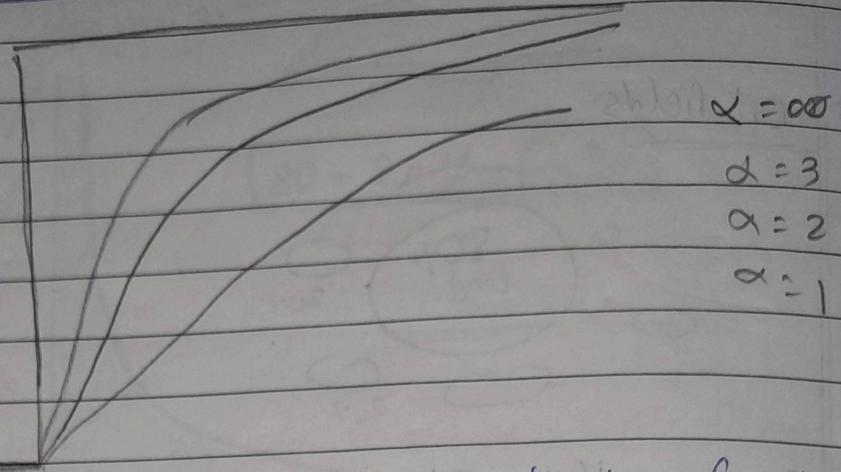
Pareto Type I probability density function for variables  $X$  with  $x_m = 1$ . As  $\alpha \rightarrow \infty$ , the distribution approaches a  $\delta(x - x_m)$  function, which is the  $\delta$ -function.

Second approach:  $f(x) = \frac{1}{x^{\alpha+1}}$ , where  $\alpha > 0$ . This is the right-hand side of a delta-function. This is a probability density function because  $\int_{-\infty}^{\infty} f(x) dx = 1$ .

Properties of the Pareto distribution:  
1. The distribution is unimodal and right-skewed.  
2. The mean is finite if and only if  $\alpha > 1$ .  
3. The variance is finite if and only if  $\alpha > 2$ .  
4. The median is given by  $x_m \cdot \sqrt{e^{\alpha}} = x_m \cdot \sqrt{e^{\alpha}} = x_m \cdot \sqrt{e^{\alpha}} = x_m \cdot \sqrt{e^{\alpha}}$ .

## Cumulative Distribution function

(2)



Passes Type I cumulative distribution functions for  
Various  $\alpha$  with  $x_m=1$

Explains It basically is ~~cumulative sum~~ cumulative Summation we just do it Cumulative Summation w.r.t the probability.

As the  $\alpha$  value  $\uparrow$ , we will be having still more come, still more come.

Hypothesis Testing :- Hypothesis testing in statistics is a way for you to test the results of a survey or experiment to see if you have meaningful results. You're basically testing whether your results are valid by figuring out the odds that your results have happened by chance. If your results may have happened by chance, the experiment won't be repeatable and so has little use.

- Null Hypothesis = is often an initial claim that is based on previous analyses or specialized knowledge.
- Alternative Hypothesis = States that a population parameter is smaller, greater or different than the hypothesized value in the null hypothesis.

How to perform Hypothesis Testing - C.I.  
Z-Test Statistics, Derive Conclusion

## Hypothesis Testing

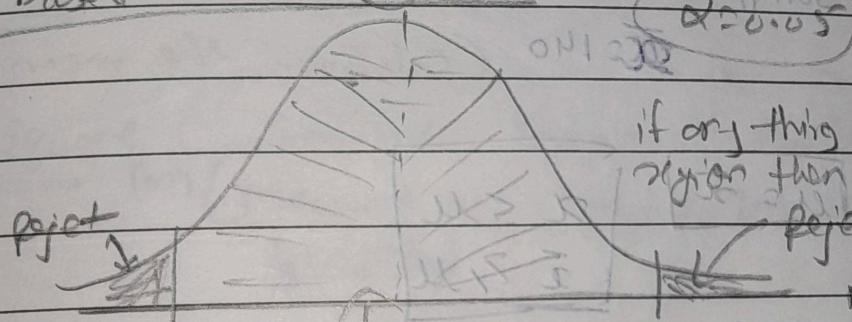
① Null Hypothesis  $\rightarrow$  Alternate Hypothesis

② Perform Experiment  $\xrightarrow{\text{with}} \alpha = 0.05$

Various Experiment can be performed like,

T-test, Z-test, chi-square test, ANOVA)

C.I based on the  $\alpha$  value, it is defined on the Domain Expertise



if anything falls in this region then we reject

Null Hypothesis

The P-value falls in this area well accept the

null hypothesis and reject the alternate hypothesis

$$0.05 > 3\sigma$$

$\rightarrow$  Test value falls in extreme end  $\rightarrow$  Reject

the Null Hypothesis, Accept the Alternate Hypothesis

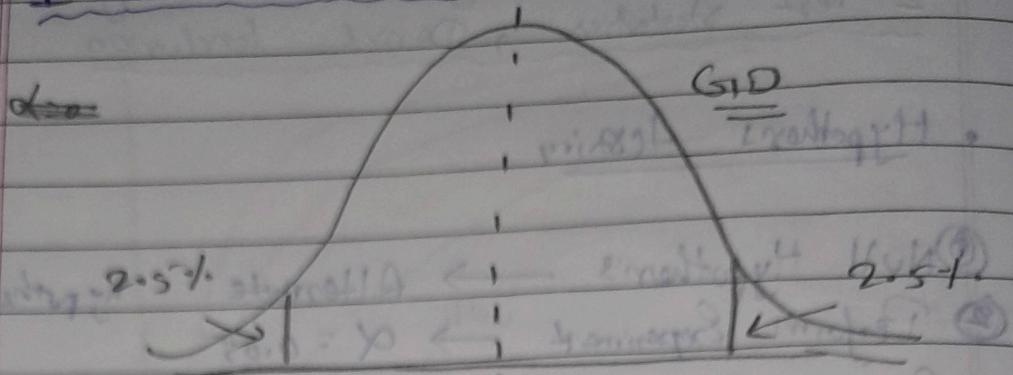
$C.I \rightarrow$  is defined based on the

\* Test value fall in Confidence Interval  $\rightarrow (\alpha)$

$\rightarrow$  Accept the Null Hypothesis and

$\rightarrow$  Reject the Alternate Hypothesis

1 tail test      2 tail Test



(e.g.  $N=100$ ,  $\mu=120$ ,  $\sigma=5$  } Population)

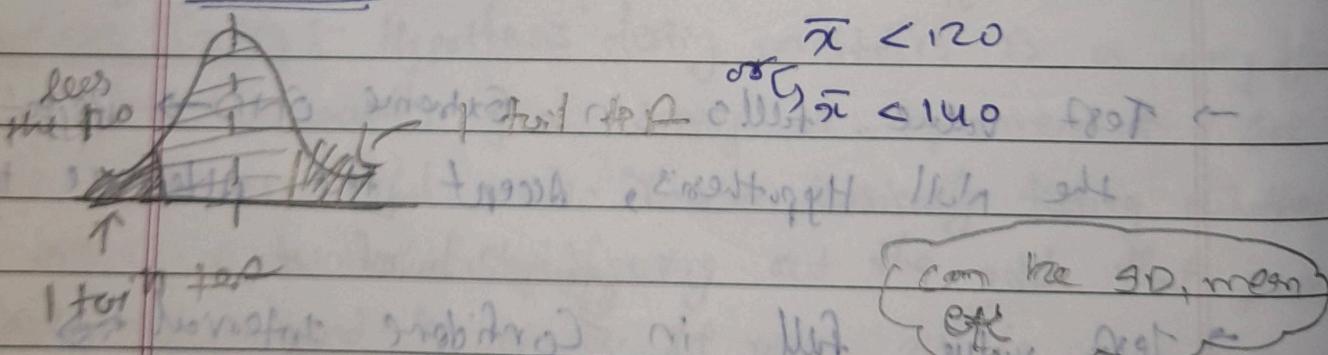
Sample  $\Rightarrow$  30

$\bar{x}=140$  (calculated the mean)

$$\therefore \boxed{\mu = \bar{x}} \quad \begin{cases} x \leq \mu \\ \bar{x} \geq \mu \end{cases}$$

$\therefore$  There kind of test is called as 2 tail test

### \* 1 Tail test



Population  $\rightarrow$  with some information

Properties

| Sample from population

Testing

Sample  $\rightarrow$  perform some experiment } Statistical testing

| Conclusion about the population