

" Statistics is the  
Grammer of  
Data Science / Data Analytics "

Hand Written Notes :

"Statistics Required for Data Analysis "

Ravindranath Kumbhar.

# \* Statistics

statistics is the science of collecting, organizing and analyzing data.

Data :- "facts or pieces of information"

e.g. Age of students in classroom

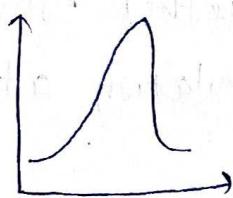
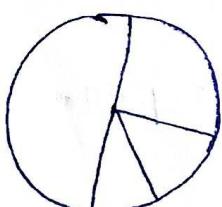
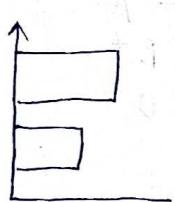
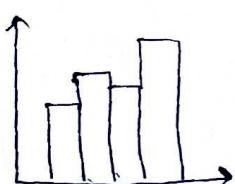
{ 24, 23, 22, 30, 32 }

statistics is divided into two types.

## Statistics

### Descriptive stats

- It consists of organizing and summarizing of the data in the form of histogram, Bar chart, Box plot, P.d.f, Pie i.e.



### Inferential stats

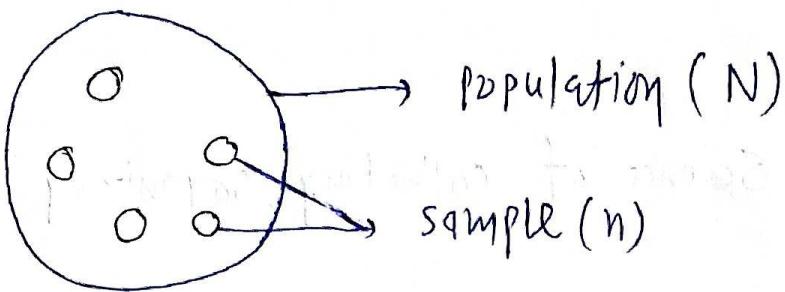
- It consists of technique to form some conclusions.  
e.g. Suppose there are 1000 students in a college and we have to measure average height of the students.

then we take some samples from this students and find the ~~average~~ average height.  
i.e.

using the sample size we can find the conclusion it belongs to inferential statistics.

- It contains to find mean, median, mode, S.D.,

## \* Population and sample :-



## \* Sampling techniques :-

① Simple random sampling :- Every member of a population (N) has an equal chance of being selected for your sample (n).

e.g.

when we are tossing a coin then chances of probability of getting head and tail is equal i.e.  $\frac{1}{2}, \frac{1}{2}$

② Stratified Sampling :-

Strata  $\Rightarrow$  Layers

It consists of different types of layers and this layers are non overlapping with each other.

e.g. ① Gender  $\rightarrow$  male

$\rightarrow$  female

② Blood group

$\rightarrow A^+$   
 $\rightarrow B^+$   
 $\rightarrow O^+$

③ Systematic sampling :-

It is the probability of sampling method where researchers select members of the population at a regular interval.

e.g. Suppose we have a list of total population then in systematic sampling we select every 5<sup>th</sup> person in the list or every 7<sup>th</sup> person in the list or so on.

④ Convenience Sampling :-  
only those people who are interested will only  
be participating.

It is used in marketing.

e.g.

Suppose a company launches a new product they send a notification only to those who are interested.

\* Variables :-

A variable is a property that can take on any value.

e.g. Age = 24

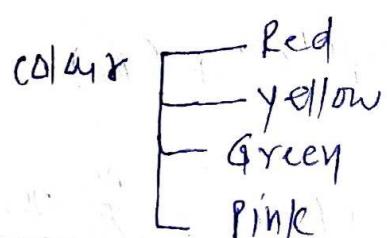
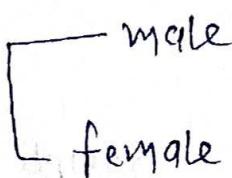
Height = 7.2

There are two kinds of variables.

① Quantitative variable :- measured numerically

② Qualitative variable :- categorical variable

e.g.: Gender



(Based on some characteristics we can group  
categorical variables)

# Quantitative Variable

↓  
Discrete variable

They are whole numbers

e.g. ① No. of bank accounts

② No. of children in family

• It does not contain fraction values

continuous variable

It is numeric

e.g. Height = 120.5 cm

weight = 50.5 kg

rainfall, temperature.

## \* Intermediate Stats :-

### \* Measure of central tendency :-

A measure of central tendency is a single value that attempts to describe a set of data by identifying the central position.

#### ① Mean (Average) :-

Population ( $N$ )

Sample ( $n$ )

$$\text{Popn mean} (\bar{\mu}) = \sum_{i=1}^N \frac{x_i}{N}$$

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

$$\text{e.g. } x = \{1, 2, 3, 4, 5\}$$

$$\bar{\mu} = \frac{1+2+3+4+5}{5}$$

$$\bar{\mu} = 3$$

## ② Median :-

let e.g.  $\{1, 2, 3, 4, 5, 100\}$

If there is any outlier is present in the set then  
steps to find the median :-

i) Sort the numbers

ii) find the central number.

① If no. of elements are even we find the average of central elements

② If no. of elements are odd we find the central element.

### \* Outlier :-

outlier is the number that is completely different than the entire distribution.

E.g. let consider the above e.g.

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

central

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

$$\therefore \text{median} = 3.5$$

## ③ Mode :-

It is the most frequent occurring element

e.g.  $\{4, 3, 1, 4, 4, 5, 6, 4, 8\}$

We select the element which has maximum frequency

i.e. 4

## \* Measure of Dispersion :-

### ① Variance ( $\sigma^2$ ) :-

It is a measure of how far a set of numbers is spread out from their average value.

for the population data (N)

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

e.g. let consider the set  $\{2, 2, 4, 4\}$

$$\therefore \text{mean } (\bar{u}) = \frac{12}{4} = 3$$

$$\therefore x_i - \bar{u} \quad (x_i - \bar{u})^2$$

$$2 \quad 3 \quad 1$$

$$2 \quad 3 \quad 1$$

$$4 \quad 3 \quad 1$$

$$4 \quad 3 \quad 1$$

$$\therefore \sigma^2 = \frac{4}{4} = 1$$

for the sample data (n)

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

\*  $n-1$  is less than  $n$ .

## A Standard deviation :-

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

It is nothing but root of variance.

The standard deviation is a measure of the amount of variation or dispersion of a set of values. A low standard deviation indicates that the value tends to be close to the mean of the set.

## \* Percentile and Quartiles :-

### ① Percentile :-

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

Suppose,

a person A got 99 percentile it means that the person has got better marks than 99 percentage of the entire students.

e.g.

dataset = {2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12}

i) what is the percentile of ranking 10?

$$\therefore \text{Percentile of rank of } x = \frac{\text{no. of value below } x}{n} \times 100$$

$$= \frac{16}{20} \times 100$$

= 80 percentile

ii) What is the value that exist at 25 percentile ?

$$\text{value} = \frac{\text{percentile}}{100} \times n$$

$$= \frac{25}{100} \times 20$$

$$\text{value} = 5 \text{ (Index)}$$

## \* 5 Number Summary :-

- ① Minimum
- ② first quartile (25%)  $Q_1$
- ③ Median
- ④ Third quartile (75%)  $Q_3$
- ⑤ Maximum

It is used for removing outliers.

Q. How to remove outliers?

Let consider a example

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

for removing the outliers we have to create lower fence and higher fence.

[lower fence  $\longleftrightarrow$  higher fence]

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

$$Q_1 = 25\%$$

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

$$Q_3 = 75\%$$

IQR =  $Q_3 - Q_1$  = difference b/w two quartile.

$$Q_1 = 25\% = \frac{\text{Percentile}}{100} \times (n+1) = \frac{25}{100} (19+1) = 5^{\text{th}} \text{ index} \\ = 3$$

$$Q_3 = 75\% = \frac{75}{100} \times (19+1) = 15^{\text{th}} \text{ index} \\ = 7$$

$$\therefore IQR = Q_3 - Q_1 \\ = 7 - 3$$

$$\underline{IQR} = 4$$

$$\therefore \text{lower fence} = Q_1 - 1.5(IQR) \\ = 3 - 1.5(4) \\ = -3$$

$$\text{higher fence} = Q_3 + 1.5(IQR) \\ = 7 + 1.5(4) \\ = 13$$

$$\therefore [-3, 13]$$

$\therefore$  All the numbers are lies b/w  $[-3, 13]$

$27$  is the greater than  $13$  so we can remove it.  
 $\{1, 2, 2, 2, 3, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9\}$

$\therefore$  ① minimum = 1

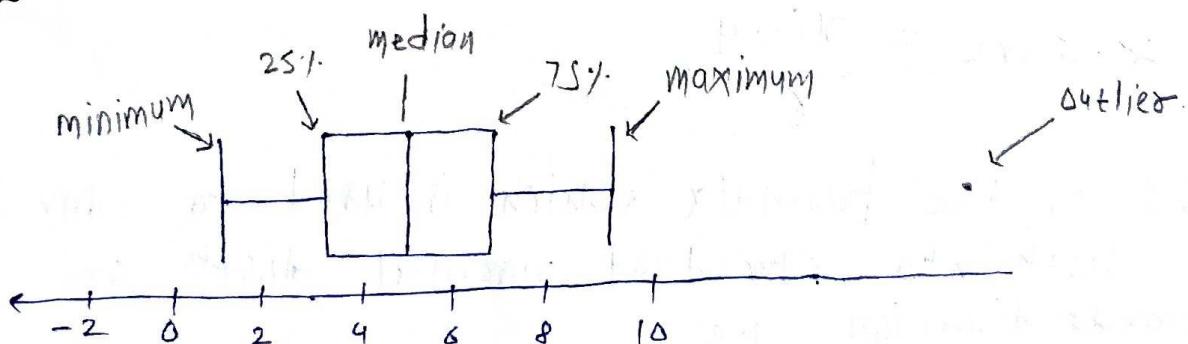
②  $Q_1 = 3$

③ median = 5

④  $Q_3 = 7$

⑤ maximum = 9

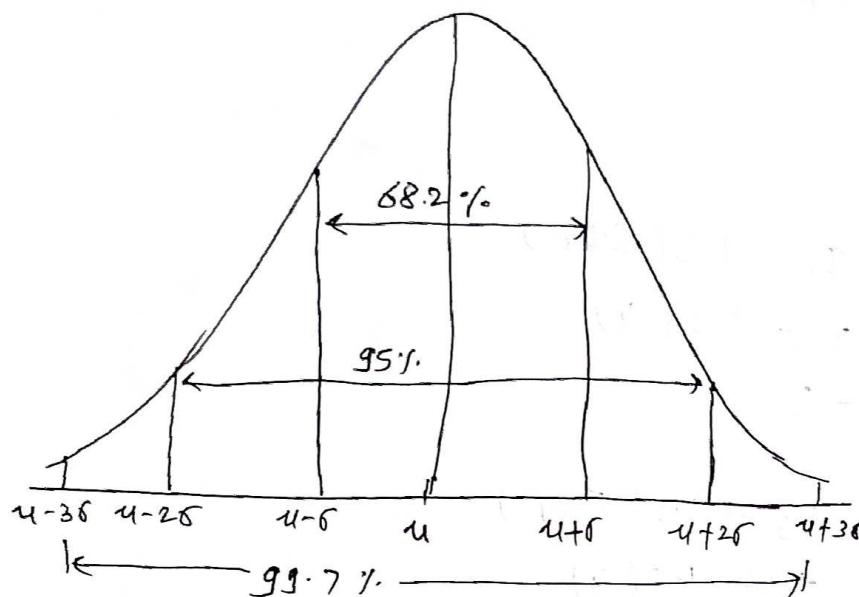
Box plot :-



## \* Normal distribution / Gaussian distn :-

It is symmetric and it has a bell shape curve.

In the normal distribution mean, median and mode are same.



e.g.

$$X = \{ \text{100 numbers in this set} \}$$

then,

68.2% of the data are present in the 1<sup>st</sup> standard deviation.

95% data are present in the 2<sup>nd</sup> S.D.

99.7% data are present in the 3<sup>rd</sup> S.D.

\* Empirical formula

$$[68-95-99.7]$$

## \* Standard normal distribution :-

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

This is the formula which is used to convert the dist into standard normal dist or standardization i.e.

$$Y \sim SND(\mu=0, \sigma^2=1)$$

## \* Why Standard Normal distn:-

Suppose in the dataset there are three diff. features and each having three different units. If we have to do some prediction then it is complicated with this data.

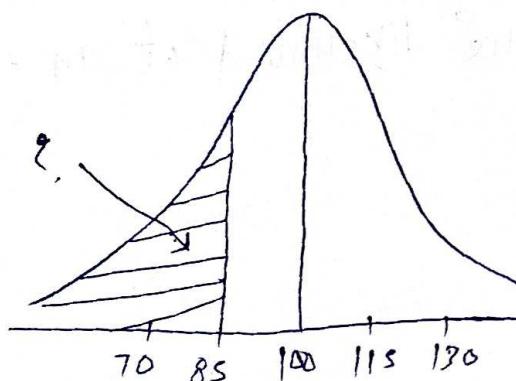
So, we can scale down of this data in some scale. The process is known as standardization.

e.g. Que.

In India the average ICP is 100 with a standard deviation of 15. What is the % of population would you expect to have an ICP lower than 85?

$$\mu = 100$$

$$\sigma = 15$$



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

Now check the value of -1 in Z-table.

The value of -1 is 0.1587.

$$\therefore \text{ICP lower than } 85 = 1 - 0.1587.$$

## \* Central limit theorem :-

If you take a sufficiently large sample size from a population with a finite level of variance  $\sigma^2$ , the mean of all samples from that population will be roughly equal to the population mean.

e.g. ~~if we take a sample of size n from a non-normal distribution, then the sample mean follows normal distribution.~~

Suppose we have a non-normal dist then we take a samples from this distribution and find the sample mean ( $\bar{x}$ ) If we plot this sample sample mean then it follows normal distribution.

## \* Probability :-

Probability is the measure of the likelihood of an event.

e.g. Tossing a coin

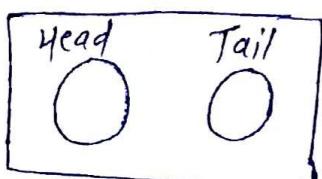
$$P(H) = 0.5$$

$$P(T) = 0.5$$

### ① Mutual exclusive events ; -

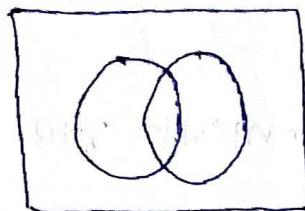
Two events are mutually exclusive events if they cannot occur at the same time.

e.g. Tossing a coin, rolling a dice.



Non mutual exclusive events :-

They occurs at the same time.



Q. What is the prob. of coin landing on heads or tails ?

$$\begin{aligned} P(A \text{ or } B) &= P(A) + P(B) \\ &= \frac{1}{2} + \frac{1}{2} \\ &= 1 \end{aligned}$$

\* Independent event :-

Two events are independent if they do not affect one another.

e.g. Tossing a coin

\* Dependent events :-

Two events are dependent if they affects from one another.

e.g. Bag of colour marbles

[ O O O O ← orange  
O O O ← yellow ]

$$\therefore P(\text{orange}) = \frac{4}{7} \text{ and } P(\text{yellow}) = \frac{3}{8}$$

## \* Permutation :-

It is the number of ways things can be ordered or arranged.

The order of the arrangement matters.

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

## \* Combination :-

order of the selection does not matter.

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

## \* Covariance :-

It is a systematic relationship between two random variables in which a change in one variable reflects a change in the other variable.

If you have a two ~~variables~~ distribution and you have to find the relation b/w them then use the covariance.

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

$$\star V(X) = \text{cov}(X, X)$$

While calculate the covariance if we get a more positive value or less negative value then we use Pearson corr coeff.

i.e. +ve  $\rightarrow \infty$       -ve  $\rightarrow \infty$ .

Pearson correlation works with linear data.

$$S_{x,y} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

Spearman rank correlation works with non-linear data, and linear.

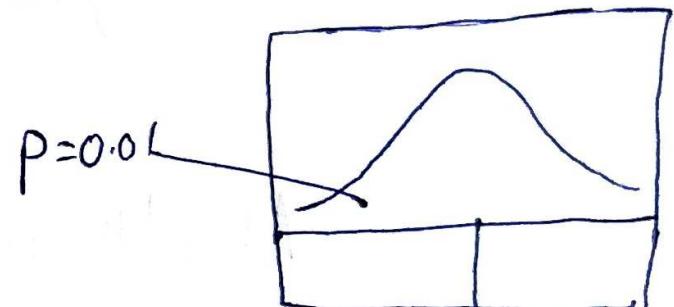
$$= \frac{\text{cov}(R(x), R(y))}{\sigma_{R(x)} \sigma_{R(y)}}$$

\* P-value :-

It is the prob. for the Null hypothesis to be true.

also called as significant value.

$P=0.01$  means that out of 100 touches 1 time you touch there.



## \* Confidence Interval :-

Confidence interval is a range of values that is used to estimate population parameters such as mean, proportion and other parameters.

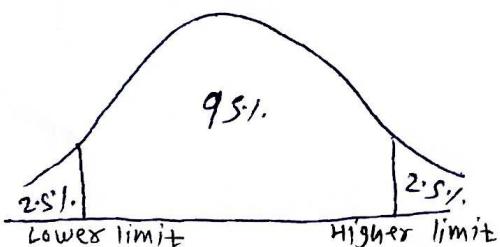
$$C.I = \text{Point estimate} \pm \text{margin error}$$

\* Point estimate :- Suppose in some experiment you are not able to find population mean in that case you have to find sample mean ( $\bar{x}$ ) and using this sample mean you are able to estimate population mean. So the sample mean  $\bar{x}$  is point estimate.

Example :- What is the average size of the shark?

$\Rightarrow$  Using the 95% confidence interval we can find it.

Let consider  $\sigma = 100$ ,  $n = 30$ ,  $\bar{x} = 500$



$$C.I = \text{Point estimate} \pm \text{margin error}$$

$$\begin{aligned} C.I &= \bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}} \\ &= 500 \pm Z_{0.05} \frac{100}{\sqrt{30}} \end{aligned}$$

$\alpha$  = Significant value

$$\alpha = 2.5 + 2.5 = 5\% = 0.05$$

$$Z_{0.025} = 1 - 0.025 = 0.976$$

$$= 1.96$$

\* Check the  $Z_{0.976}$

in the z-table

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

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

$\therefore$  At the 95% CI the average ~~size~~ of shark lies between the range 386 - 613.

## \* Hypothesis testing :-

Hypothesis testing is a statistical method used to make inferences about the population parameters based on the sample data.

This process involves formulating two hypothesis

1) Null Hypothesis ( $H_0$ ) :- This is the assumption or claim that there is no significant effect, relationship or difference.

2) Alternative Hypothesis ( $H_1$ ) :-

This is the assumption that there is a significant effect, relationship or difference.

for example,

Suppose we have to measure average height of the people in a particular city then manually it is not possible. Then with the help of sampling technique we can take a sample data and perform some experiment i.e. statistical testing. Then we come up with some conclusion about the population. This entire process is called as Hypothesis testing.

Use case :-

- (1) Hypothesis testing provides structured and objective framework for making decisions based on data.
- (2) Researchers use hypothesis testing to determine whether there is enough evidence to support or reject a specific research hypothesis.

## \* Chi-Square test :-

It is used to check the dependency of two variables.

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

i.e.  $f_o$  = Observed values

$f_e$  = Expected values

$$f_e = \frac{f_e f_r}{f_n}$$

e.g. 500 elementary school boys and girls are asked which is their favourite colour : Blue, green, or pink

Result are shown below.

	Blue	Green	Pink	
Boy	100	150	20	270
Girl	20	30	180	230
	120	180	200	

using  $\alpha = 0.05$  would you conclude that there is a relationship between gender and favourite colour.

$\Rightarrow$  Null hypothesis :-

$H_0$  : Gender and colour are related

$H_1$  : Gender and colour are not related.

$$\chi^2 = 0.05, C.I = 95\%$$

Degree of freedom :-

$$= (\text{Row} - 1)(\text{Column} - 1)$$

$$= (2 - 1)(3 - 1)$$

$$= 2$$

Check the value of d.f 2 and  $\alpha = 0.05$  in the chi-square table.

Test statistics = 5.991

chi-square test :-

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

Given table

Observed	Blue	Green	Pink	
Boy	100	150	20	270
Girl	20	30	180	230
	120	180	200	500

$$(Boys, Blue) = (120 \times 270) / 500 = 64.8$$

$$(Boys, Green) = (180 \times 270) / 500 = 97.2$$

$$(Boys, Pink) = (200 \times 270) / 500 = 108$$

$$(Girls, Blue) = (120 \times 230) / 500 = 55.2$$

$$(Girls, Green) = (180 \times 230) / 500 = 82.8$$

$$(Girls, Pink) = (200 \times 230) / 500 = 92$$

Expected	Blue	Green	Pink
Boy	64.8	97.2	108
Girl	55.2	82.8	92

$$\chi^2 = \frac{(100 - 64.8)^2}{64.8} + \frac{(150 - 97.2)^2}{97.2} + \frac{(20 - 108)^2}{108} + \\ \frac{(20 - 55.2)^2}{55.2} + \frac{(30 - 82.8)^2}{82.8} + \frac{(180 - 92)^2}{92} \\ = 19.12 + 28.68 + 71.70 + 22.44 + 33.66 + 84.17 \\ \chi^2 = 259.79$$

$$\chi^2 = 259.79 > 5.99$$

Reject the  $H_0$ .

i.e. They are not related.

## \* Statistical Testing :-

- ① Z test - (comparison of mean) - sample s.d is given (s)
- ② t test - (comparison of mean) - popn s.d is given (r)
- ③ ANOVA - (Analysis of variance)
- ④ F test - (comparison of variance)
- ⑤ Chi-square test (comparison of categorical variable)

## \* One Sample Z test :-

Z test is a statistical hypothesis test which is used to determine whether the means of two datasets are different. This test is useful when you have a large sample size and know the population standard deviation.

There are two main types of Z test

- ① One sample Z test
- ② Two sample Z test

### ① One sample Z test :-

This test is used whether the sample mean is significantly different (greater than, less than or not equal) than a population mean when the population standard deviation is known.

formula - 
$$Z = \frac{\bar{x} - u}{\sigma / \sqrt{n}}$$

e.g.

In a popl<sup>n</sup> the average IQ ( $\mu = 100$ ) and  $\sigma = 15$  then the doctor tested new medication to find out whether it increases the IQ or decreases IQ. After medication, after one month sample of 30 participants are taken ( $n = 30$ ) and this 30 had a mean ( $\bar{x} = 140$ ) so, did medication affect intelligent?  $\alpha = 0.05$

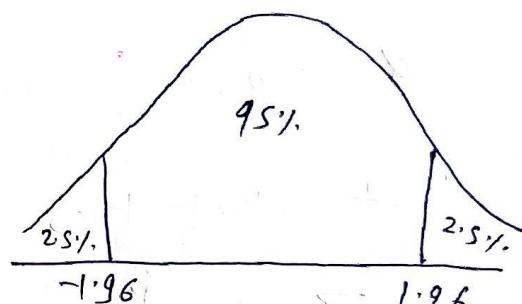
$\Rightarrow$  ① Null hypothesis  $H_0: \mu = 100$

Alternate hypothesis  $H_1: \mu \neq 100$

It is a two tail test

② State alpha  $\alpha = 0.05$

③ State decision rule



If the value lies between the  $-1.96$  to  $1.96$  then accept the null hypothesis. and

the value  $[< -1.96] \text{ or } [> 1.96]$  then reject the null hypothesis.

④ Test statistics

$$Z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{140 - 100}{15/\sqrt{30}} = 14.60$$

⑤ Result

$$14.60 > 1.96$$

$\therefore$  reject the null hypothesis and accept the alternate hypothesis.

\* One sample Z-test with proportion :-

e.g.

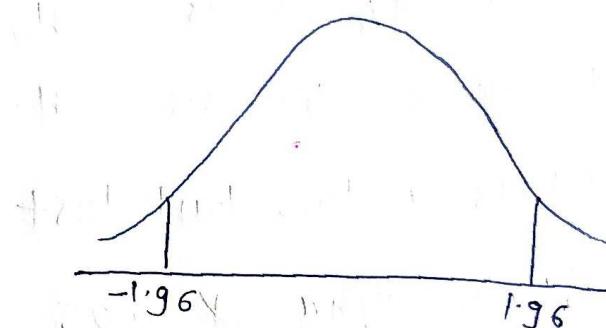
A survey claims that 9 out of 10 doctors recommend Aspirin for their patients with headache. To test this claim a random sample of 100 doctors is taken. Out of this 100 doctors, 82 indicates that they recommend Aspirin. Is this claim is accurate?  $\alpha = 0.05$

$\Rightarrow$

$$H_0: p = 0.90$$

$$H_1: p \neq 0.90$$

$$\alpha = 0.05$$



test statistics

$$Z_0 = \frac{\hat{p} - p_0}{\sqrt{\frac{p_0(1-p_0)}{n}}}$$

$$\hat{p} = 0.80$$

$$p_0 = 0.90$$

$$n = 100$$

$$= \frac{0.82 - 0.90}{\sqrt{\frac{0.90(0.10)}{100}}}$$

$$= -2.667$$

$$\therefore -2.667 \notin [-1.96, 1.96]$$

$\therefore$  Reject the null hypothesis and accept the alternate hypothesis

$\therefore$  Claim is inaccurate

## ② Two Sample Z-test :-

This test is used to compare the means of two independent samples to determine if they are different or not.

$$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

e.g.

Suppose you want to determine if there is a significant difference in the mean strength of two diff. types of materials (material A and B) that you use in your products. You have collected samples of each material and conducted strength test on them. You want to know if the mean strength of material A is significantly different from the mean strength of material B.  $\alpha = 0.05$

Material A:      Material B:

$$n_1 = 50 \quad n_2 = 60 \quad n = \text{sample size}$$

$$\bar{x}_1 = 75 \quad \bar{x}_2 = 80 \quad \bar{x} = \text{sample mean}$$

$$s_1 = 10 \quad s_2 = 12 \quad s = \text{sample standard deviation.}$$

$\Rightarrow$

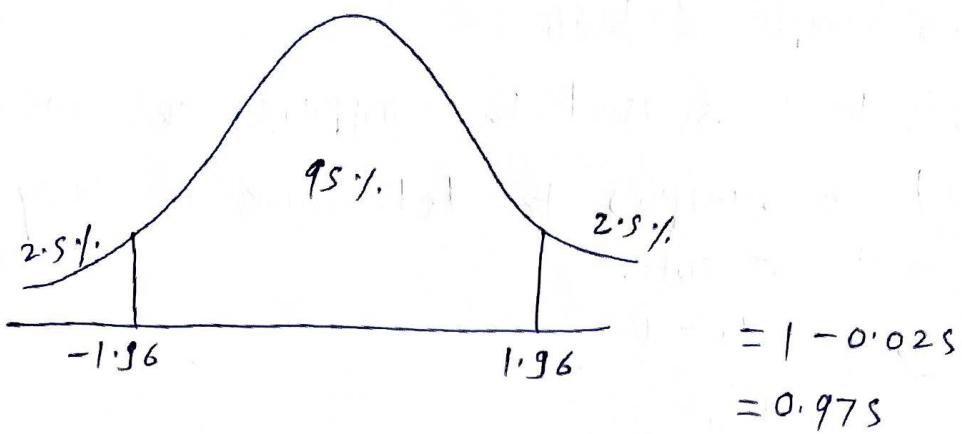
Null Hypothesis :  $H_0: \mu_1 = \mu_2$

Alternate Hypothesis :  $H_1: \mu_1 \neq \mu_2$

significant value  $\alpha = 0.05$

Test Statistics

$$Z = \frac{\bar{x}_1 - \bar{x}_2}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}} = \frac{75 - 80}{\sqrt{\frac{10^2}{50} + \frac{12^2}{60}}} = -2.369$$



Result :-

$$-2.369 \notin [-1.96, 1.96]$$

∴ we Reject the Null hypothesis.

### \* t-test :-

A t test is a statistical test that is used to compare the means of two groups. To determine if there is a significant difference between the means of two groups.

\* In this the population S.D is not given and n should be less than 30.

e.g.

In the population the average IQ is 100, a team of scientist want to test a new medication to see if it has a +ve or -ve effect on intelligence or no effect at all.

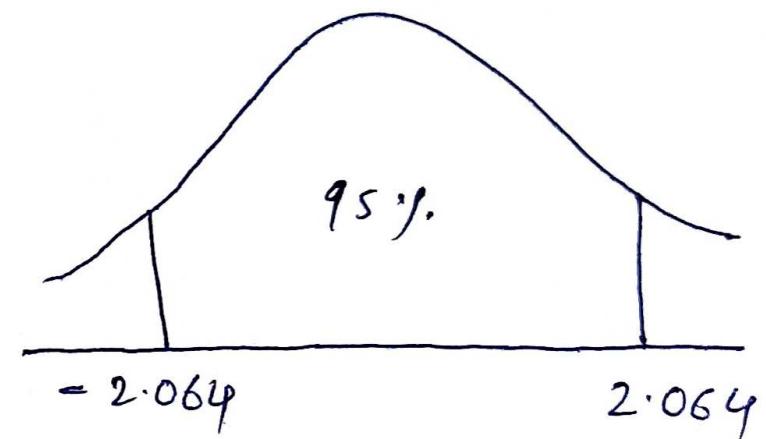
A sample of ~~25~~<sup>25</sup> participants have taken the medication and has a mean IQ of 140 with standard deviation 20. did medication affect the intelligence.  $\alpha = 0.05$

$$> ① H_0 : \mu = 100$$

$$H_1 : \mu \neq 100$$

② Degrees of freedom

$$n - 1 = 24$$



$$t \text{ test} = \frac{\bar{x} - \mu}{s/\sqrt{n}} = \frac{140 - 100}{20/\sqrt{25}} = 10$$

$$\therefore 10 > 2.064$$

$\therefore$  Reject the null hypothesis.

$\therefore$  Medication has increased the intelligence.