

Introduction to stats

Descriptive stats

- * Measure of Central tendency
- * Measure of Dispersion
- * Summarizing the data
Histograms, Pdf, Cdf,
Probability, Permutations
Mean, Median, Mode, Variance
and Standard deviation
- * Gaussian Distribution
- * Lognormal Distribution
- * Binomial distribution
- * Standard Normal distribution
- * Transformation & Standardization
- * Q-Q plot

What is statistics?

Statistics is the science of collecting, organizing and analyzing the data. { Better decision making }

Data :- fact or piece of information that can be measured

Eg:- The IQ of a class

{ 10, 12, 16, 19, 23, 25 }

Types of statistics:

(i) Descriptive stats :— It consists of organizing and summarizing.

Eg:- What is the avg of the student in the class.

(ii) Inferential stats :— Technique where in we used the data we have measured to form conclusion.

Eg:- Are the marks of the student of the classroom similar to the age of the matric classroom in the college.

Inferential stats

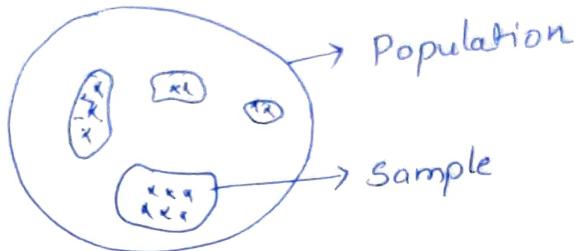
- * Z-test
- * t-test
- * ANOVA
- * CHI Square
- * Hypothesis testing
- * Confidence Intervals
- * Z-table, t-table

Population & Sample :-

Elections → Goa, UP

Population (N)

Sample (n)



- (i) Simple Random Sampling :- Every member of the population (N) has equal chance of being selected for your sample (n).
- (ii) Stratified Sampling :- Where the population (N) is split into non-overlapping groups.

Eg:- Gender [Male
 Female]

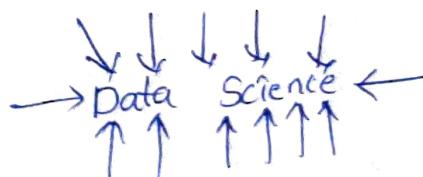
Age Groups
(10-20) (20-30) (30-40)

- (iii) Systematic Sampling :- Consider every n^{th} individual from Population.

$N \rightarrow n^{\text{th}}$ individual

Eg:- Shopping Mall → Covid Survey
 ↳ 8th person enter mall.

- (iv) Convenience Sampling :-
Sample from the only had knowledge on the domain.
(Domain expert knowledge) → only those people

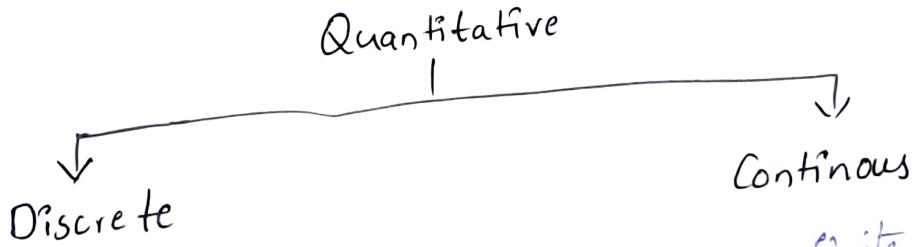


variable :- A variable is a property that can take on any value.
Eg:- $\{128, 173, 142, 162, 168\}$
Height, weight, age, marks

Two kinds of variable :

- i) Quantitative Variable :- Measured Numerically. Where we can perform arithmetic operations (+, -, *, /)
- ii) Qualitative variable / Categorical :- Based on some characteristic we can derive categorical variables.

Eg:- Gender $\begin{cases} \text{Male} \\ \text{Female.} \end{cases}$



Eg:- Whole number

- i) No. of Bank Account
- ii) No. of children in family

Eg:- non-finite numbers

Height = $\{121.6, 137.01, 198.06\}$

Weight = $\{98.6, 77.7, 88.2\}$

Rainfall = $\{0.2, 0.8, 1.1\}$

Variable Measurements :-

→ 4 types of measured variable

- ① Nominal :- Categorical data classes, color, Gender, type of flower
- ② Ordinal :- Order of the data matter, value does not matter.
- ③ Interval :- Order matters, value also matters, natural zero is present.
- ④ Ratio :-

Eg 6

<u>student (marks)</u>	<u>Rank</u>
100	1
96	2
57	3
85	4

} ordinal data

Temperature in Fahrenheit

70-80, 80-90, 90-100 →

Frequency Distribution—

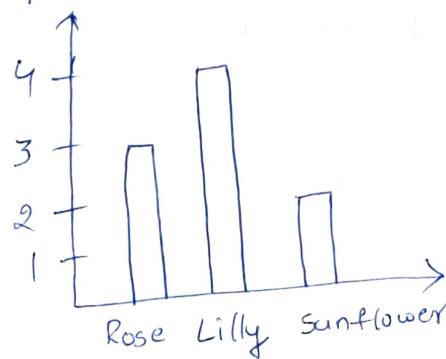
Sample Dataset :—

Rose, Lilly, sunflower, Rose, Lilly, sunflower, Rose, Lilly, Lilly

<u>Flower</u>	<u>frequency</u>	<u>cumulative frequency</u>
Rose	3	3
Lilly	4	7
Sunflower	2	9

Bar graph—

→ Bar graph is used on discrete data



Histogram :—

→ Histogram is used on continuous data

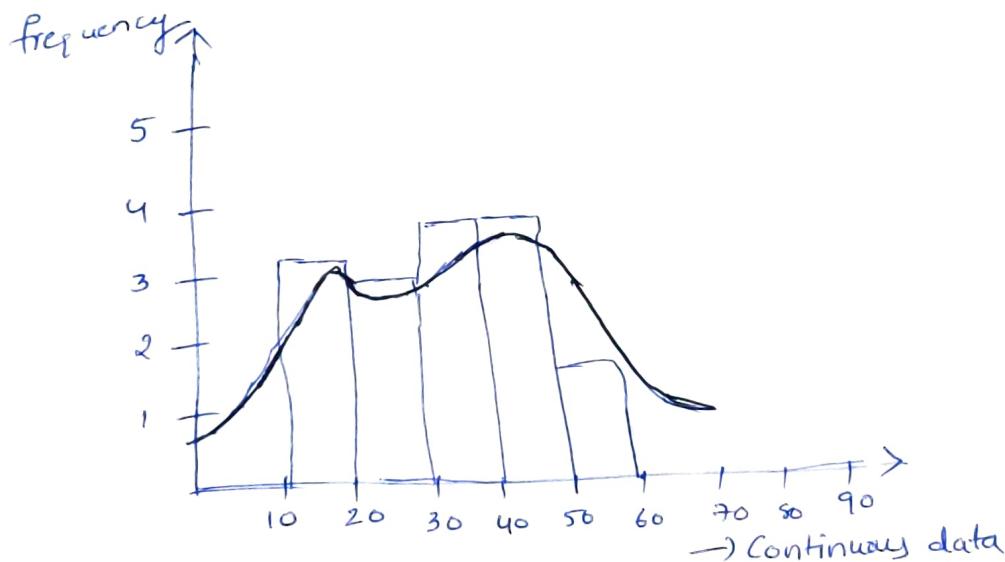
Age :- {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51}

Assume, Bins = 10

Pdf : probability distribution function

L, used to smoothing of histogram

(Kernel density ~~function~~ Estimator)



Bar vs Histogram
↓
Discrete ↓
Continuous

- ① Measure of central tendency
- ② Measure of Dispersion
- ③ Gaussian Distribution
- ④ Z-score
- ⑤ Standard normal distribution

(i) Arithmetic mean for population & Sample :-

Mean (Average)

Population (N)

Sample (n)

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

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

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

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

$$= \frac{32}{10} = 3.2$$

Central tendency :- Refers to the measure used to determine the centre of the distribution of data.

① Mean ② Median ③ Mode

(ii) Median

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

Suppose add 100 to the data

$$x = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100\}$$

Mean = $\frac{132}{11} = 12$, Mean becomes 12, it is very vary from before (3.2)

\rightarrow 100 is outlier \rightarrow In such cases use Median

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

Steps :-

(i) Sort the numbers (ASC/DSC)

(ii) Consider middle value (if n = odd)
or

take avg of two middle values { if n = even }

Median = 3

if $x = \{ 1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100, 112 \}$

$$\text{Median} = \frac{3+4}{2} = \frac{7}{2} = 3.5$$

→ Median works well with outlier

* outlier : The value completely different from distribution.

Mode :-

→ Most frequent element in data

$x = \{ 1, 2, \underbrace{2, 3, 4, 5, 6, 6, 6, 7, 8, 100, 200 } \}$

Mode = 6

Where specifically used :-

Mode works both numerical and categorical data. But it suits most to categorical.

Eg 5-

Type of flower	Petal length	Sepal length
----------------	--------------	--------------

Rose

Lily

Sunflower

—

—

Petal length

Sepal length

{ Missing value → Most frequent value (Mode) }

→ Specifically works with Categorical.

Measure of Dispersion

↳ Dispersion means spread.

(i) Variance

(ii) Standard deviation

* How 2 distributions are distributed, at that point of time, we use variance.

Population Variance
μμμ μμμ

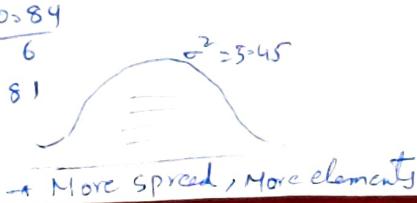
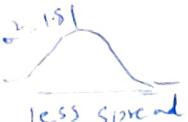
Sample Variance
n n n n n

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

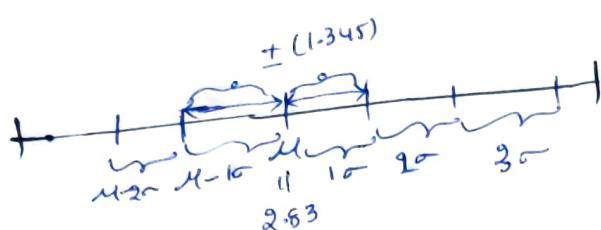
$$s^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n-1}$$

x	μ	$x-\mu$	$(x-\mu)^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
3	2.83	0.17	0.03
4	2.83	1.17	1.37
5	2.83	2.17	4.71

$$\sigma^2 = \frac{10.84}{6} = 1.81$$



$$\sigma = \sqrt{10.81} = 3.45$$



Variance = Spread

S.D = What range of value falls.

④ Percentiles & Quartiles \Rightarrow [Find outliers]

Percentage : 1, 2, 3, 4, 5

% of the numbers that are odd?

$$\% = \frac{\# \text{ of numbers that are odd}}{\text{Total Numbers}}$$

$$= 3/5 = 0.6 = 60\%$$

Percentile :-

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

Dataset : 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

(Q) what is the percentile ranking of 10?

$$\text{percentile Rank of } 10 = \frac{\# \text{ of values below } 10}{n} \times 100$$

$$= 16/20 \times 100$$

$$= 80\%$$

→ This means 80% of the data lie below 10.

(Q) percentile rank of 11 = $\frac{17}{20} \times 100$
= 85%. [85% elements are below 11]

(Q) What is value exists at percentile ranking of 25%?

$$\boxed{\text{Value} = \frac{\text{Percentile}}{100} \times (n+1)}$$
$$= \frac{25}{100} \times (21)$$
$$= 5.25 \Rightarrow [\text{Index position}]$$

Value at index 5 ip {5}

(Q) What is value at percentile of 75%?

$$\text{Value} = \frac{75}{100} \times (20+1)$$
$$= 15.75$$

Value at 16th position ip {9}.

Five number Summary :-

- (i) Minimum
- (ii) First Quartile (Q_1)
- (iii) Median
- (iv) Third Quartile (Q_3)
- (v) Maximum

Removing the outlier :-

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

less than values $\leq [Lower\ fence, Higher\ fence]$ \Rightarrow greater than values are
are outliers \Rightarrow outlier

$$Lower\ fence = Q_1 - 1.5(IQR)$$

$$Higher\ fence = Q_3 + 1.5(IQR)$$

Interquartile Range (IQR) :

$$IQR = Q_3 - Q_1$$

$$Q_3 = 75\%$$

$$Q_1 = 25\%$$

$$value\ at\ 75\%.\ percentile\ (Q_3) = \frac{75}{100} \times (20)$$

$$Q_3 = 15 \Rightarrow index = \{7\}$$

$$value\ at\ 25\%.\ percentile\ (Q_1) = \frac{25}{100} \times (20)$$

$$Q_1 = 5 \text{ index} = \{3\}$$

$$IQR = Q_3 - Q_1 = 7 - 3 = 4$$

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

$$= 3 - 1.5(4)$$

$$= -3$$

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

$$= 7 + 1.5(4)$$

$$= 7 + 6$$

$$= 13$$

outlier $< [-3, 13] >$ outlier [27 is outlier]

Remaining data = {1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9}

$$\text{Minimum} = 1$$

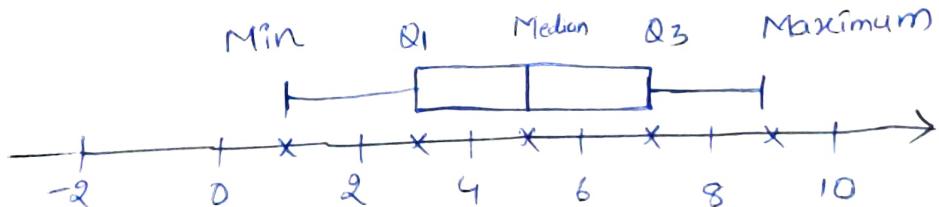
$$Q_1 = 3$$

$$\text{Median} = 5$$

$$Q_3 = 7$$

$$\text{Maximum} = 9$$

Box plot :-



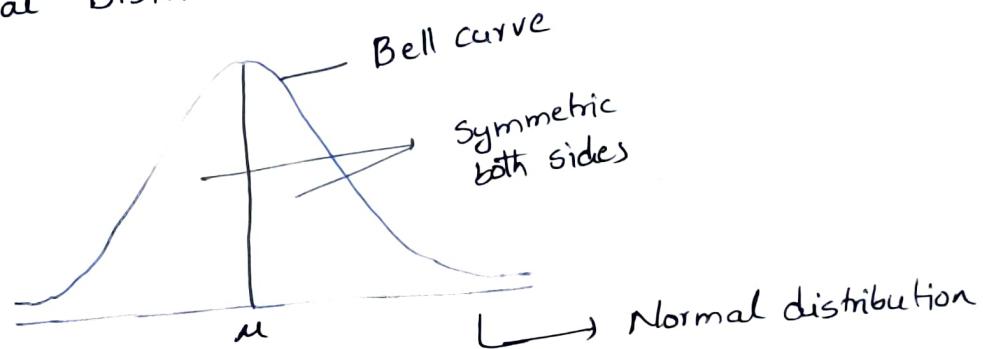
Distribution

- (i) Normal Distribution
- (ii) Standard Normal Distribution
- (iii) Z-score
- (iv) Log normal Distribution
- (v) Bernoulli's Distribution
- (vi) Binomial Distribution

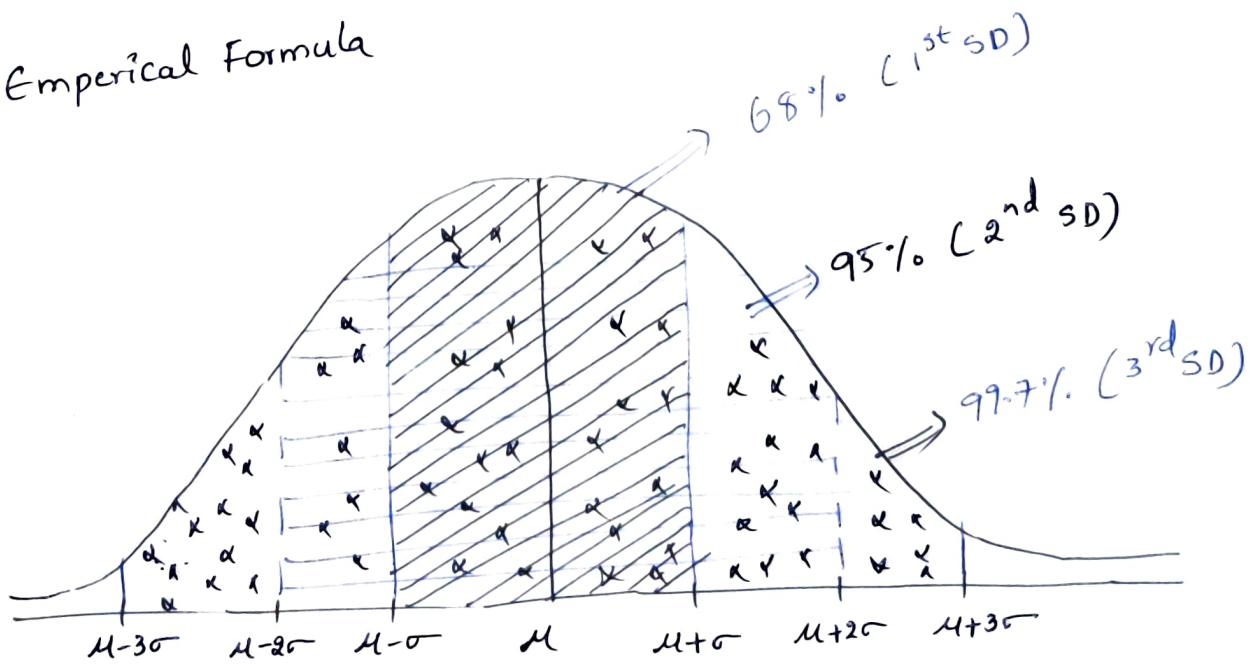
Practicals

- (i) Mean, Median, Mode
- (ii) Variance, SD
- (iii) Histogram, Pdf, Bar plot, violin plot
- (iv) IQR
- (v) Log normal Distribution

Gaussian / Normal Distribution



Empirical Formula



$$\text{formula} \Rightarrow 68 - 95 - 99.7\%$$

Eg 6

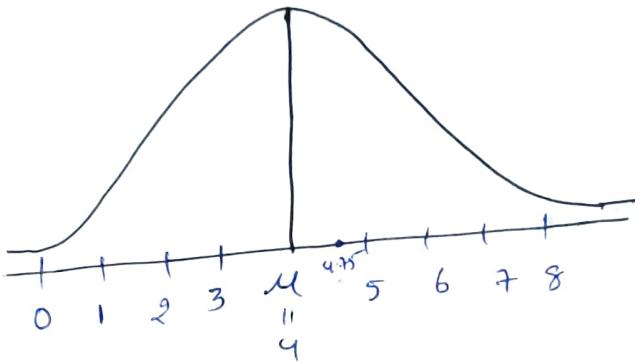
Height \rightarrow Normally distributed

Weight \rightarrow "

IRIS dataset \rightarrow Normally distributed

Example :-

$$\mu = 4$$
$$\sigma = 1$$



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

$$\text{if } x = 4.75$$

$$z = \frac{4.75 - 4}{1}$$

$z = 0.75$ sd to right

$$\text{if } x = 3.75$$

$$z_{\text{score}} = \frac{3.75 - 4}{1}$$

$= -0.25$ [SD to left]

Example :-

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



$$\mu = 4 \quad \sigma = 1$$

$$z_1 = \frac{1-4}{1} = -3$$

$$z_2 = \frac{2-4}{1} = -2$$

$$z_3 = \frac{3-4}{1} = -1$$

$$z_4 = \frac{4-4}{1} = 0$$

$$z_5 = \frac{5-4}{1} = 1$$

$$z_6 = \frac{6-4}{1} = 2$$

$$z_7 = \frac{7-4}{1} = 3$$

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

↓
Z-score
↓

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

↳ Standard Normal
distribution

which
satisfy $\Rightarrow \mu = 0, \sigma = 1$

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

↳ A Random variable y specifically satisfy $\Rightarrow \{\mu = 0, \sigma = 1\}$

Dataset :-

Age (years)	units	Salary (RS)
21		40k
26		30k
31		10k
26		55k
11		60k

units	weight (kg)
	41 kg
	25 kg
	33 kg
	56 kg
	69 kg

\Rightarrow Apply
standard deviation
 \downarrow
 $\mu = 0, \sigma = 1$
 \downarrow

This process
called
standard
deviation

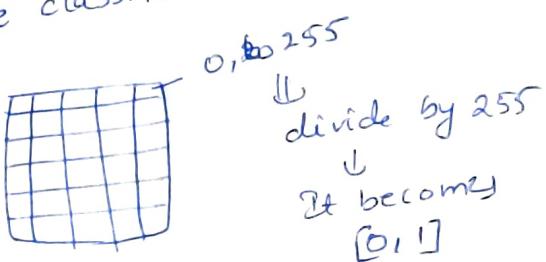
Normalization :- (Range =) $(0, 1)$

L) change all values between range $(0, 1)$
or
Interval.

MinMax scalar used to perform normalization

L) It changes all values between $(0, 1)$

CNN \rightarrow Image classification



Example 6

(Q) ODI Series

Series average in 2021 = 250

S.D of score = 10

Rishabh pant avg score = 240

Series average score 2020 = 260

S.D of score = 12

Rishabh pant avg score = 245

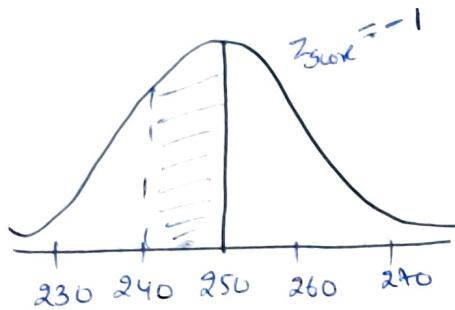
a) Compare both series in which year Rishabh pant score was better?

$$\begin{aligned} \text{2021} \\ z &= \frac{x_i - \mu}{\sigma} \\ &= \frac{240 - 250}{10} \\ &= -1 \end{aligned}$$

$$\begin{aligned} \text{2020} \\ z &= \frac{245 - 260}{12} \\ &= -1.25 \end{aligned}$$

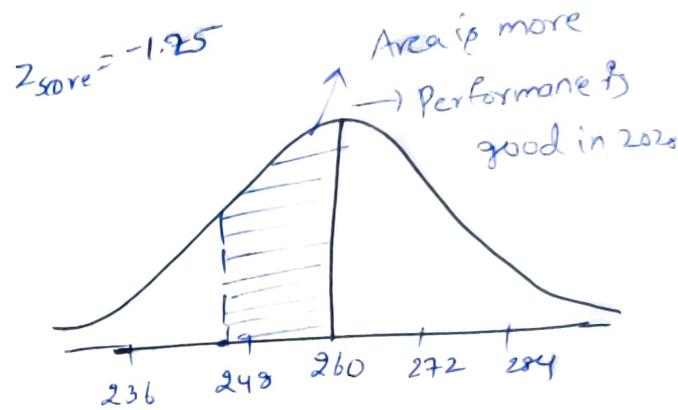
In 2021

$$\mu = 250, \bar{x}_i = 240, \sigma = 10$$



In 2020

$$\mu = 260, \bar{x}_i = 245, \sigma = 12$$

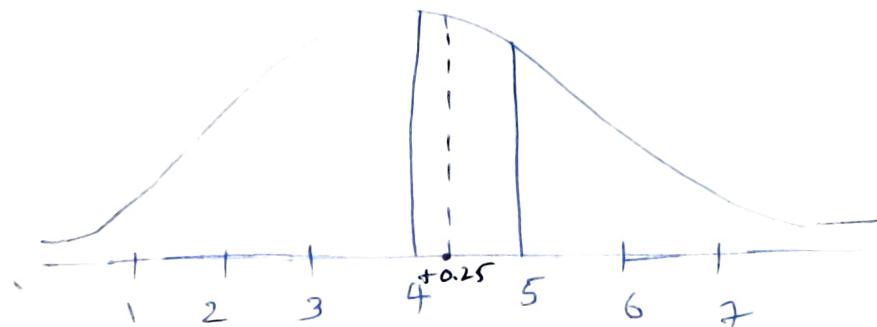


Stat Interview Questions

(Q) What is percentage of score fall above 4.25?

$$\mu = 4, \sigma = 1$$

$$Z_{\text{score}} = \frac{\bar{x}_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25 \quad [\text{s.d to left}]$$



④ Look value at 0.25 in z-table

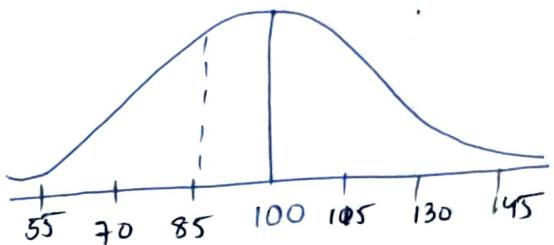
(Q) In India, the avg IQ is 100, with standard deviation of 15, what percentage of the population would you expect to have an IQ lower than 85?

Here given

$$\mu = 100$$

$$x_i = 85$$

$$\sigma = 15$$



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

$$= \frac{85 - 100}{15}$$

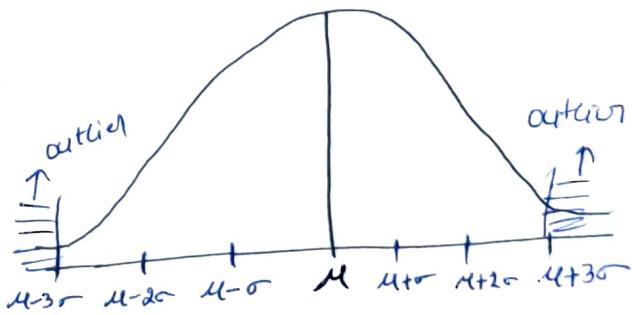
$$= -1 \Rightarrow \text{look at z-table}$$

$$\text{Value at } -1 \text{ in z-table} = 0.84$$

$$= 84\%$$

Day 4

- i) IQR python
- ii) Probability
- iii) Permutation & Combinations
- iv) Confidence Interval
- v) P-value
- vi) Hypothesis testing



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

- * using z-score and Empirical formula, Any value outside 3rd standard deviations are outliers.

Probability :-

Probability is the measure of the likelihood of an event.

Roll a dice = {1, 2, 3, 4, 5, 6}

$$pr(x) = \frac{\# \text{ of way an event can occur}}{\text{Total no. of outcomes}}$$

$$P(6) = \frac{1}{6}$$

Toss a coin = {H, T}

$$P(H) = 1/2$$

Addition Rule:- (or)

Mutual Exclusive Events : Two events are mutual exclusive if they cannot occur at the same time

Eg:- Roll a dice
~~Toss a coin~~

Non-Mutual Exclusive Events :- Multiple events can occur at the same time.

Eg:- Deck of cards

(Q) If I toss a coin, what is the probability of the coin landing on head or tail?

* They are mutual exclusive

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

Roll a dice = {1, 2, 3, 4, 5, 6}

$$\begin{aligned} P(1 \text{ or } 3 \text{ or } 6) &= P(1) + P(3) + P(6) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \\ &= \frac{3}{6} \\ &= \frac{1}{2} \end{aligned}$$

(Q) you are picking a card randomly from a deck. what is the probability of choosing a card that is queen or heart?

* This non-mutually exclusive events

$$\begin{aligned} P(A \text{ or } B) &= P(A) + P(B) - P(A \cap B) \\ P(Q \text{ or } \heartsuit) &= P(Q) + P(\heartsuit) - P(Q \cap \heartsuit) \\ &= \frac{4}{52} + \frac{13}{52} - \frac{1}{52} \\ &= \frac{4+13-1}{52} = \frac{16}{52} \\ &= \frac{4}{13} \end{aligned}$$

Multiplication Rule :-

Independent Events :-

Each event is independent

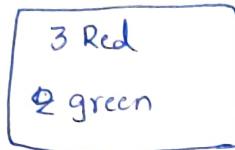
Roll a dice = {1, 2, 3, 4, 5, 6}

↳ each is independent

Dependent Events :-
 Events are dependent on next occurring events.

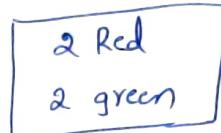
\Rightarrow After drawing red ball

Box



$$\begin{aligned} \Rightarrow P(R) &= \frac{3}{5} \\ \text{Total} &= 3+2 \\ &= 5 \end{aligned}$$

Remaining



$$\Rightarrow P(g) = \frac{2}{4}$$

(Q) Example for Independent Events

(Q) What is the probability of rolling '5' and then '4' in dice?

* This is Independent events

use multiplication Rule here

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

$$\begin{aligned} P(5 \text{ and } 4) &= P(5) * P(4) \\ &= \frac{1}{6} * \frac{1}{6} \end{aligned}$$

$$= \frac{1}{36}.$$

Example for Dependent Events :-

(Q) What is the probability of drawing a 'queen' and then 'ace' from a deck of cards?

\rightarrow These are dependent events \rightarrow Conditional probability

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

$$\begin{aligned} P(Q \text{ and } A) &= P(Q) * P(A/Q) \\ &= \frac{4}{52} * \frac{4}{51} \end{aligned}$$

Permutations & Combinations

Permutations

chocolate factory = { Dairy, 5 star, Milky bar, Eclairs, Silk }
 Total = 5

student need fill three places saw

— — —
 1st place 4 chance 3 chance
 5 chance

$$= 5 \times 4 \times 3$$

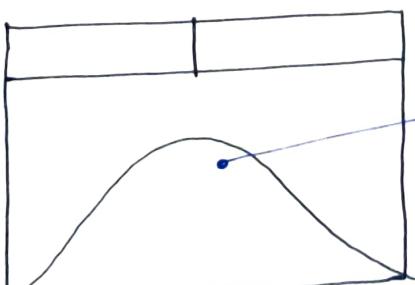
$$\text{Permutations} = \boxed{n_{Pr} = \frac{n!}{(n-r)!}} = \frac{5!}{2!} = \frac{5 \times 4 \times 3 \times 2 \times 1}{2 \times 1} = 60$$

Combinations

$$\boxed{n_{Cr} = \frac{n!}{(n-r)! r!}}$$

$$\text{Eg 6} \quad 5_{C_3} = \frac{5!}{(2!)(3!)} = \frac{5 \times 4 \times 3 \times 2 \times 1}{2 \times 1 \quad 2 \times 1} = 10$$

P-value :-



Every 100 times we touch 80% here.

Mouse pad

i.e. $p = 80\%$.

Hypothesis testing, Confidence Interval And Significance value ?

(Q) coin → Test whether this coin is fair or not by performing 100 tosses?

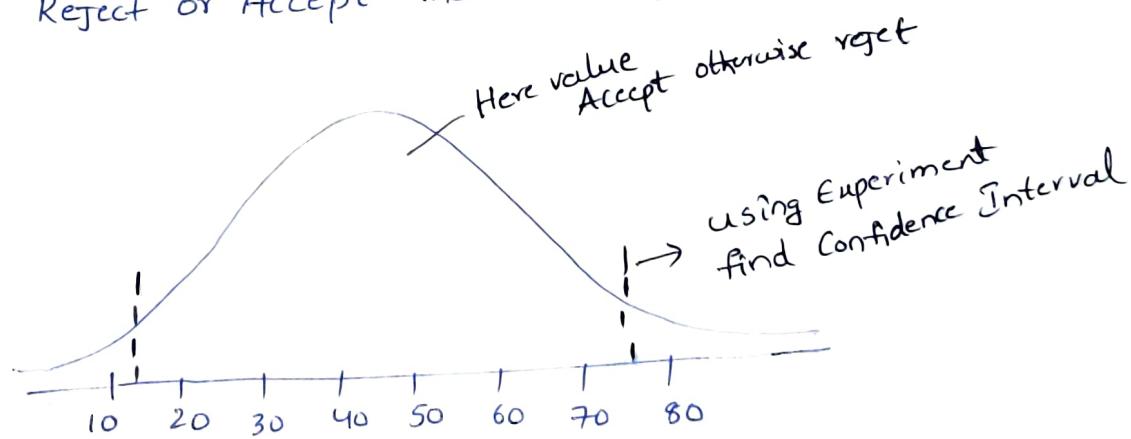
Hypothesis testing :-

i) Null Hypothesis : Coin is fair

ii) Alternative Hypothesis : Coin is unfair

iii) Experiment [z-test / t-test / f-test ---]

iv) Reject or Accept the Null hypothesis



Significance value = 0.05

↳ decided by Domain experts

Day - 5^o

- ① Type I or Type II Error
- ② One tail and 2-tail test
- ③ Confidence Interval
- ④ z-test, t-test, chi-square test

(i) Type I or Type II Errors :-

Null hypothesis (H_0) = Coin is fair

Alternative hypothesis (H_1) = Coin is unfair

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~~ reject the null hypothesis, when in reality it is false. \hookrightarrow good.

outcome 2 :-

* We reject the null hypothesis, when in reality it is True. \Downarrow

[Type I Error]

outcome 3 :-

* we retain or accept Null hypothesis, when in reality it is false. \Downarrow

[Type II Error]

outcome 4 :-

* we accept Null hypothesis, when in reality it is True. \hookrightarrow good.

Confusion Matrix :-

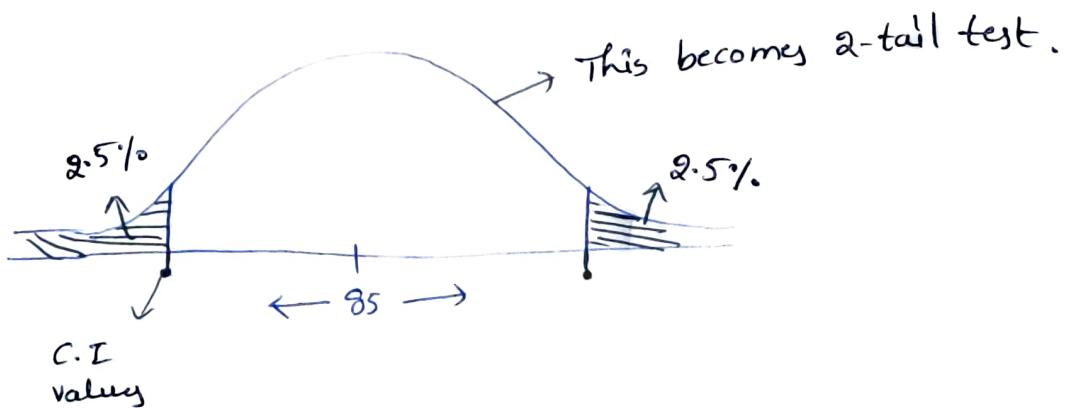
	P	N
T	TP	TN
F	FP	FN
	Type I Error	Type II Error

1-tail and 2-tail test :-

Eg 6 Colleges in Karnataka have an 85% placements rate. A new college was recently opened and it was found that a sample of 150 students had a placements rate of 88%. with a S.D 4%. Does this college has a different placement rate?

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

$$\begin{aligned}\text{Confidence Interval} &= 1 - \alpha \\ (\text{CI}) &= 1 - 0.05 \\ &= 0.95 \\ &= 95\%.\end{aligned}$$

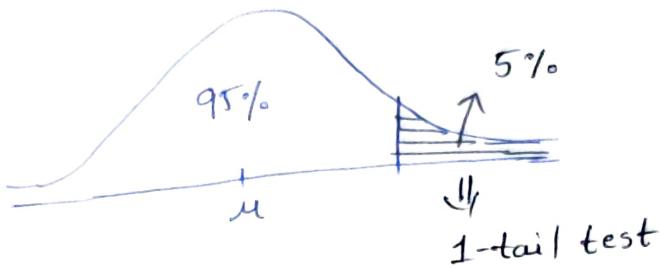


Let say

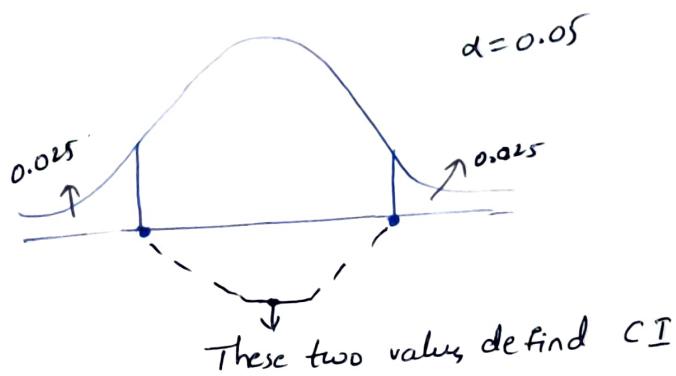
Does this college have a placements greater than 85%?

$$\alpha = 0.05$$

$$CI = 95\%$$

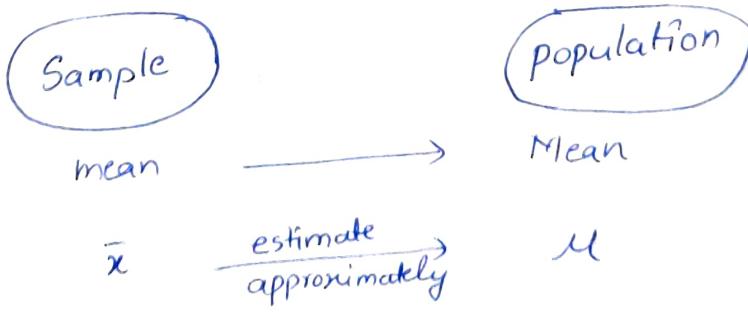


Confidence Interval for



Point Estimate \hat{x} — The value of any statistics that estimate the value of a parameter.

Inferential stats for



Eg :-
 $\bar{x} = 2.9$

$$\Rightarrow \mu = 3$$

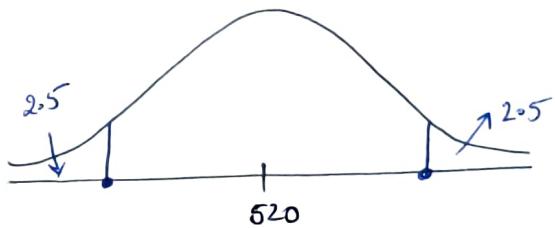
Confidence Interval = point estimate \pm Margin of error

(Q) On the quant test of CAT exam, the population standard deviation is known to be 100. A sample of 25 test takers has a mean of 520 score. Construct 95% CI about the mean?

Sol:

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

$$CI = 95\% \Rightarrow \alpha = 0.05$$



- (i) population SD is given
 - (ii) $n \geq 30$
- } [Z-test]

Point Estimate \pm Margin of Error

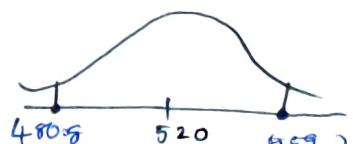
$$\boxed{\bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}}$$

$\therefore \left[\frac{\sigma}{\sqrt{n}} \Rightarrow \text{standard error} \right]$

$$\begin{aligned} \text{Upper Bound} &= \bar{x} + Z_{0.05} \frac{\sigma}{\sqrt{n}} \\ &= 520 + (1.96) \left(\frac{100}{\sqrt{25}} \right) \\ &= 559.2 \end{aligned}$$

$$\begin{aligned} \text{Lower Bound} &= \bar{x} - Z_{0.05} \frac{\sigma}{\sqrt{n}} \\ &= 520 - (1.96) \left(\frac{100}{\sqrt{25}} \right) \\ &= 480.8 \end{aligned}$$

$$\text{Range} = [480.8, 559.2]$$



(Q) find the avg size of the shark throughout the world?

Assume

$$\alpha = 0.05$$

$$\bar{x} = 1050$$

$$n = 50$$

$$\begin{aligned}\bar{x} \pm z_{\alpha/2} \frac{\sigma}{\sqrt{n}} &= 1050 + 2.005 \cdot \frac{1050}{\sqrt{50}} \\ &= 1050 + 291.045 \\ &= 1341.045\end{aligned}$$

if population SD is not given then use t-test.

(Q) On the quant test of CAT exam, a sample of 25 test takers has a mean of 520, with S.D of 80. Construct 95% CI about the mean?

Sol 6

Here population SD is not given. This situation use t-test.

$$n = 25 \quad \bar{x} = 520 \quad s = 80$$

$$\alpha = 0.05$$

formula 6

Point Estimate \pm Margin of Error

$$\bar{x} \pm t_{\alpha/2} \frac{s}{\sqrt{n}}$$

$$\text{Upper Bound} = \bar{x} + t_{\alpha/2} \frac{s}{\sqrt{n}}$$

$$\text{Lower Bound} = \bar{x} - t_{\alpha/2} \frac{s}{\sqrt{n}}$$

$$\begin{aligned}\text{Degree of freedom} &= n-1 \\ &= 25-1 \\ &= 24\end{aligned}$$

$$\frac{t_{0.05}}{2} = 2.064$$

$$\begin{aligned}UB &= 520 + 2.064 \left(\frac{80}{5} \right) \\ &= 553.024\end{aligned}$$

$$\begin{aligned}LB &= 520 - 2.064 \left(\frac{80}{5} \right) \\ &= 486.97\end{aligned}$$

$$\Rightarrow [486.97, 553.024]$$

one sample z-test:-

(Q) In the population, the avg IQ is 100 with SD of 15. Researcher wants to test a new medication to see if there is positive or negative effect on intelligence, or no effect at all of sample 30 participants who have taken the medication has a mean IQ of 140. Did the medication effect the intelligence?

Sol 5

i) Define Null hypothesis

$$\mu = H_0 = 100$$

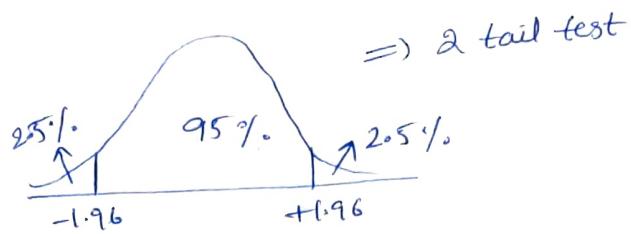
ii) Alternative Null hypothesis

$$H_1 = \mu \neq 100$$

iii) state α value

$$\alpha = 0.05$$

iv) State decision rule



\Rightarrow 2 tail test

$$Z_{\alpha/2} = \frac{Z_{0.05}}{2} = Z_{0.025} = Z_{1-0.025} = Z_{0.975} \\ = 1.96$$

v) Calculate z-statistic

$$z = \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} \quad \text{for one sample} \quad \sqrt{n}=1$$

$$= \frac{140 - 100}{\frac{15}{\sqrt{30}}} \\ = \frac{40}{15} \times \sqrt{30} \\ = 14.60$$

vi) State out decision

$$14.60 > 1.96$$

If z value is less than -1.96 or greater than 1.96 \Rightarrow reject null hypothesis.

Medication improved Intelligence? \Rightarrow yes improved.

one sample t-test

Population avg IQ = 100

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

Did the medication effect intelligence?
 $\alpha=0.05$

(i) Null hypothesis (H_0) = $\mu = 100$

(ii) Alternative hypothesis (H_1) = $\mu \neq 100$

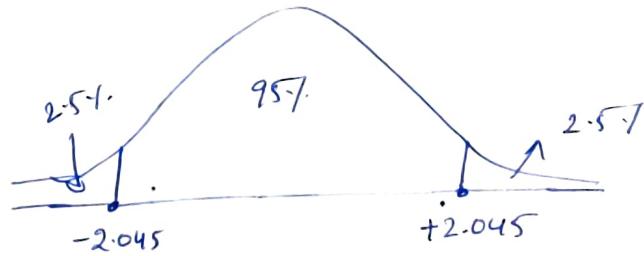
(iii) calculate degree of freedom

$$= n-1$$

$$= 30-1$$

$$= 29$$

(iv) State Decision rule



(v) Calculate t-statistics

$$\begin{aligned} t &= \frac{\bar{x}-\mu}{\frac{s}{\sqrt{n}}} = \frac{140-100}{\frac{20}{\sqrt{30}}} \\ &= 2 \times \sqrt{30} \\ &= 10.95 \end{aligned}$$

$10.95 > 2.045 \Rightarrow$ Reject null hypothesis

↳ Improved intelligence.

Day-6 B

- ① CHI square
- ② Covariance
- ③ Pearson Correlation Coefficient
- ④ Spearman Rank Correlation
- ⑤ practical Implementation
z-test, t-test, chi square test
- ⑥ f-test (ANOVA)

CHI square test :-

* Chi square test claims about population proportions

→ It is a non parametric test that is performed on categorical (nominal or ordinal) data.

(Q) In the 2000 Indian census, the ages of the individual in a small town were to be the following :

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

Sol:
In 2010, Ages of sample $n=500$ individuals were sampled. Below are the results

<15	18-35	>35
121	288	91

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

Sol 6

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

<18	$18-35$	>35	$\Rightarrow n=500 \text{ sample}$
121	288	91	$\rightarrow \text{observed}$
$500 \times 0.2 = 100$	$500 \times 0.3 = 150$	$500 \times 0.5 = 250$	$\rightarrow \text{This is expected.}$

① H_0 = The data meets the distribution 2000 census

H_1 = The data does not meet the distribution 2000 census

$$② \quad \alpha = 0.005 \Rightarrow CI = 95\%.$$

$$\begin{aligned} ③ \quad \text{Degree of freedom} &= n-1 \\ &= 3-1 \\ &= 2 \end{aligned}$$

④ Decision Boundary \rightarrow check chisquare table

$$\chi^2_{0.005} = 5.991$$

⑤ calculate test statistics

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

$$= \frac{(121-100)^2}{100} + \frac{(288-150)^2}{150} + \frac{(91-250)^2}{250}$$

$$= 232.94.$$

Covariance

\boxed{X}	\boxed{Y}
<u>In height</u>	<u>Height</u>
50	160
60	170
70	180
75	181

\Rightarrow Here relation ?
 $\boxed{X} \uparrow$ $\boxed{Y} \uparrow$

<u>No. of hour</u> <u>study</u>	<u>play</u>
2	6
3	4
4	3

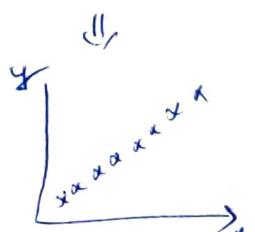
$\boxed{X} \uparrow$ $\boxed{Y} \downarrow$

$$\text{Covariance} = \boxed{\text{Cov}(x, y) = \sum_{i=1}^N \frac{(x_i - \bar{x})(y_i - \bar{y})}{N}} \Rightarrow \text{population}$$

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

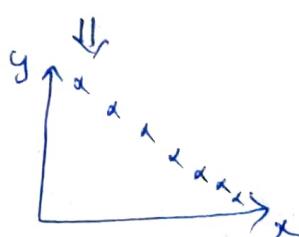
= {+ve} or {-ve} or {zero}

+ve
 $x \uparrow y \uparrow$
 $x \downarrow y \downarrow$



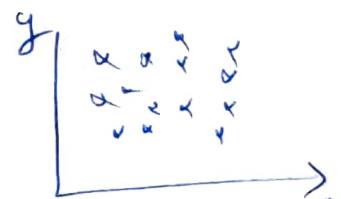
Positive
co-relation

-ve
 $x \uparrow y \downarrow$
 $x \downarrow y \uparrow$



Negative
co-relation

0
 \Downarrow No relation



Disadvantage of Co-variance :-

(1) There is no fix range of magnitude.

Pearson Coefficient :-

→ because of disadvantage of covariance, we will use pearson coefficient.

it ranges { -1 to +1 }

- * The more towards +1 is more positively correlated.
- * The more towards -1 is more negatively co-related.

$$P(x,y) = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} \quad \therefore \{-1, 1\}$$

- ④ Pearson Coefficient not works well with non-linear properties.

Spearman Co-relation Coefficient :-

- ④ Here Non-linear properties also works very well.

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

X	Y	R(x)	R(y)
170	75	2	2
160	62	3	3
150	60	4	4
145	55	5	5
180	85	1	1

$\{ p \leq \alpha \Rightarrow \text{Reject Null} \}$

Here p means probability.

Day 7 :-

(1) P-value and Significance value

(2) Distribution

(3) Central limit theorem

(4) Bernoulli's distribution

(5) Binomial distribution

(6) Poisson distribution

(7) Log normal distribution

(8) Poisson distribution

(9) Power law.

P-value & Significance value :-

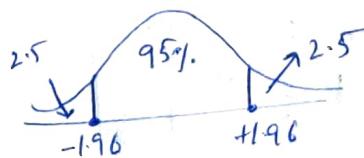
(Q) The avg weight of all residents in Bangalore city is 168 pound with standard deviation 3.9. we take sample 36 individual and mean is 169.5 pounds. CI = 95%.

$$\mu = 168, \sigma = 3.9, \bar{x} = 169.5, n = 36, \alpha = 0.05$$

$$(i) H_0 = \mu = 168$$

$$(ii) H_1 = \mu \neq 168$$

(iii) Decision boundary



(iv) Calculate Z-stat

$$Z = \frac{\bar{x} - \mu}{\sigma / \sqrt{n}} = \frac{169.5 - 168}{3.9 / \sqrt{36}} = \frac{1.5}{0.65} = 2.307$$

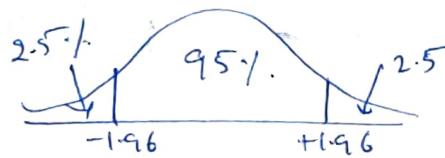
(v) Decision Rule

$Z = 2.307 > 1.96$, so we Reject Null hypothesis.

(Q) Average age of a college is 24 years with a std 1.5
 we take sample 36 students and mean is 25 and with
 $\alpha = 0.05$.

$$\mu = 24, \bar{x} = 25, \sigma = 1.5, \alpha = 0.05, n = 36$$

- (i) $H_0 = \mu = 24$
- (ii) $H_1 = \mu \neq 24$
- (iii) Decision boundary



- (iv) calculate z-statistics

$$\begin{aligned} z &= \frac{\bar{x} - \mu}{\frac{\sigma}{\sqrt{n}}} = \frac{25 - 24}{\frac{1.5}{\sqrt{36}}} \\ &= \frac{1.0}{0.25} \\ &= 4 \end{aligned}$$

- (v) Decision Rule

$|z| > 1.96$ then we Reject the null hypothesis.

Log Normal distribution



→ This is log normal distribution

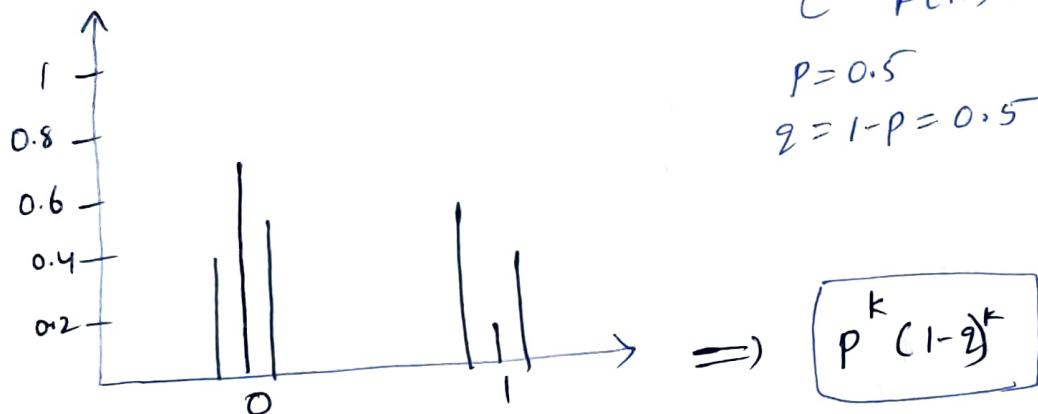
→ Apply $\log(y)$ to become Normal distribution.

Bernoulli's Distribution :-

Here 2 outcomes either {0 or 1}

Assume { Tossing a coin
 $P(H) = 0.5$

$$P=0.5 \\ q=1-P=0.5$$



→ probability mass function

↓
specifically for categorical variables.

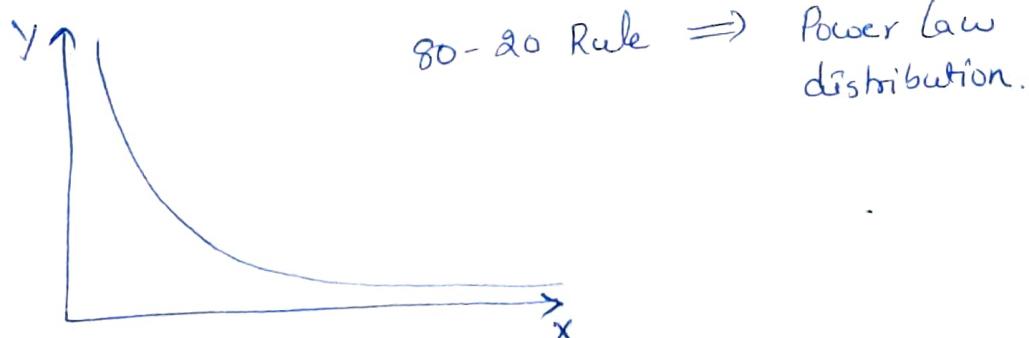
Binomial Distribution :-

With multiple trials

$$\binom{n}{k} p^k q^{n-k}$$

Pareto distribution :-

It is non-Gaussian distribution



Eg: 80% of the wealth is distributed by 20% of people

④ 80% of the company project by 20% of people in the team.

Central limit theorem 6

if sample data ≥ 30
 $n \geq 30$



* Start picking multiple samples.

$$\begin{array}{l} s_1 \rightarrow \bar{x}_1 \\ s_2 \rightarrow \bar{x}_2 \\ s_3 \rightarrow \bar{x}_3 \\ \vdots \\ s_n \rightarrow \bar{x}_n \end{array}$$

Sample mean
↓

