

Math.

Q: What is statistic?

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

Data: fact or information that can be measured.

Type of Stats:

1) DESCRIPTIVE STATS: It consists of organizing and summarizing data.

2) Inferential Stats:

where Techniques where we used the data that we have measured to form conclusion.

Population = N Sample = n

Sampling Techniques.

1) Simple Random Sampling: ~~where~~ Every member of population (N) has an equal chance of being selected to your sample (n).

2) Stratified Sampling: where the population (N) is split into non-overlapping groups (strata)

Gender -
 → Male
 → Female

Survey

3) Systematic Sampling:

From Population N we take only n^{th} individual

Eg: If I am out side of the mall
and try to Survey people

I took 8th / 1st person of the mall

4) Convenience Sampling:

Survey related to specific topic.

Eg: A survey on Govt. Job People.

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

1) Quantitative Variable: Age, weight, height,

measured Numerically (even +, -, \times , \div)

2) Quantitative

2) Qualitative / Categorical Variable:

Gender $\rightarrow M$ { Based on some characteristic
 $\rightarrow F$ we can define

Categorical Variable?

Quantitative



Discrete Variable

when number no. of
children in a family

E.g. - 1, 2, 3, ...

Continuous Variable

Height = 1.72.5, 1.62.
weight > 100kg

What kind of variable:

Gender \Rightarrow Categorical Marital status \Rightarrow Categorical

River length \Rightarrow Continuous Population of state \Rightarrow Discrete

Song length \Rightarrow "

Blood type \Rightarrow "

Variable Measurement Scales:

① Nominal: Categorical data.

② Ordinal: Order of the data matter but value doesn't.

(Student marks \Rightarrow 100 90 80 70
 \rightarrow Rank \Rightarrow 1st 2nd 3rd 4th)

③ Interval: Order matter, value matters, zero is not pos.

e.g. Age limit \Rightarrow (1 - 10) (10 - 20) (20 - 30).

④ Ratio: Naturally numerical, order data with absolute zero., zero matters

length, weight, height, Temp.

Frequency :-

Data Set : A, B, C, A, A, B, C, A, B, B, B

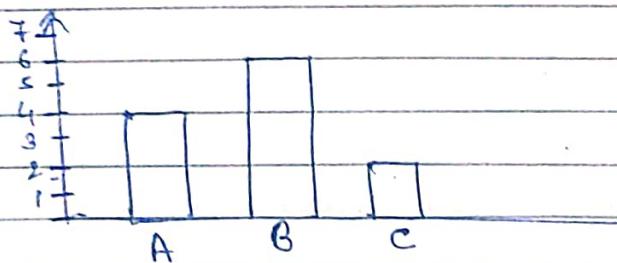
Alphab.	Frequency Distribution	Commulative Frequency
---------	------------------------	-----------------------

A	4	4
B	6	10
C	2	12

Data base -> postal
some mail items

Bin width -
some items
4,

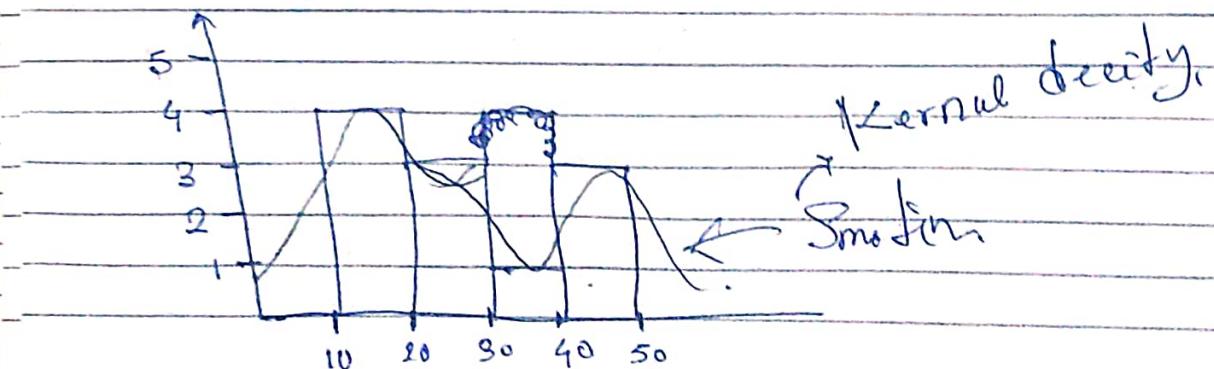
① BAR Graph



② Histogram:

$$Age = \{10, 12, 13, 20, 25, 26, 31, 40, 44, 45, 50\}$$

$$\text{Bins} = 10 \quad (\text{वर्त वर्गीयों की संख्या})$$



Aithmetic Mean for population & Sample.

MEAN: (Averages) \bar{M}

$$\bar{M} = \sum_{i=1}^N \frac{x_i}{N}$$

Population (N)

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

$$\bar{M} = \sum_{i=1}^N \frac{x_i}{N}$$

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

$$= 3.2$$

Sample (n)

$$\bar{X} = \sum_{i=1}^n \frac{x_i}{n}$$

$$= 3.2$$

► Central Tendency:

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

(i) Mean (ii) Median (iii) Mode.

$$\text{If } X = \{1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$$

$$\text{mean}(\bar{y}) = \sum_{i=1}^N \frac{x_i}{N} = \frac{1+1+2+2+3+3+4+5+5+6+100}{11} \\ \therefore \underline{\underline{132}} = 12$$

Median:

it shows the numbers = {1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100} \rightarrow 11 nm

► find out the median number: 3

but even \Rightarrow {1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100, 111)
med. $\left(\frac{3+4}{2}\right) = 3.5$

Mode : most frequent element

$$\{1, \underbrace{2, 2}_{2}, 3, 4, 5, \underbrace{6, 6, 6}_{3}, 7, 8, 100\}$$

$$\text{mode} = 6$$

* MEASURE AND DISPERSION

(i) Variance

only we use this

(ii) Standard Deviation

$$\{1, 1, 2, 3, 4\} \bar{x} = \frac{10}{5} = 2$$

Q Variance :

$$\{2, 2, 2, 2, 2\} \bar{x} = \frac{10}{5} = 2$$

both mean is same. But distribution
are diff so we use those

Q Variance :

(Population Variance (σ^2))

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

Sample variance (s^2)

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

Variance মাত্র দৃঢ়িতে

বেশি হওয়া, কম বেশি spread.

Ex-

n

μ

$n-\mu$

$(n-\mu)^2$

1

-1.88

3.34

2

-0.88

0.6889

3

2.83

+0.83

0.6889

4

0.18

0.03

5

1.17

1.37

~~total~~

5

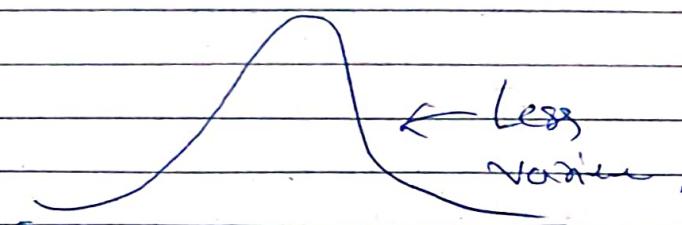
2.17

4.71

$\Sigma z^2 = 42.83$

10.89

$$\sigma^2 = \frac{10.89}{6} = 1.81$$



* Standard deviation:

$$= \sqrt{\text{variance}}$$

$$= \sqrt{\sigma^2} \rightarrow \text{with the help of formula}$$

now what is the range.

IV) Percentile & Quartiles.

With the help of those we can find the outlier.

Percentiles: If it is the value below which a certain percentage of observation lie.

Percentile Rank of $n = \frac{\# \text{ of Value below } x}{n} \times 100$

\Rightarrow number of

Q Data: 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11

What is the Percentile ranking of 10?

Percentile Ranking of $n = \frac{\# \text{ of values below } 10}{n} \times 100$

$n = 20$ ← number of ~~values~~ values

Percentile Ranking of 10 = $\frac{16}{20} \times 100$
= 80%.

Q What value will be at Percentile Ranking of 25%?

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

$$\text{Index position} = \frac{25}{100} \times 21 = 5.25$$

$$\begin{array}{l} 5^{\text{th}} \text{ index} = 5 \\ 6^{\text{th}} \text{ index} = 5 \end{array} \quad \frac{6^{\text{th}} + 5^{\text{th}}}{2} = \frac{5+5}{2} = 5$$

↑
Value

FIVE NUMBER Summary

* Removing the outliers

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

value we want to remove outlier

we need [lower fence \leftrightarrow higher fence].

$$n = 19$$

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

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

$$\text{IQR} = \text{Interquartile Range} = Q_3 - Q_1$$

$$Q_3 = 75\% \text{ (Percentile)} \quad | \quad Q_1 = 25\% \text{ (Percentile)}$$

$$= \frac{75}{100} \times (19+1)$$

$$= 15 \text{ fence}$$

$$= 7$$

$$= \frac{25}{100} \times (19+1) = 5 \text{ index}$$

$$= 3$$

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

$$\text{lower fence} = Q_1 - 1.5 \times 4 = -3$$

$$\text{Upper fence} = Q_3 + 1.5 \times 4 = 13$$

$$[-3 : \longleftrightarrow 13]$$

anything is greater than 13 is outlier
less than -3 is outlier

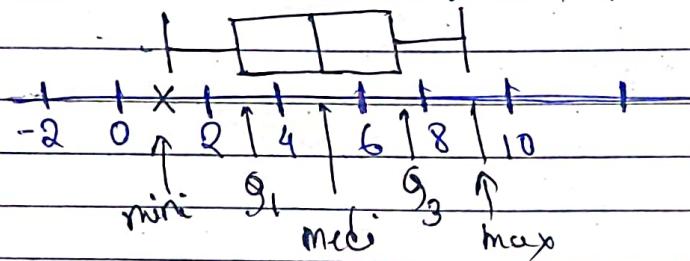
After Removal outlier

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

minimum = 1 $Q_1 = 3$

median = 5 $Q_3 = 7$ max = 9

→ Box Plot.



for the outliers

Q Cauchy Sample varij is $n-1-1$ Remove

DISTRIBUTIONS:

→ Normal Distribution

→ Standard normal Distribution

→ Z Score

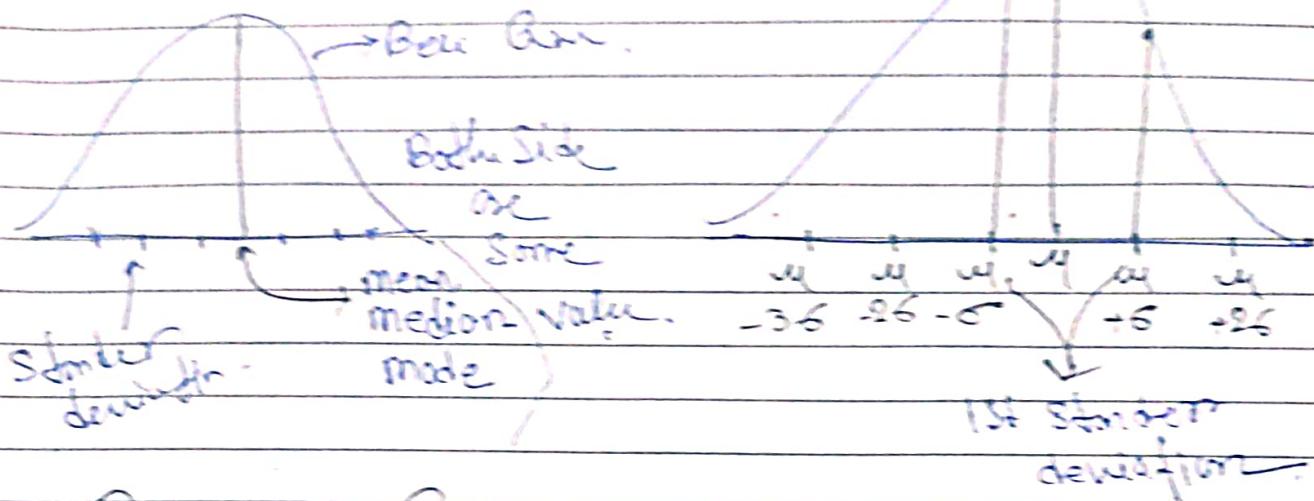
→ Log normal Distri

→ Bernoulli Dis ~

→ Binomial Dist

Ex- Ages = {24, 26, 27, 28, 30, 32}

① Gaussian / Normal Distribution



* Empirical Formula:

68-95-99.7% Rule

68% of data in the red.

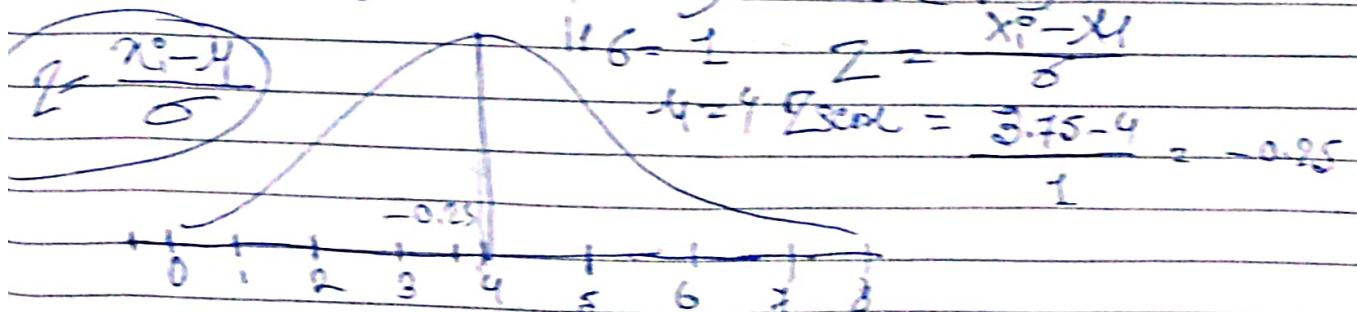
\Rightarrow 1st standard deviation = $\pm \sigma \leftrightarrow +\sigma$

95.4% distribution = 2nd standard dev. $\Rightarrow -2\sigma \leftrightarrow +2\sigma$

99.7% distribution = 3rd standard dev. $\Rightarrow -3\sigma \leftrightarrow +3\sigma$

* Z Score :

When we want to find out the corner standard deviation (like 0.25, 0.75) we use the



After we apply Z score in all value $\Rightarrow \{-3.2, -1.0, 1.3, +3.9\}$

* Standarized normal distribution ($\mu=0, \sigma^2=1$)

{1, 2, 3, 4, 5, 6, 7} \rightarrow Normal Distribution ($\mu=4, \sigma^2=1$)

↓
Apply Z score

{-3, -2, -1, 0, 1, 2, 3} \rightarrow Standard Normal distribution ($\mu=0, \sigma^2=1$)

= X -

Beneficial application

Age Salary Weight

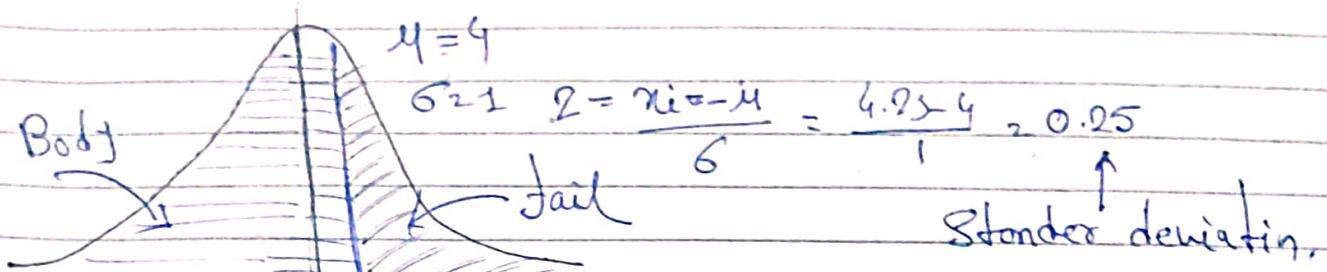
24	40k	70
25	80k	80
26	60k	55
27	70k	45

Normalization:

In machine learning we want to convert all values between 0 to 1 that is called Normalization.

* Std Dev Review Ques.

What Percentage of Score fall above 4.25?



This Curve is Symmetric
Some can thought
that whole area of this Curve
is = 1

$$\begin{aligned} \text{So Tail Area} &= 1 - \text{Left Area} \rightarrow \text{Standard deviation} \\ &= 1 - 0.5987 \rightarrow (0.25 \rightarrow 250 \text{ m}) \\ &= 0.4013 \rightarrow (250 \text{ m}, \text{ Left}) \\ &\approx 40\% \end{aligned}$$

→ we get the
Value from
Z table (graph)

Q In India the average IQ is 100, with a Standard deviation of 15. What Percentage of the population would you expect is have an IQ lower than 85?

$$\sigma = 15 \quad \mu = 100$$

$$Z = \frac{x_i - \mu}{\sigma} = \frac{85 - 100}{15} = -1$$

Probability

Probability is a measure of the likelihood of an Event.

Roll die $\{1, 2, 3, 4, 5, 6\}$ Probability of getting 6

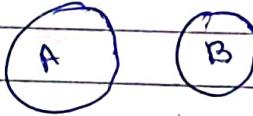
$$Pr(6) = \frac{\# \text{ of ways an event can occur}}{\# \text{ of possible outcomes}}$$

$$= \frac{1}{6}$$

* Addition Rule

► Mutual Exclusive Event : Two events are mutual exclusive if they can't occur at the same time.

Eg: Rolling die $\{1, 2, 3, 4, 5, 6\}$ it put 1 value at one time.



► Non Mutual Exclusive : multiple events can occur at the same time.



Q If I toss the a coin, what is the Probability of the coin landing on heads or tails?

$$Pr(A \text{ or } B) = P(A) + P(B) = \frac{1}{2} + \frac{1}{2} = 1$$

OR \rightarrow $\frac{6}{6} + \frac{1}{6}$

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

$$= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{3}{6} = \frac{1}{2} \approx 0.5$$

Q You are picking a Card randomly from deck
 what is the probability of choosing a Card
 that is queen or heart? (Non mutual)

$$\text{total Card} = 52$$

One type card = 13 Card

$$\text{total type} = 4$$

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

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

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

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

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

* Multiplication Rule:

① Independent Event:

② Dependent Event:

* Independent events

what is the probability of rolling a '5' and
 then a '4' in die?

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

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

$$= \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

* Dependent

which is the probability of drawing a Queen and then a King from a deck?

$$P(A \text{ and } B) = P(A) \times P\left(\frac{B}{A}\right)$$
$$= \frac{4}{52} \times \frac{3}{51}$$

* Permutation & Combination :-

↓
Arrangement
of

↓
Selection . $P_{C_p} = \frac{n!}{r!(n-r)!}$

Formation $P_{P_p} = \frac{n!}{(n-r)!}$

* P Value : is a number describe how likely it is that your data would have occurred under the null hypothesis of your statistical test.

β How many time ए पाति कराये = 25 & value to अवलम्बन
करा दिया गया

Hypothesis Testing

lets, a Guest base on Rain.

there is some test steps

1) Null Hypothesis : Coin is Fair

2) Alternative Hypothesis : Coin is unfair

3) Experiment

4) Reject or Accept the null hypothesis

How we find out the Coin is fair or unfair?

with the help of Significant Value (α)

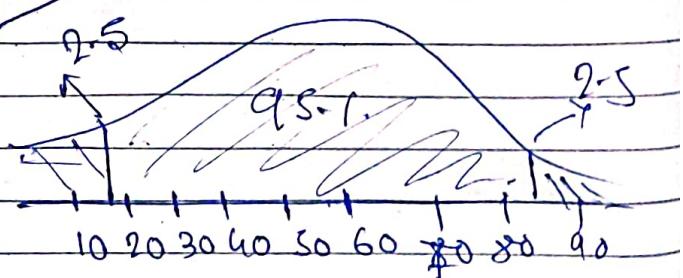
if $\alpha = 0.05 = 5\%$

$\therefore CI = \text{Confidence Interval} = 100\% - \alpha$

$$= 100\% - 5\%$$

$$= 95\%$$

If value fall in the
CI, that is
fair value.



$$2.5 + 2.5 = 5\% \\ (-) (+)$$

* ERROR.

Type I and Type II error

Null Hypothesis (H_0) = Coin is fair

Alternative Hypothesis (H_1) = Coin is unfair

Two type of error

Reality Check: Null Hypothesis is True or Null Hypothesis is False

Decision:

Null Hypothesis is True or Null Hypothesis is False.

Outcome 1:

We reject the null Hypothesis, when in reality it is false \rightarrow good decision

Reality

Outcome 2:

We reject the null Hypothesis, when in reality it is true \rightarrow good decision

\hookrightarrow Good decision (Type I error)

Outcome 3:

We accept / retain the Null Hypothesis when in reality it is false

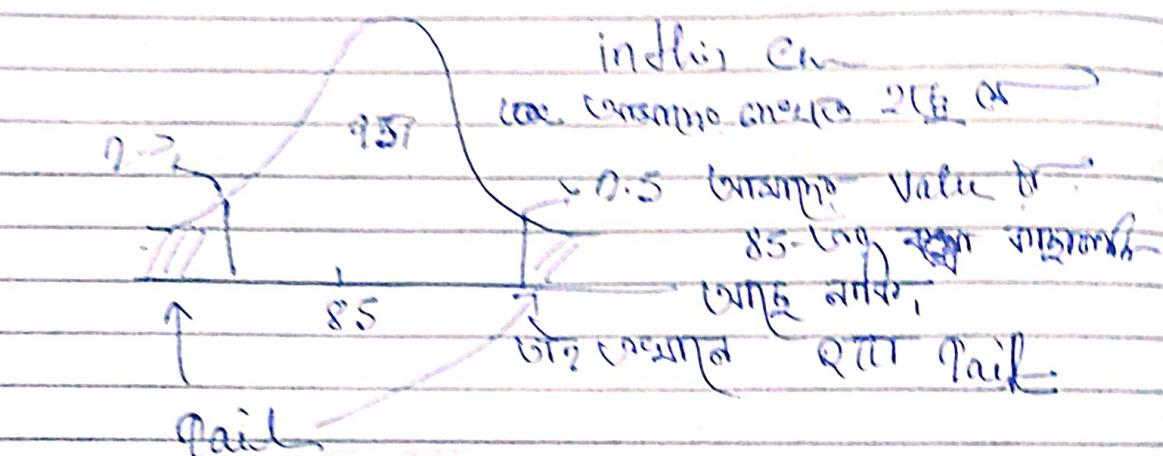
\hookrightarrow Not Good decision (Type II error)

Outcome 4:

We accept the Null Hypothesis when in reality it is true. \rightarrow Good.

X & 1 Tail and 2 Tail Test:

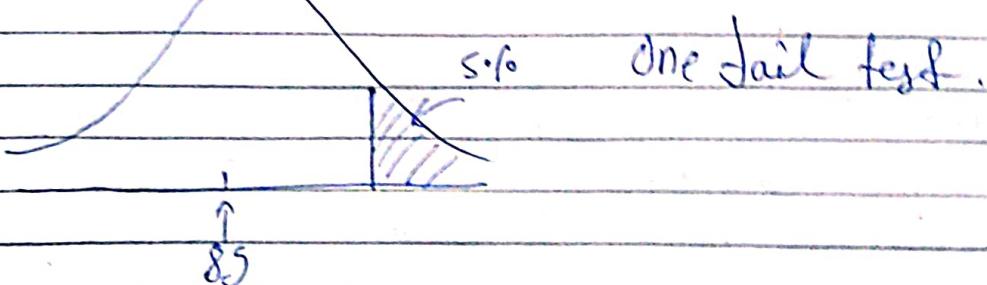
Q^Y Colleges in Karnataka have an 85% placement date. A new college was recently opened and it was found that a sample of 100 students had a placement date of 88% with a standard deviation 4%. Does this College has different placement date? Let $\alpha = 0.05$



\therefore So this is two tail test.

But If Q^Y Does this College have placement date greater than 85%?

\rightarrow See fig,



* Confidence Interval (CI)

Point Estimate : The value of any statistic that estimate the value of a parameter.

$$\bar{x} \rightarrow \mu \rightarrow \text{Estimate the mean}$$

$\bar{x} = 9.9$ 3.

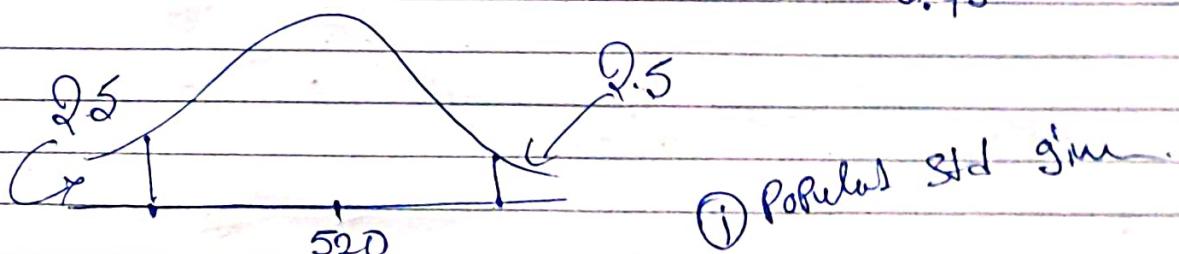
CI : Point Estimate \pm Margin of Error.

Q) On the Quant test of CAT Exam, Standard deviation is known to be 100. A Sample of 25 test takers has a mean of 520 score. Construct a 95% CI about the mean?

Ans

$$\sigma = 100, n = 25 \quad \alpha = 0.05, \bar{x} = 520$$

$$CI = \bar{x} \pm z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} = 520 \pm z_{0.05} \cdot \frac{100}{\sqrt{25}}$$



CI = Point est \pm margin

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

$$\text{Upper} = \bar{x} + z_{0.05/2} \cdot \frac{\sigma}{\sqrt{n}} = 520 + z_{0.025} \cdot \frac{100}{\sqrt{25}} = 559.2$$

$$\text{Lower} = \bar{x} - z_{0.05/2} \cdot \frac{\sigma}{\sqrt{n}} = 520 - z_{0.025} \cdot \frac{100}{\sqrt{25}} = 480.8$$

One Sample Z-Test:

- (i) Population Std is given
 - (ii) Sample size $n=30$

In the population, the average IQ is 100 with a SD of 15. Researchers want to test a new medication to see if there is a positive or negative effect on intelligence or no effect at all. A sample of 30 patients takes the medication has mean of 110. Did the medication affect IQ? $\alpha=0.05$

Am

$$CI = 95\%$$

i) Define Null Hypothesis

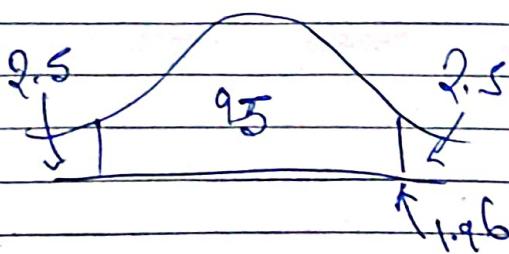
$$H_0 = \mu = 100$$

2) Arithmetic Hypothesis $H_1: \mu \neq 100$

3) Stat Alpha $\alpha = 0.05$

q) Star dicidin ful

C1295-1



3) Calculate Z statistic ~~and~~

$$2 \times \bar{n} - 4 = \bar{n} - 4 \Rightarrow 140 - 100 \Rightarrow 14.60$$

~~$\frac{6}{\sqrt{n}}$~~ ~~$\frac{6}{\sqrt{n}}$~~ ~~15~~ ~~$\sqrt{30}$~~

↓ ↓ ↓ ↓

standard error

* One Sample t-test .

use one 2-test when Population S.d.

use " t-test " unknown population

9 population the mean $\bar{x} = 100$

$$n=30, \bar{x}=140, s=20$$

$$\alpha = 0.05$$

Ans

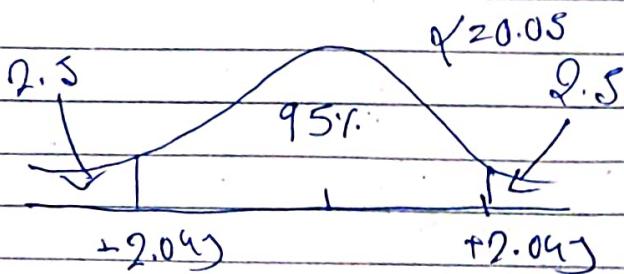
$$\textcircled{1} H_0: \mu = 100$$

$$\textcircled{2} H_1: \mu \neq 100$$

(iii) Calculate degree of freedom

$$n-1 = 30-1 = 29$$

(iv) State Decision Rule



(v) T - Test

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

* CHI SQUARE TEST (क्षेत्र विशेष)

Chi Square test claims about Population Proportions.

It is a non Parametric test that is performed on Categorical (nominal or ordinal) data.

In ~~the~~ the 2001 Indian Census, the age of Individual in a Small town were found to be following:

Less than 18	18-35	35+
20%	30%	50%

In 2010, age of $n=5000$ individual were the Sample

<18	18-35	35+
121	288	91

Using $\alpha = 0.05$, would you conclude the population distribution of ages has changed in last 10 years?

Ans:

≤ 18	$18-35$	> 35	Q_000
90%	30%	30%	← Expected

≤ 18	$18-35$	> 35	
121	288	91	$n=500$
$\frac{500 \times 90}{100}$	$500 \times \frac{30}{100}$	$500 \times \frac{30}{100}$	→ Observed
= 100	= 150	= 250	← Expected

≤ 18	$18-35$	> 35	→ Observed (f_o)
121	288	91	

100	150	250	→ Expected (f_e)
-----	-----	-----	----------------------

1) $H_0 = \text{The data meets the distribution } Q_000 \text{ test}$

$H_1 = \text{" " " does not meet, " " " , " " " .}$

2) $\alpha = 0.05$ (95% CI)

3) Decision Boundary



4)

$$\chi^2 = \text{CHI-SUM}$$

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

$$= 282.94$$

Covariance: (there is no fixed value)

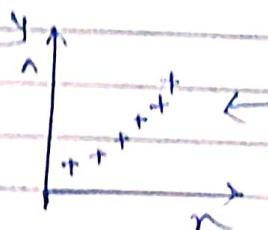
Satisfy relationship between X & Y

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

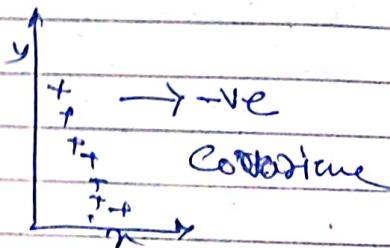
= +ve or -ve or 0

+ve \rightarrow [X ↑ Y ↑ or X ↓ Y ↓]

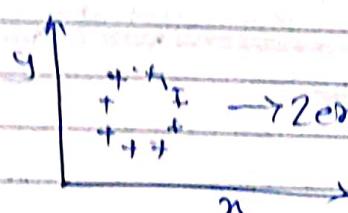
-ve \rightarrow [X ↓ Y ↑ or X ↑ Y ↓]



← +ve Correlation
Covariance



→ -ve Correlation



→ Zero Covariance

* Pearson Correlation Coefficient.

(i) It value between (-1 to +1)

The more towards +1 more positively Correlated.

The more towards -1 more negatively Correlated

$$f(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

* Spearman's Rank Correlation.

It Capture the non linear factors

$$\text{Spearman}(x, y) \rightarrow \frac{\text{Cov}(R(x), R(y))}{\text{R}_{xy} \cdot \text{R}_{yy}} = \frac{\text{Cov}(R(x), R(y))}{\text{R}_{xy} \cdot \text{R}_{xy}}$$

$R \rightarrow$ Rank of

\rightarrow It give the ranks of value

X	Y	$R(X)$	$R(Y)$
130	75	2	2
160	60	3	3
150	60	4	4
145	55	5	5
180	85	1	1

Reject Hypothesis
(P-value < Significance level)
H₀: H_0 is true
H_a: H_a is true

* P value & Significance. ($P\text{value} < \text{Significance}$)

Q The mean weight of all residents in Bangalore city is 168 pound with a standard deviation 3.9. we take a sample of 36 individuals and the mean is 169.5 pounds. C.I. = 95% ?

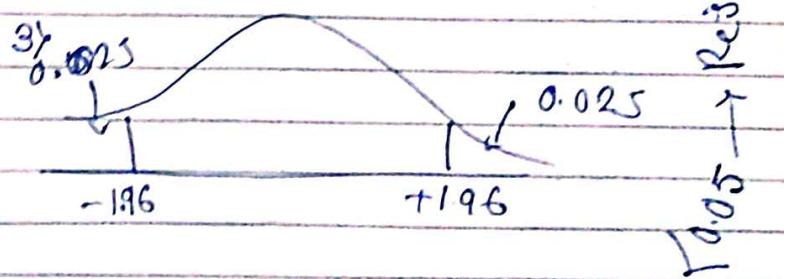
Ans

$$\mu = 168 \quad \sigma = 3.9 \quad \bar{X} = 169.5 \quad n = 36 \quad \alpha = 0.05$$

$$H_0: \mu = 168$$

$$H_1: \mu \neq 168$$

$$\text{or } \alpha = 0.05$$



if Z-test

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}} = 2.307 > 1.96 \quad \begin{matrix} \text{reject the} \\ \text{null hypothesis} \end{matrix}$$

if p-value

$$z(2.307) = 0.99111$$

$$1 - 0.99111 = 0.00889$$

