

Introduction to

Statistics

Date 06 / 03 / 23

Defn:- Statistics is the science of collecting, organizing and analyzing the data. (Decision making process)

Data → facts or pieces of information.

e.g. → Height of students in the class

{ 175 cm, 160 cm, 165 cm, 190 cm, 185 cm. -- }

I.Q of students { 85, 100, 95 -- }

⇒ Types of Statistics :-

① Descriptive Stats

→ It consists of organising and summarising the data.

→ Measure of central tendency

[Mean, Median, Mode]

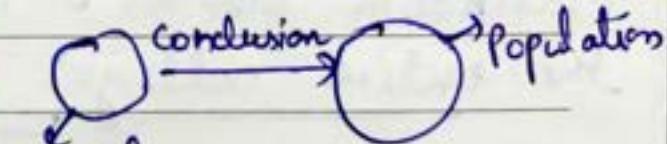
→ Measure of Dispersion

[Variance, Std]

→ Histogram, Bar chart, pie chart

② Inferential Stats

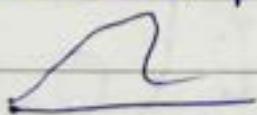
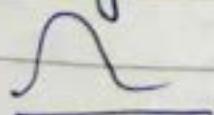
→ It consists of using data you have measured to form conclusion/output



sample

→ z-test, t-test

→ Hypothesis testing, p value



Diff. charts

  
Me Tu We Th Fr Sa Su

Population → The total collection of all the elements that we are interested in is called population.

Eg → Taking data of students in Maths class

- Total stud = 50
- Height of stud = [175cm, 180, 185, - 145,-]

Descriptive Question

Q What is the avg height of students.

Solⁿ Mean (Avg) = $\frac{\text{Summation of height}}{\text{Total students}}$

Q What is the common height of the students.

Solⁿ Let us say 140 cm is the most frequent height. = 140 cm (Mode)

Inferential Question - , (sample data)

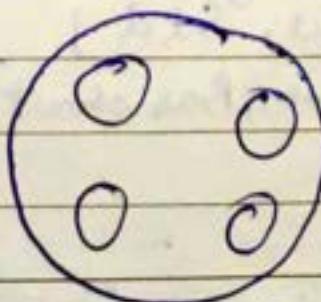
Q Are the ^{avg} height of the students in the classroom similar to what you expect in the entire college? , (population)

→ Sample Data and Population Data.

Eg → Exit poll

Party A gets 56%
Vote

Party B get 46%
Vote



Sample → A subgroup of the population that will be studied in detail is called sample.

→ In this eg. Sample data is collected from many parts of the country to know some rough data that who will win the upcoming election result.

Then conclusion will be party A will win because they got more no. of votes acc. to sample data so it is population data.

③ Types of Data :-

① Quantitative

Discrete
whole No

Continuous
Any num. value

Nominal

Ordinal

② Qualitative

eg Gender,
(M/F)

Blood Gr.,

Color.

eg Ranking
data.

- Like review
system

Good ①

Satisfactory ②

Bad ③

→ Quantitative data are numerical form data in which all the arithmetic formulae applies e.g. +, -, ×, ÷;

→ Discrete are whole nos.
e.g. No. of Bank a/c,



Mo Tu We Th Fr Sa Su

Date / /

- No. of children in a family. eg 2, 3, 4, 5
not like 1.5, 2.5, 3.5
 - In Nominal data we just categorise data.
 - In Ordinal data we assign the rank to the data acc. to preferences.
- Continuous data contain any values.
eg Weight, Height, temp., speed.

Ques What type of variable Marital Status is?

Ans Nominal Data because it is just categorical data like Married, Unmarried, Single,

Ques What kind of variable Nile River length?

Ans Continuous

Ques What type of variable is movie duration?

Ans Continuous

④ Scale of Measurements of Data :-

i) Nominal Scale Data →

- Qualitative/Categorical Variable
- e.g. Gender, Colors, Labels
- Order does not matter

Eg → Survey b/w 10 people about colour preferences -

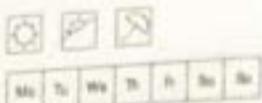
- { Red → 5 person like it → 50% opted
- Blue → 3 " " " " → 30% opted
- Yellow → 2 " " " " → 20% opted
- In this order/rating doesn't matter but only data matters in numerical form.

ii) Ordinal Scale Data →

- Ranking and order matter
- Difference can't be measured

Eg → Qualification analysis

PHD	1 st Rank	Difference can't be measured based on rank.
BE	3 rd Rank	
Masters	2 nd Rank	
B.com	5 th "	
BSC	4 th "	



Date

eg	Marks	Rank	
	100	1 st	→ It is not a
	90	2 nd	purely ordinal
	65	3 rd	data because
Dif. can be measured	70	4 th	diff. b/w the marks can
	30	5 th	be measured.

iii) Interval Scale Data

- The Rank & order matters
- Difference can be measured (excluding ratio).
- It doesn't have '0' (zero) starting value.

Eg → Temperature

If it is not a starting value

30 F	$\rightarrow 60 - 30 = 30$	$\frac{30}{60} = \frac{1}{2}$
60 F	$\rightarrow 80 - 60 = 20$	$\frac{60}{80} = \frac{3}{4}$
80 F		

Of

-10 F	90 F	$\frac{80}{90} \rightarrow \frac{8}{9}$
-------	------	---

→ Ratio can't be

-30 F relatable.

$$\frac{1}{2} \neq \frac{3}{4} \neq \frac{8}{9}$$

→ Temp. can have any negative value.



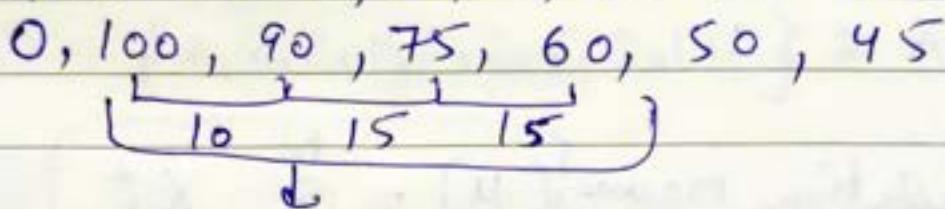
Mo Tu We Th Fr Sa Su

Date / /

i) Ratio Scale Data →

- Order and Range matter
- Diff. and ratio are measurable
- It does have a '0'(zero) starting

eg Grades → 100, 90, 60, 55, 45, 50, 75



$$\frac{100}{50} \Rightarrow 2 : 1 \rightarrow \text{So } 50 \text{ is 2 times}$$

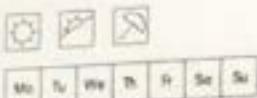
less than 100 or we can say that marks obtained by 100 grade is too greater than 50.

→ So Ratio can be relate according the mark distribution.

→ It can be start from 0 because there will be no negative grade.

Ques → length of diff. rivers in the world

Ans Ratio Scale Data



use in EDA & feature engineering
Date / /

→ Measure of Central Tendency :-

① Mean - It is the avg value of the data.

Population (N)

Sample (n)

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

$$\text{Population mean } (\mu) = \frac{\sum_{i=1}^N X_i}{N}$$

$$\text{Sample mean } (S) = \frac{\sum_{i=1}^n X_i}{n}$$

$$\mu = \frac{(1+1+2+2+3+\dots)}{10} \quad (N=10)$$

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

$$S = \frac{(1+1+2+2+\dots)}{10} \quad (n=10)$$

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

→ In some condition N, n are different and values will also be different.



Mo Tu We Th Fr Sa Su

Date / /

② Median → Middle most value of the series.

4, 5, 2, 3, 2, 1

Sort → 1, 2, 2, 3, 4, 5.

If n is even,

$$1, 2, \underline{2}, \underline{3}, 4, 5$$
$$\frac{2+3}{2} \Rightarrow \frac{5}{2} = 2.5.$$

If n is odd,

1, 2, 2, 3, 4, 5, 7

median → 3

→ Why Median is used?

{ 1, 2, 3, 4, 5 } y

$$S = \frac{1+2+3+4+5}{5} \Rightarrow \frac{15}{5} = 3.$$

Median → { 1, 2, 3, 4, 5 } y

Median = 3.

→ If we add an outlier to this distribution

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

100 is not a member of this distribution
or doesn't look like it belongs to
this distribution.

$$\text{Mean} (\bar{x}) = \frac{1+2+3+4+5+10}{6} \rightarrow \frac{115}{6} \approx 18$$

$$\text{Median} = \left\{ 1, 2, 3, \underbrace{4, 5, 10} \right\} \\ \frac{3+4}{2} \rightarrow \frac{7}{2} \rightarrow \underline{\underline{3.5}}$$

\therefore In this median is more accurate
 as there is no influence of an outlier.
 → Therefore when we have an outlier
 to our distribution it is always advisable
 to get Median it is more accurate.

③ Mode = Most frequent number

$$\{ 2, 1, 1, 4, 5, 7, 8, 9, 10 \}$$

$$\boxed{\text{Mode} = 1}$$

Type of flower

Age

Lilly

10

Rose

3

→ — (Rose)

5

Sunflower

— (Mean or
median)
in case of
outlier

Rose

8

→ Blank space will be replaced by Mode.

⇒ Measure of Dispersion :-

① Variance → Spread of the data.

Population Variance

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

x_i → Data points

μ → population mean

N → Population size

Sample Variance

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

\bar{x} → sample mean

n → population size

Eg → {1, 2, 3, 4, 5} → sample

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2, \quad \bar{x} = \frac{1+2+3+4+5}{5}$$

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

$$\bar{x} = \frac{15}{5} = 3$$

$$1 \quad 3 \quad 4$$

$$2 \quad 3 \quad 1$$

$$3 \quad 3 \quad 0$$

$$4 \quad 3 \quad 1$$

$$5 \quad 3 \quad 4 \\ \hline 10$$

$$S^2 = \frac{10}{5-1} = \frac{10}{4} = 2.5$$

$$\boxed{S^2 = 2.5}$$

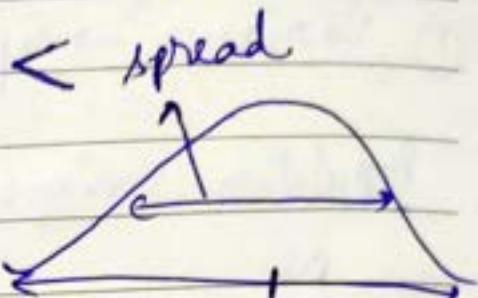
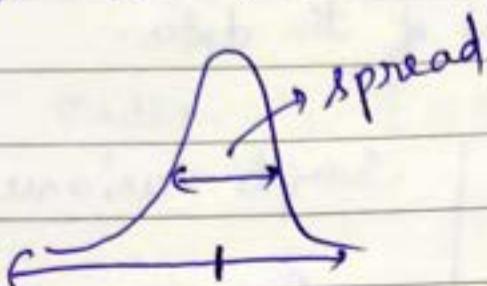


Mo Tu We Th Fr Sa Su

Date / /

$$\text{eg } s \approx 2.5$$

$$s \approx 6.5$$



② Standard Deviation :-

Population Std

$$\sigma = \sqrt{\text{Variance}}$$

Sample Std

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

$$s \approx 2.5$$

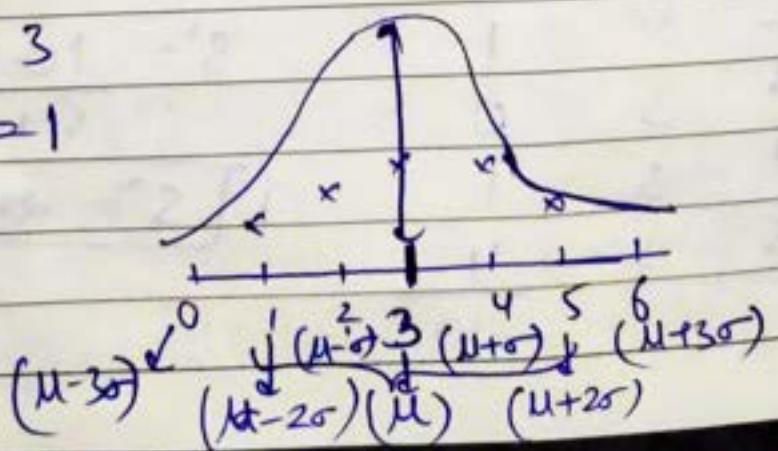
$$\sqrt{s^2} = \text{sample std.}$$

~~Q8 Consider~~

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

$$\rightarrow \mu = 3$$

$$\rightarrow \sigma = 1$$





Mo Tu We Th Fr Sa Su

Date / /

③ Random Variable :- It is a process of mapping the output of a random process or experiment to a number.

e.g. Tossing a coin

Rolling a dice

Measure the temp. for next day

$X = \begin{cases} 0, & \text{if head} \\ 1, & \text{if tail} \end{cases}$ Quantifying a Random process

$Y = \{ \text{sum of the rolling of dice 7 times} \}$

$$\{3, 4, 5, 6, 1, 2, 2\} \rightarrow Y = 23$$

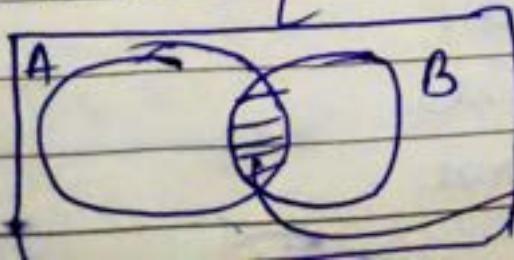
④ Sets

$$A = \{1, 2, 3, 4, 5, 6, 7, 8\}$$

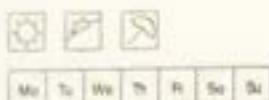
$$B = \{3, 4, 5, 6, 7\}$$

i) Intersection (\cap) → most common element

$$A \cap B = \{3, 4, 5, 6, 7\}$$



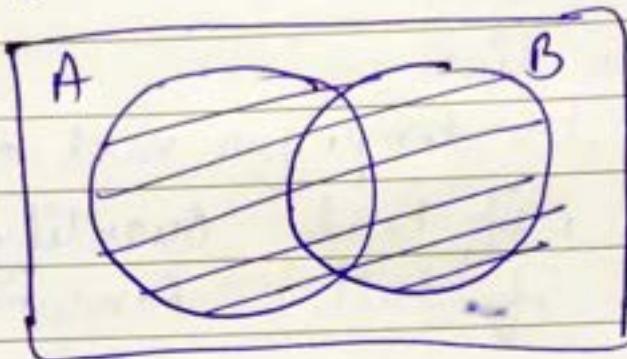
→ intersection



Date / /

ii) Union (U) → Combination of all the elements present in the both sets

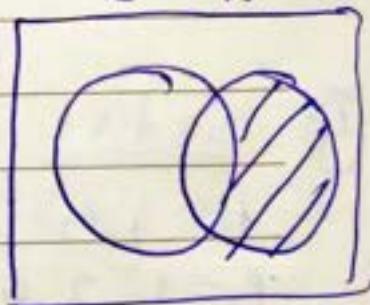
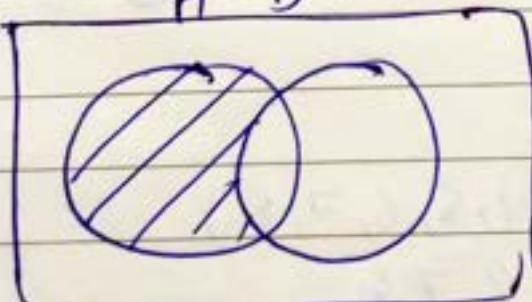
$$A \cup B = \{1, 2, 3, 4, 5, 6, 7, 8\}$$



iii) Difference (-) → Removing the element

$$A - B = \{1, 2, 8\}, B - A = \{1, 2, 0\}$$

$A - B$ $B - A$



iv) Subsets → All the elements in one set has to be present in other set to be called as subset.

$$A \rightarrow B \Rightarrow \text{false}$$

$$B \rightarrow A \Rightarrow \text{True}$$

v) Supersets → One set has all the elements which another set have.

$A \rightarrow B \Rightarrow \text{True}$

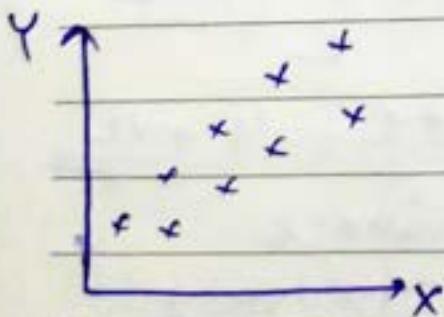
$B \rightarrow A \Rightarrow \text{False}$

⇒ Covariance and Correlation :-

	X	Y
2	3	
4	5	
6	7	
8	9	

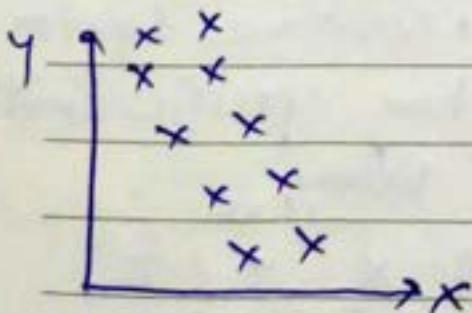
$X \uparrow Y \uparrow$
$X \downarrow Y \downarrow$
$X \uparrow Y \downarrow$
$X \downarrow Y \uparrow$

Possible relation



$X \uparrow Y \uparrow$
$X \downarrow Y \downarrow$

+ve
cov



$X \uparrow Y \downarrow$
$X \downarrow Y \uparrow$

-ve
cov



Mo Tu We Th Fr Sa Su

Date / /

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

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

spread of the data $\leq \text{cov.}(X, X) = \sum_{i=1}^n \frac{(x_i - \bar{x})(x_i - \bar{x})}{n-1}$

X	Y
2	3
4	5
6	7
$\bar{x} = \frac{12}{3}$	$\bar{y} = \frac{15}{3}$

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

$$= \frac{(2-4)(3-5) + (4-4)(5-5) + (6-4)(7-5)}{2-1}$$

$$= \frac{-2 \times -2 + 2 \times 2}{2}$$

$$= \frac{4+4}{2} = 4 \text{ +ve}$$

X & Y having a positive variance.

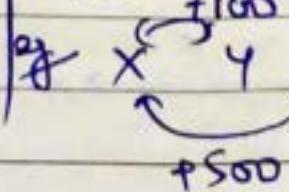
Advantage

→ Relationship b/w X & Y
+ve or -ve value

Disadvantage

→ covariance doesn't have specific limit

Value



→ Pearson Correlation Coefficient :-

Range $[-1 \text{ to } +1]$

$$f_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y} = [-1 \text{ to } +1]$$

- The more the value towards +1 the more the correlated (+ve) it is (x,y)
- The more the value towards -1 the more -ve correlated it is (x,y)

e.g. Dataset :- 1000 features (ML Models) ↑
 o/p or dependent feature

Size of house	No. of rooms	Location	No. of people staying	Price
---------------	--------------	----------	-----------------------	-------

Independent feature

- In above example there is a positive correlation b/w size of house, no. of rooms, location and price as one increases other will also be increases. but in case of no. of people staying there will be no effect on the price value. i.e. no relation. and if the house is haunted therefore it is -ve correlated with price.



Mo Tu We Th Fr Sa Su

Date / /

- feature selection \rightarrow near to 0 \rightarrow Drop no. of people staying

\rightarrow Spearman Rank Correlation \rightarrow It is better than Pearson correlation.

$$\gamma_s = \frac{\text{Cov}(R(X), R(Y))}{\sigma_{R(X)} \times \sigma_{R(Y)}}$$

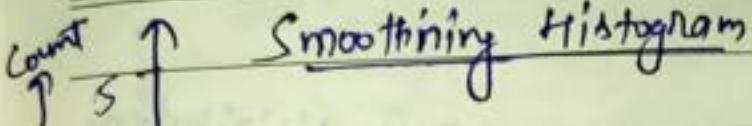
X	Y	R(X)	R(Y)
1	2	5	5
3	4	4	4
5	6	3	3
7	8	2	1
0	7	6	2
8	1	1	6

\Rightarrow Histograms [Frequency] \Rightarrow Count

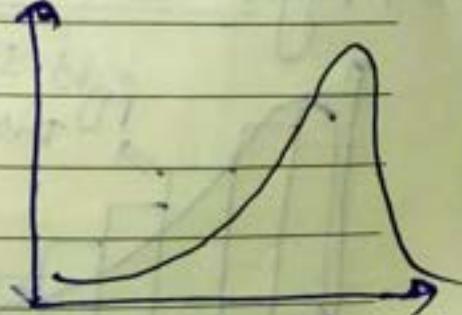
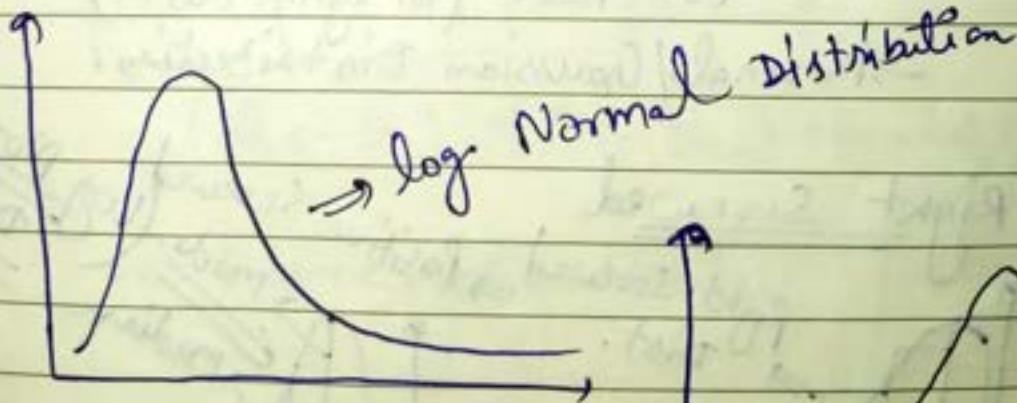
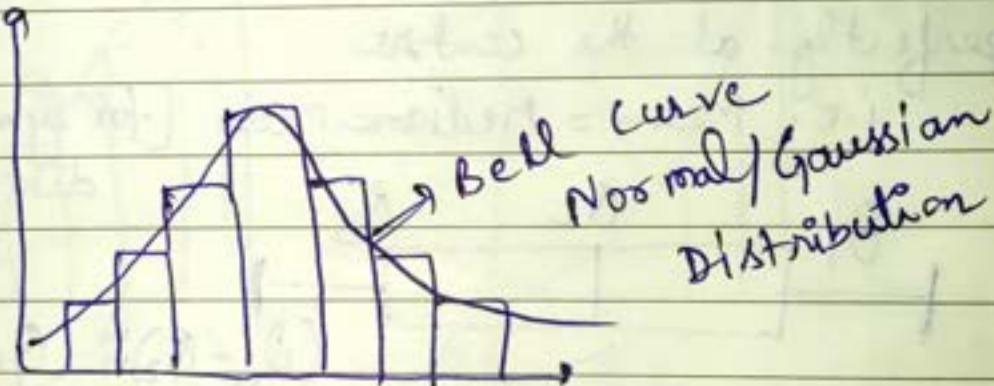
$$\text{Ages} = \{19, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 43, 50, 51\}$$

$$\text{Bins} = 10 \text{ (buckets)}$$

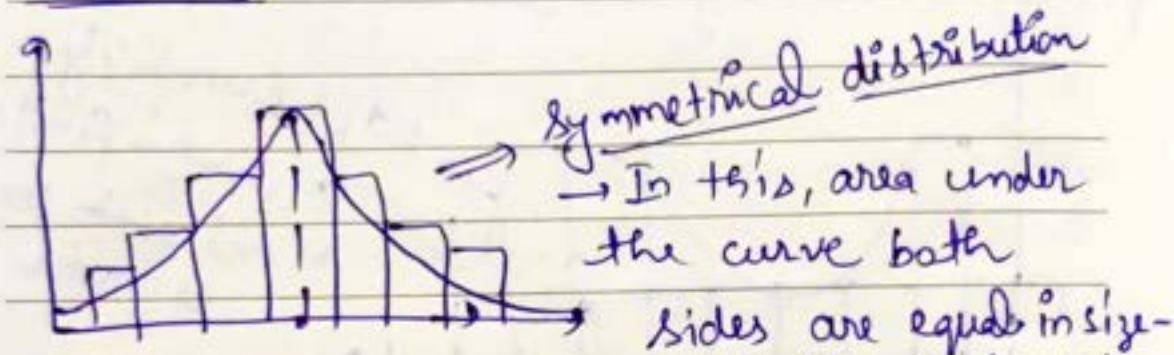
$\text{bin size} = \frac{50 - 5}{10}$



pdf [Probability Distribution Function]
 Kernel estim Density
 estimator (10 bins)

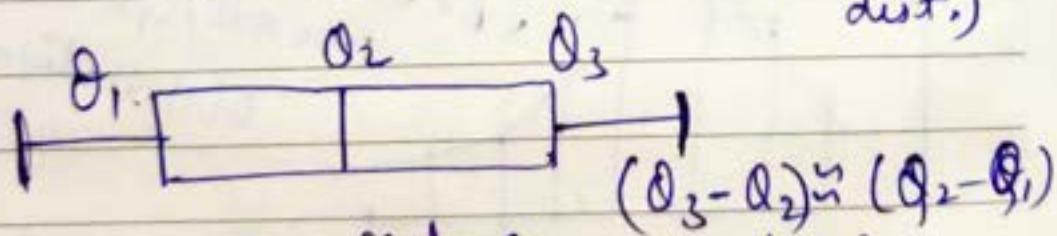


→ Skewness



- No skewness in symmetrical distribution
- The mean, median and mode all are perfectly at the centre

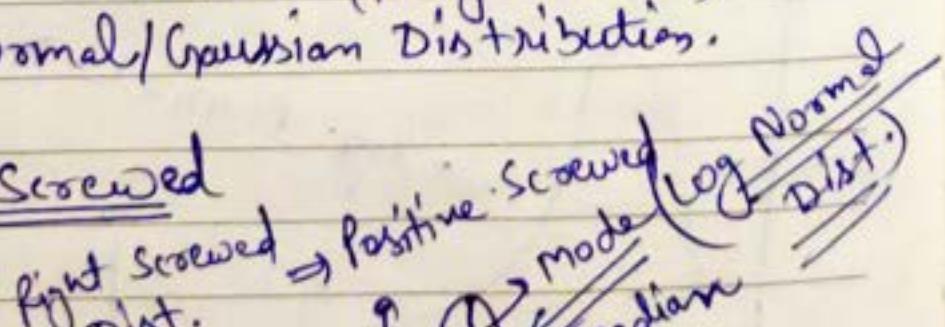
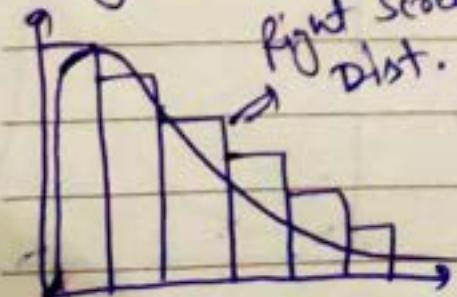
i.e. Mean = Median = Mode (for symmetrical dist.)



Box Plot (for symm. dist.)

→ Normal/Gaussian Distribution.

→ Right skewed

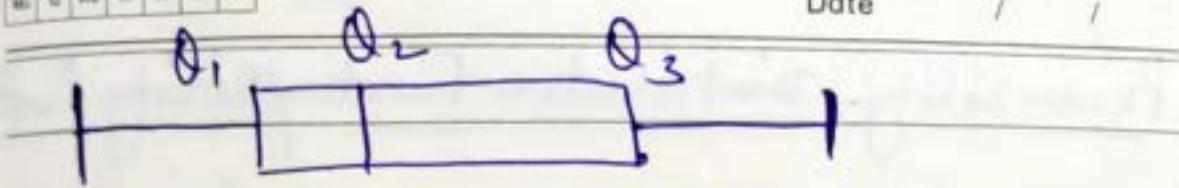


log Normal Dist.



Mo Tu We Th Fr Sa Su

Date / /

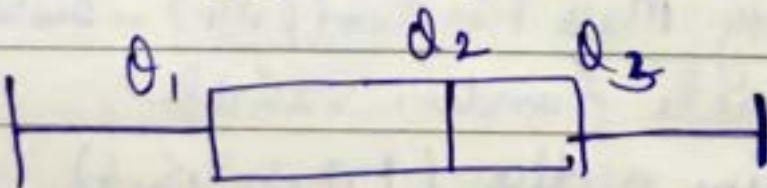
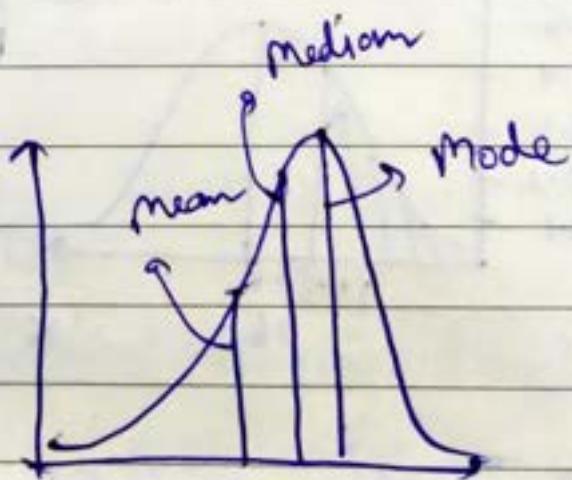
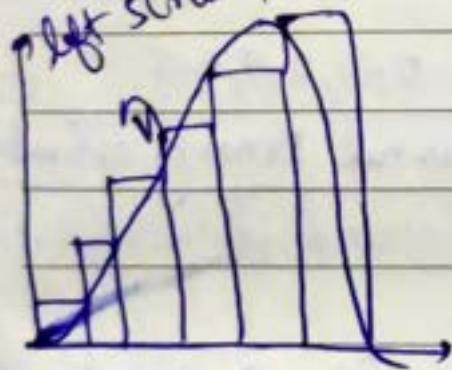


i.e. $(Q_3 - Q_2) > (Q_2 - Q_1)$ \rightarrow for right skewed

i.e. Mean > Median > Mode.

\rightarrow left skewed

left skewed dist.

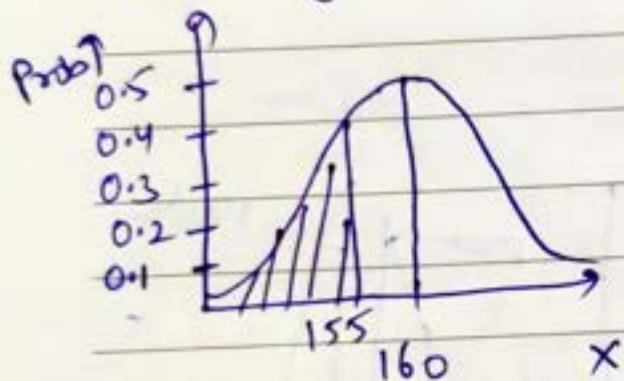


$(Q_2 - Q_1) > (Q_3 - Q_2)$ \rightarrow for left skewed

i.e. Mean < Median < Mode

→ Probability Distribution Function/Density Function

① Probability Density func" (PdF) → Distribution of the continuous data.
 eg Height of students [0 to 1]



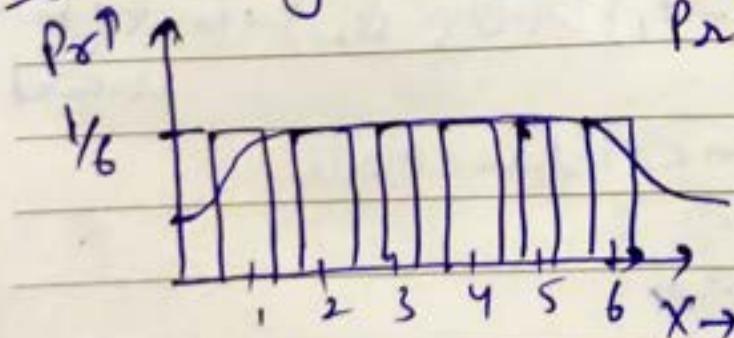
$\Pr(x < 155)$ = Area under the curve.

0.3, 0.4, 0.5

→ Kernel Density Estimator

② Probability Mass Function (PMF) → Distribution of a discrete random variable

eg Rolling a dice (1, 2, 3, 4, 5, 6)

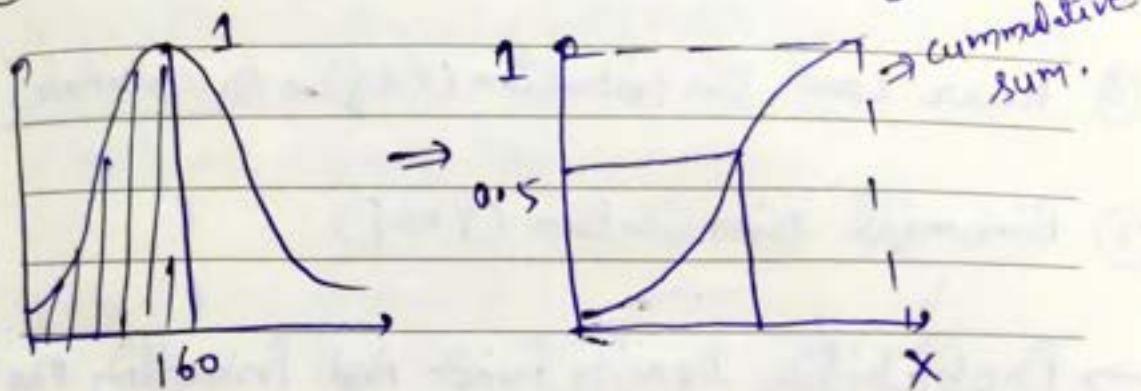


$$\Pr(x \leq 4) = \Pr(x=1) + \Pr(x=2) + \Pr(x=3) + \Pr(x=4)$$

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

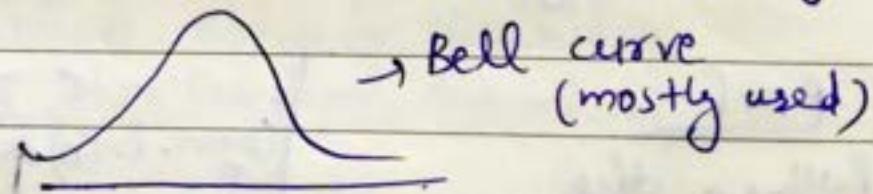
$$= \frac{4}{6} = \frac{2}{3}$$

③ Cumulative Distribution func (cdf) -



→ Types of Probability Distribution:-

① Normal/Gaussian Distribution → (in Pdf)



② Bernoulli Distribution → (in Pmf) → outcomes
are of success or failure

③ Uniform Distribution

④ Log Normal Distribution (Pd f)



Mo Tu We Th Fr Sa Su

Date / /

⑤ Poisson Distribution (Pmf)

⑥ Power Law Distribution (Pmf) \rightarrow 80-20% Rule

⑦ Binomial Distribution (Pmf)

\rightarrow Probability Density funcⁿ And Probability Mass Funcⁿ
Cumulative Distributive Funcⁿ (CDF)

① PMF :-

\rightarrow Used for discrete random variable

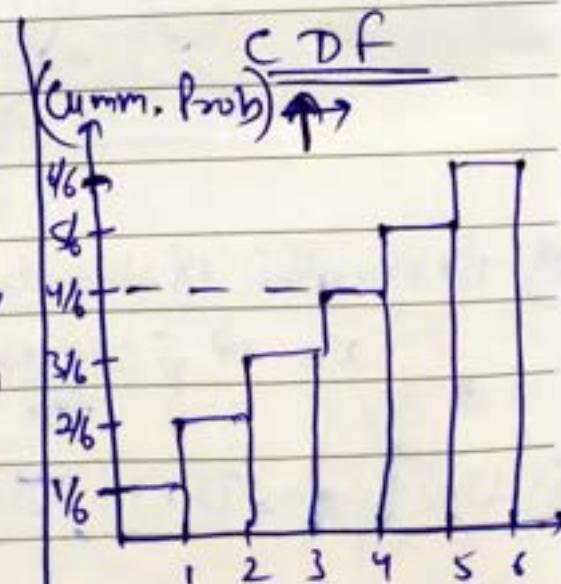
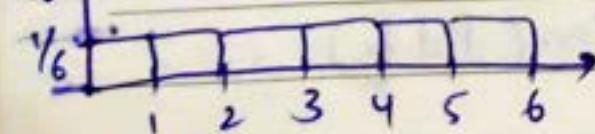
PMF

\rightarrow Rolling a dice

$$\{1, 2, 3, 4, 5, 6\}, P(1) = \frac{1}{6}, P(2) = \frac{1}{6} \dots P(6) = \frac{1}{6}.$$

$\frac{1}{6}$ $\frac{2}{6}$ $\frac{3}{6}$ $\frac{4}{6}$ $\frac{5}{6}$ $\frac{6}{6}$ \rightarrow uniform distribution

pmf



\rightarrow Rolling

$$P(X=1) = \frac{1}{6}, P(X=2) = \frac{1}{6}$$

$$\begin{aligned} P(X \leq 4) &= P(X=1) + P(X=2) + \\ &\quad P(X=3) + P(X=4) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} = \frac{4}{6} = \frac{2}{3}. \end{aligned}$$

$$\begin{aligned} P(X \leq 6) &= P(X=1) + P(X=2) + \\ &\quad \cdots \cdots \cdots P(X=6) \\ &= \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} + \frac{1}{6} \\ &= \frac{6}{6} = 1 \end{aligned}$$

But in CDF,

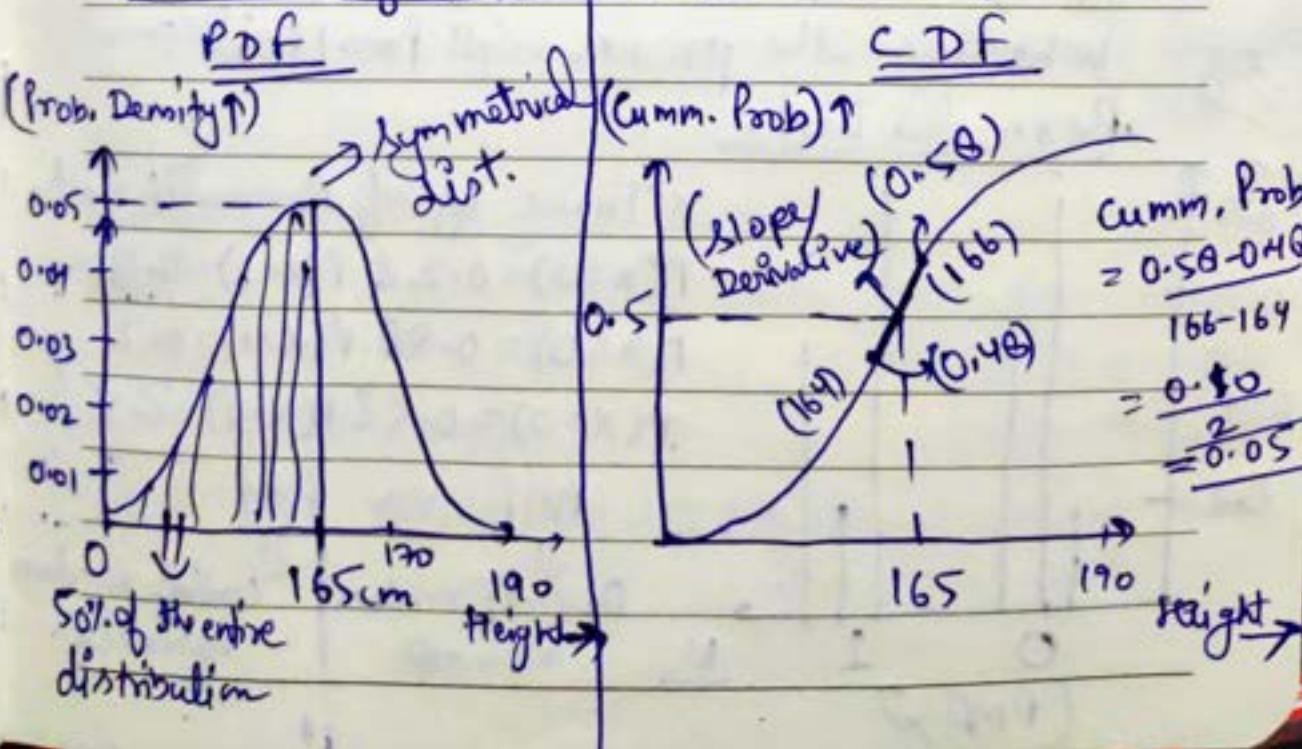
$$P(X \leq 4) = \frac{4}{6} = \frac{2}{3}.$$

$$P(X \leq 6) = \frac{6}{6} = 1.$$

② Probability Density func (PDF) :-

→ Continuous random variables

Prob. Density \rightarrow Gradient descent of CDF





Mo Tu We Th Fr Sa Su

Date / /

→ Bernoulli Distribution :- [Binary O/c]

In probability theory and statistics, the Bernoulli distribution, named after Swiss mathematician Jacob Bernoulli, is the discrete probability distribution of a random variable which takes the value 1 with prob. p and the value 0 with prob. $q=1-p$.

i.e. Outcomes are Binary {0, 1}

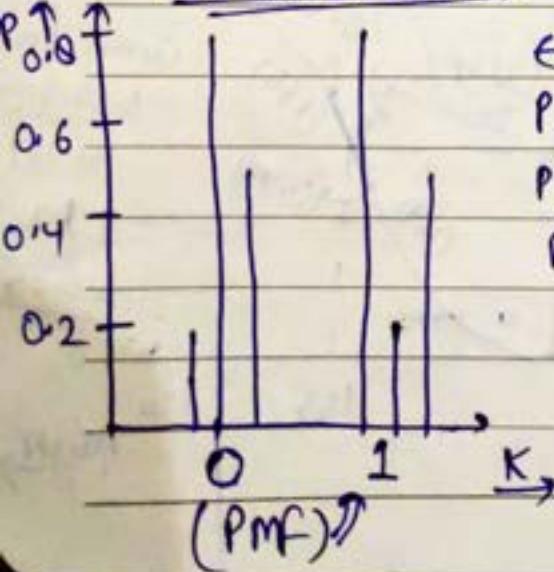
e.g. Tossing a coin. { $\frac{0}{H}$, $\frac{1}{T}$ }

$$P(T) = \frac{1}{2} = 0.5 = p$$

$$P(H) = 1 - p = 1 - 0.5 = 0.5 = q.$$

e.g. Whether the person will Pass/Fail →

Bernoulli Dist.



← Three eg of Bernoulli dist.

$$P(X=0) = 0.2 \text{ & } P(X=1) = 0.8$$

$$P(X=0) = 0.8 \text{ & } P(X=1) = 0.2$$

$$P(X=0) = 0.5 \text{ & } P(X=1) = 0.5$$

PMF vs Pdf

Discrete random variable

Continuous random variable.

PMF $\rightarrow K=0 \text{ or } 1$

$$\boxed{P(X=K) = P^K (1-P)^{1-K}}$$

$$\textcircled{1} \quad P(X=1) = P^1 (1-P)^{1-1} \\ = P^1 (1-P)^0 = P //$$

$$\textcircled{2} \quad P(X=0) = P^0 (1-P)^{1-0} \\ = (1-P)^1 = q //$$

Simplified way of showing PMf \rightarrow

$$\text{pmf} = \begin{cases} q = 1-p & , \text{if } K=0 \\ p & , \text{if } K=1 \end{cases}$$

\rightarrow Mean, Variance & Std. Deviation of Bernoulli Dist.

Mean $E(K) = \sum_{i=1}^K K \cdot p(K)$

$K=1 \text{ or } 0$

$$1 \times 0.6 + 0 \times 0.4 \\ = 0.6 = p //$$

$$P(K=1) = 0.6 = p //$$

$$P(K=0) = 1 - p = 1 - 0.6 = 0.4 = q //$$

Median of Bernoulli Dist.

$$\text{Median} = \begin{cases} 0 & \text{if } p < \frac{1}{2} \\ (0, 1) & \text{if } p = \frac{1}{2} \\ 1 & \text{if } p > \frac{1}{2} \end{cases}$$

Variance & Std Deviation →

$$\text{Variance} = p(1-p) = p \cdot q$$

$$\text{Std} = \sqrt{p \cdot q}$$

→ Binomial Distribution → The binomial dist. with parameters n and p is the discrete probability distribution of the no. of successes in a sequence of n independent experiments, each asking a yes-no question, and each with its own Boolean-valued outcome: success (with prob. p) or failure (with prob. $q = 1 - p$). .

→ A single success/failure exp. is also called a Bernoulli trial or Bernoulli exp. ($n=1$).

eg Tossing a coin {Bernoulli Distribution}

$$P(H) = 0.5 = 1-p$$

↓
0

$$P(T) = 0.5 = p$$

↓
1

eg Tossing a coin for 10 times.

1st time 2nd 3rd 4th

$$\begin{cases} P(T) = p & \cdot p & p & p \\ P(H) = 1-p & 1-p & 1-p & 1-p \end{cases}$$

Binomial Distribution

① Parameters → $n \in \{0, 1, 2, 3\}$ → no. of trials or exp.

$p \in [0, 1]$ → success or prob. for each trial

$$q = 1-p$$

② PMF $P(K) = {}^n C_K p^K \cdot (1-p)^{n-k}$, $K \rightarrow$ no. of successes

③ Mean of Binomial Dist.

$$\text{mean} = n \cdot p$$

④ Variance and Std.

$$\text{Variance} = npq$$

$$\text{Std} = \sqrt{npq}$$



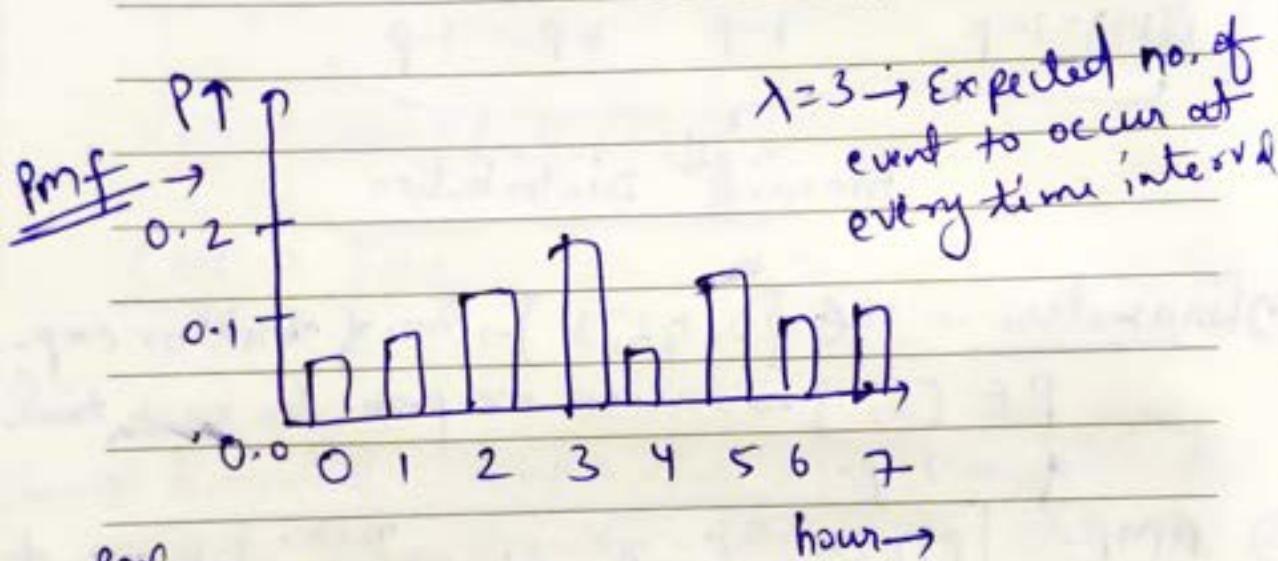
Mo Tu We Th Fr Sa Su

Date / /

→ Poisson Distribution :-

- Discrete Distribution (Pmf)
- Describes the no. of events occurring in a fixed time interval.

e.g. No. of people visiting hospital every hour.
bank _____
airport _____.



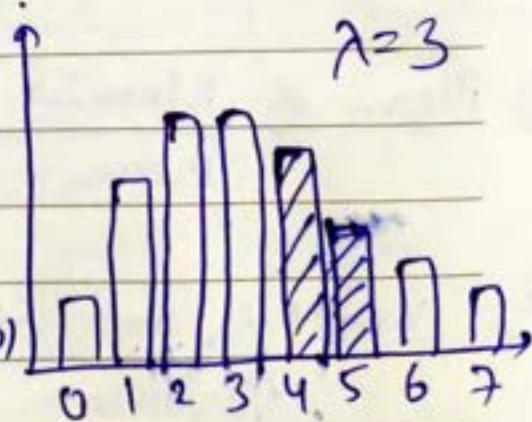
$$\text{pmf}$$

$$P(X=s) = \frac{e^{-\lambda} \lambda^s}{s!}$$

$$= \frac{e^{-3} 3^5}{5!} = \frac{0.101}{10.11} = 0.101$$

$$P(X=4) + P(X=5) = \frac{e^{-3} 3^4}{4!} + 0.101$$

$$= \underline{\hspace{2cm}}$$



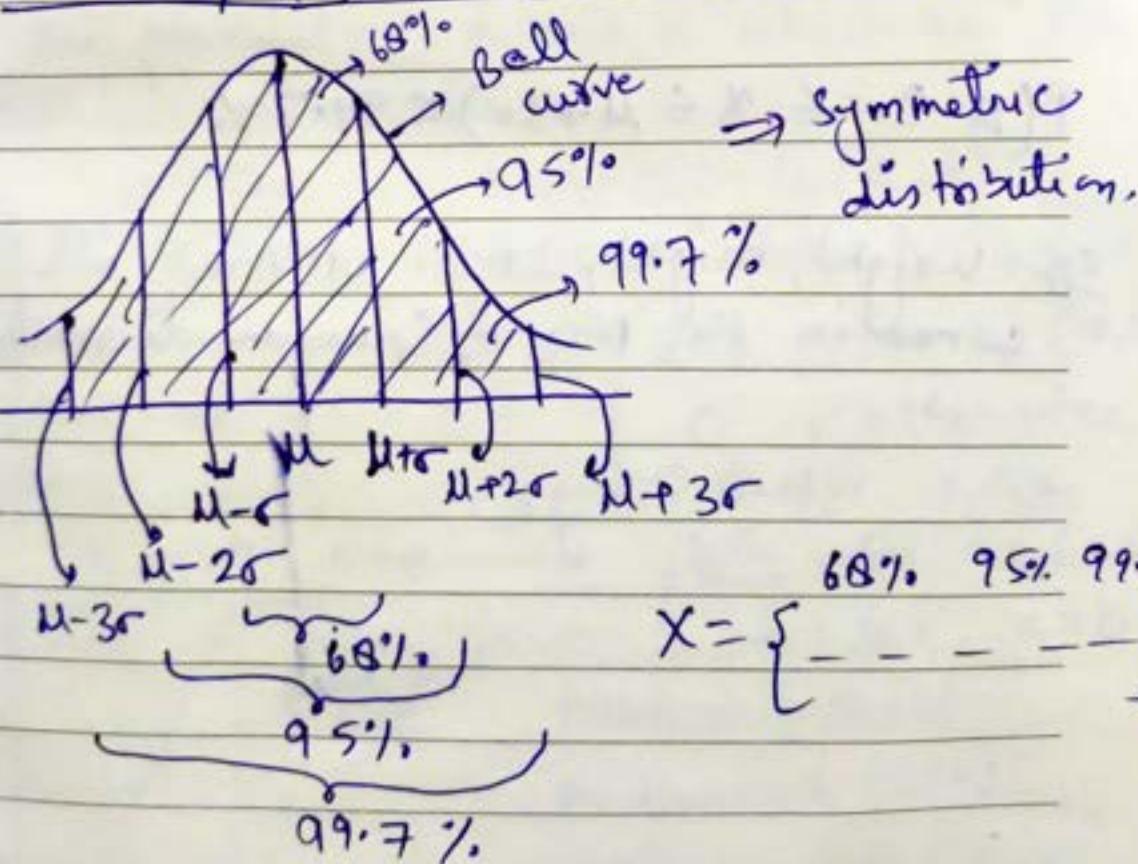
Mean & Variance of Poisson Distr.

$$\text{Mean} = E(X) = \mu = \lambda \times t$$

$$\text{Variance} = E(X) = \mu = \lambda \times t$$

$\lambda \rightarrow$ Expected no. of events to occur at every time interval
 $t \rightarrow$ time interval

→ Normal/Gaussian Distribution :-





Mo Tu We Th Fr Sa Su

Date / /

- Empirical Rule [3-Sigma Rule]

68% - 95% - 99.7%

Q-Q Plot ^{To check} → Whether a distribution is Gaussian/Normal Distribution.

Probability

$$P(\mu - \sigma \leq X \leq \mu + \sigma) \approx 68\%$$

$$P(\mu - 2\sigma \leq X \leq \mu + 2\sigma) \approx 95\%.$$

$$P(\mu - 3\sigma \leq X \leq \mu + 3\sigma) \approx 99.7\%.$$

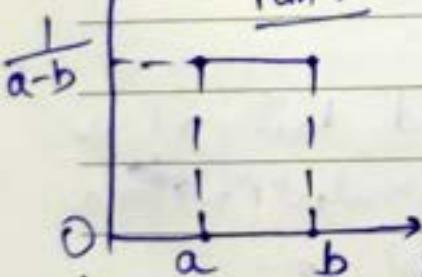
e.g. weight, height, IRIS DATASET all come in the Normal/Gaussian distribution.

→ Uniform Distribution :-

- ① Continuous Uniform Dist. (pdf)
- ② Discrete " " (pmf)

① Continuous Uniform Dist. → (Cont. Random Variable)
 Cont. Unif. Dist or rectangular dist.
 is a family of symmetric prob. distributions.
 The distribution describes an exp. where
 there is an arbitrary outcome that lies
 b/w certain bounds. The bounds are defined
 by the parameters a and b which are the
 min. and max. values.

Uniform
 $f(x)$ Prob. Density Funⁿ.



Notation :- $U(a, b)$

Parameter $\rightarrow -\infty < a < b < \infty$

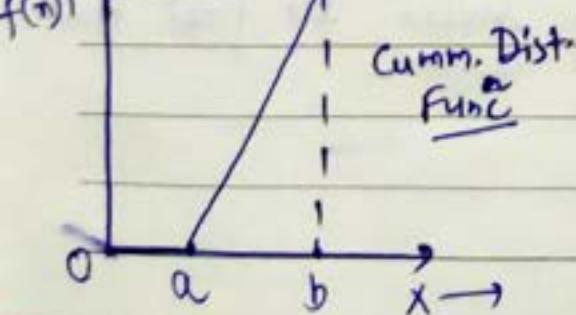
$$\text{Pd.f} = \begin{cases} \frac{1}{b-a}, & \text{for } x \in [a, b] \\ 0, & \text{otherwise} \end{cases}$$

$$\text{Cdf} = \begin{cases} 0 & \text{for } x < a \\ \frac{x-a}{b-a}, & \text{for } x \in [a, b] \\ 1, & \text{for } x > b \end{cases}$$

$$\text{Mean} = \frac{1}{2}(a+b)$$

$$\text{Median} = \frac{1}{2}(a+b)$$

$$\text{Variance} = \frac{1}{12}(b-a)^2$$





Mo Tu We Th Fr Sa Su

Date / /

Uniform Distribution →

e.g. Rolling a dice $\{1, 2, 3, 4, 5, 6\}$

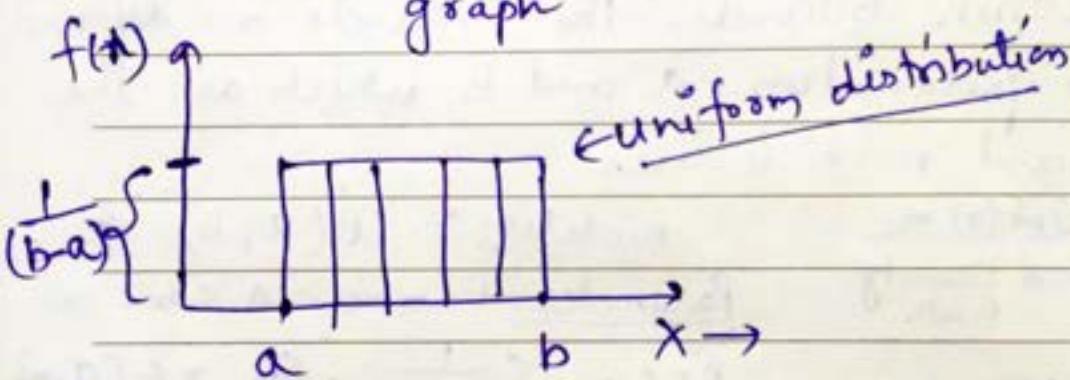
$$P(1) = \frac{1}{6}, P(2) = \frac{1}{6}, P(3) = \frac{1}{6}$$

In uniform dist. the probability of getting the o/c is equal.

A variable X is said to be uniformly distributed.

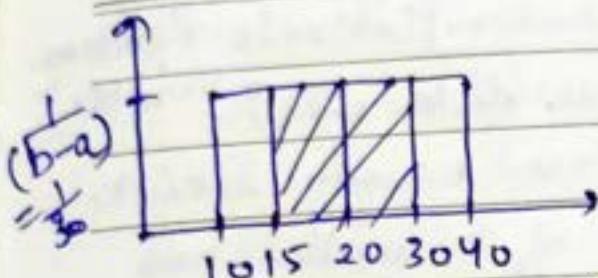
$$f(x) = \frac{1}{b-a} \text{ for } -\infty \leq a \leq b < \infty$$

↓
graph



- ✳ The no. of candies sold daily at a shop is uniformly distributed with a max. of 40 and minimum of 10.

i) Probability of daily sales to fall b/w 15 and 30.



$$x_1 = 15, a = 10$$

$$x_2 = 30, b = 40$$

$$\frac{1}{b-a} = \frac{1}{40-10} = \frac{1}{30}$$

$$P(15 \leq x \leq 30) = (x_2 - x_1) * \frac{1}{b-a}$$

$$= (30 - 15) * \frac{1}{30}$$

$$= 15 * \frac{1}{30} \Rightarrow 0.5 = \underline{\underline{50\%}}$$

$$P(x \geq 20) = (40 - 20) * \frac{1}{30}$$

$$= 20 * \frac{1}{30} = \frac{2}{3} = 0.666$$

$$= \underline{\underline{66.6\%}}$$

Mean and Variance of Uniform Distributions

$$\text{Mean} \rightarrow E(x) = [a+b]/2 = (40+10)/2 = \frac{50}{2} = 25.$$

$$\text{Variance} \rightarrow V(x) = (b-a)^2/12 = \frac{(40-10)^2}{12} \Rightarrow \frac{30^2}{12}$$

$$\Rightarrow \frac{900}{12} = \underline{\underline{75}}$$

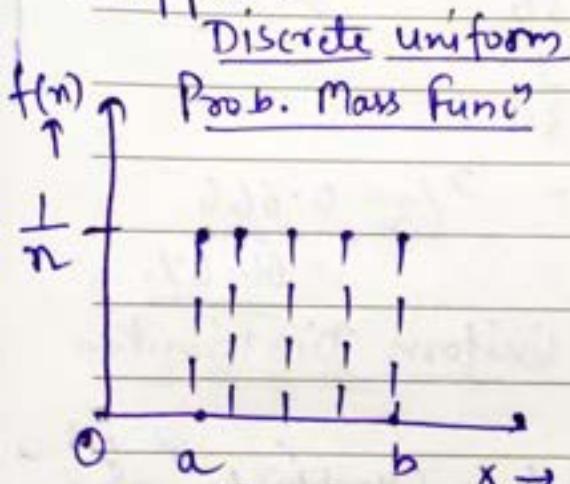


Mo Tu We Th Fr Sa Su

Date / /

② Discrete Uniform Distribution (Discrete Random Variable)

→ It is a symmetric prob. dist., wherein a finite no. of values are equally likely to be observed; every one of n values has equal prob. $1/n$. Another way of saying 'discrete uniform dist' would be 'a known, finite no. of outcomes equally likely to happen'.



$$n=5, \text{ when } n=b-a+1$$

Eg rolling a dice

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

$$P(1)=\frac{1}{6}, P(2)=\frac{1}{6}, P(3)=\frac{1}{6}$$

$$a=1, b=6$$

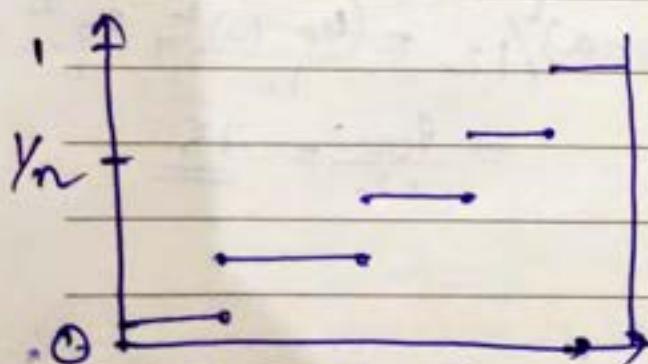
$$n=b-a+1=6 //$$

Notation $\rightarrow u(a, b)$

Parameters $\rightarrow a & b, b \geq a$

$$\text{PMF} = \frac{1}{n}$$

$$\text{Mean} = \frac{a+b}{2}$$



Cumm. Dist. func'



Mon Tue Wed Thu Fri Sat Sun

(SND)

Date / /

→ Standard Normal Distribution → when we convert Normal/Gaussian distribution mean (μ) and variance (σ^2) as $\mu=0$ and $\sigma^2=1$ called as Standard Normal Distribution. Conversion of this possible by Z-score.

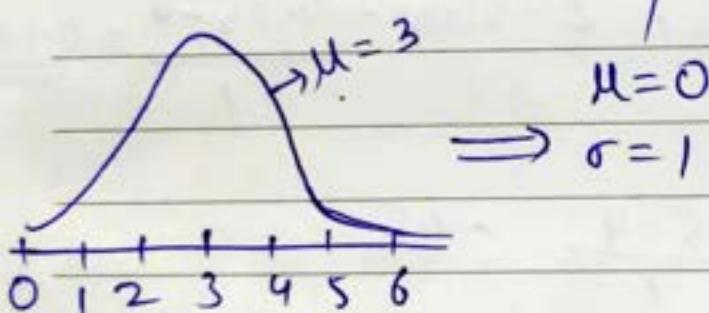
$$X = \{1, 2, 3, 4, 5\}$$

$$\mu = 3$$

$$\sigma = 1.414 \approx 1$$

Notation

$$X \sim SND(\mu=0, \sigma=1)$$



$$X = \{1, 2, 3, 4, 5\}, Z\text{-score} = \frac{x_i - \mu}{\sigma}, y$$

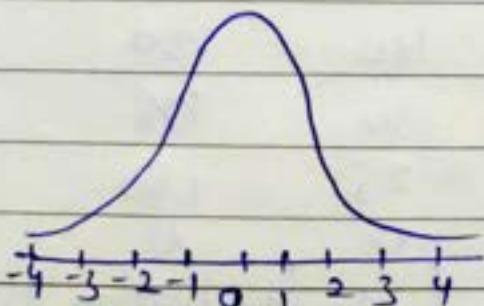
$$1 \rightarrow \frac{1-3}{1} = -2 \quad y = \{-2, -1, 0, 1, 2\}$$

$$2 \rightarrow \frac{2-3}{1} = -1$$

$$3 \rightarrow \frac{3-3}{1} = 0$$

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

$$5 \rightarrow \frac{5-3}{1} = 2$$



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

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

$$z \rightarrow 1$$

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

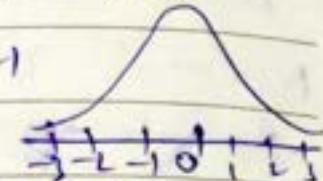
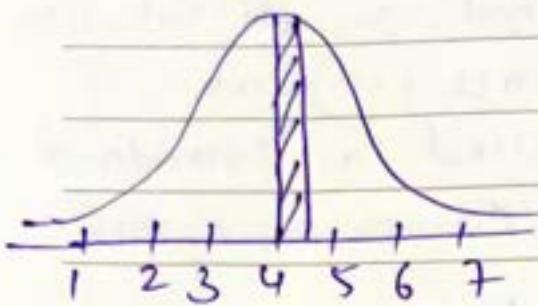
$$6 \rightarrow 2$$

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

$$7 \rightarrow 3$$

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

$$4 \rightarrow 0$$



Q How many Std. deviation 4.25 is away from the mean?

Ans $x_i = 4.25$

using z-score, $z\text{-score} = \frac{x_i - \mu}{\sigma} = \frac{4.25 - 4}{1} = 0.25$

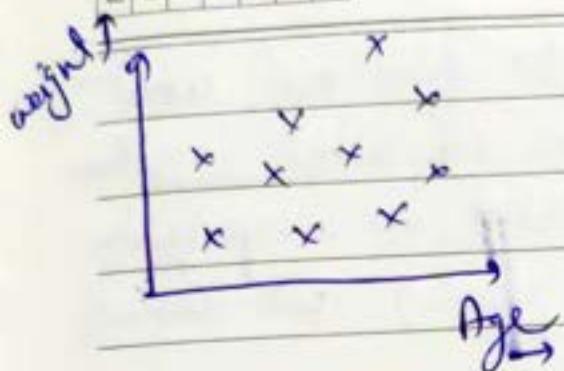
Q $x_i = 2.5$

Ans $z\text{-score} \Rightarrow \frac{2.5 - 4}{1} = -1.5$

Eg \rightarrow Dataset (kg) (cm) (INR)

Age (yr)	Weight	Height	Salary
24	70	175	40K
25	60	160	50K
26	55	150	60K
27	40	130	30K
30	30	175	20K
31	25	180	70K
11	11	11	11

(Standardization) using z-score



$$z\text{-score} = \frac{x_i - \mu_{\text{age}}}{\sigma}, \quad \frac{x_i - \mu_{\text{weight}}}{\sigma}$$

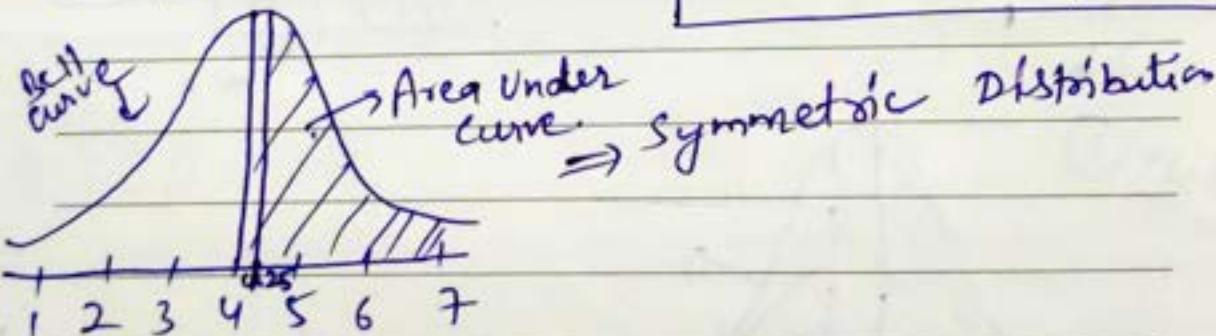
→ problem statement on z-score :-

Z-table

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

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

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



Q) What % of scores fall above 4.25?

$$z\text{-score} = \frac{x_i - \mu}{\sigma} \Rightarrow \frac{4.25 - 4}{1}$$

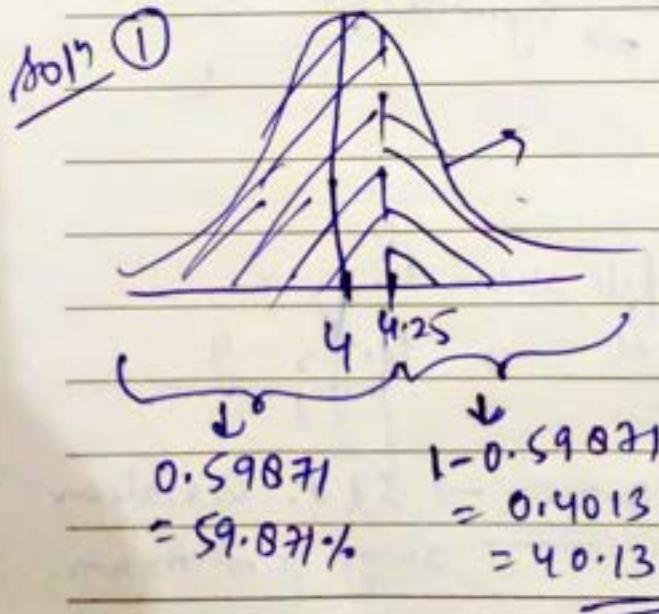
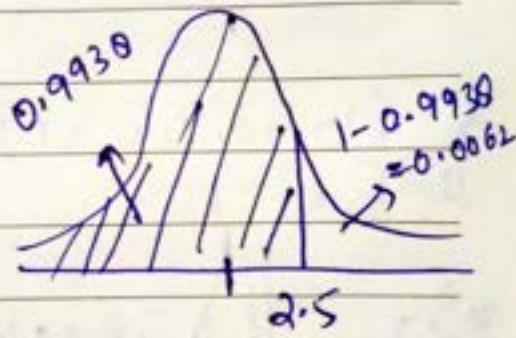
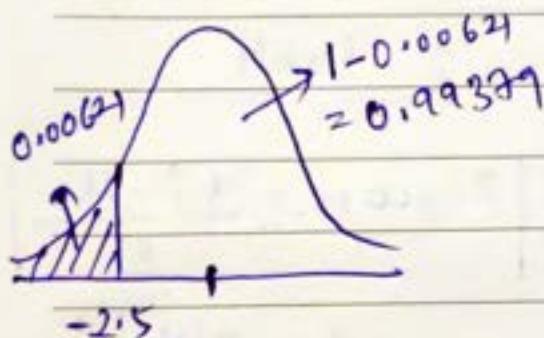
= 0.25 → std. Deviation away from mean.

Note → When we need to find out what is the % of distribution or what is the area under the curve based on a specific std. deviation (derive from z-score) we use z-table.

→ There are two types of z-table →

① Negative z-table

② Positive z-table



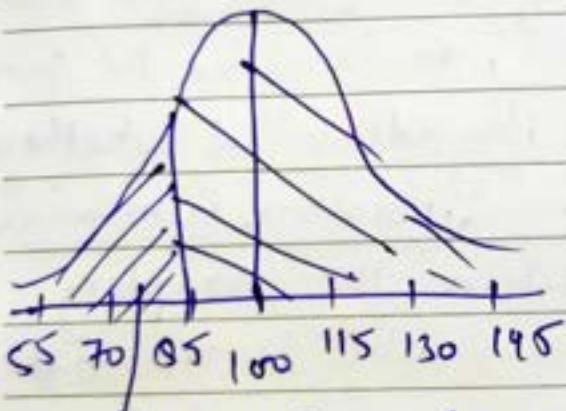


Maths English Science Social Studies

Date / /

Q) In India the avg IQ is 100, with a std. deviation of 15. What is the % of the pop^ would you expect to have an IQ lower than 85?

Sol) $\mu = 100, \sigma = 15$



Z-score

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

$$= \frac{85 - 100}{15} \Rightarrow \frac{-15}{15} \\ = -1$$

on Z-table using z-score = -1.0
we get = $0.15866 = \underline{\underline{15.866\%}}$

Q) $IQ > 85$

$$\Rightarrow 1 - 0.15866 = 0.84134 \Rightarrow \underline{\underline{84.134\%}}$$

Q) ~~75~~ $75 \geq IQ \leq 100$

$$Z\text{-score} = \frac{75 - 100}{15} = \frac{-25}{15} \Rightarrow -1.666, 0.0485 \\ = 4.85\%$$

$$So - 4.85\% = \underline{\underline{45.15\%}}$$



Mo Tu We Th Fr Sa Su

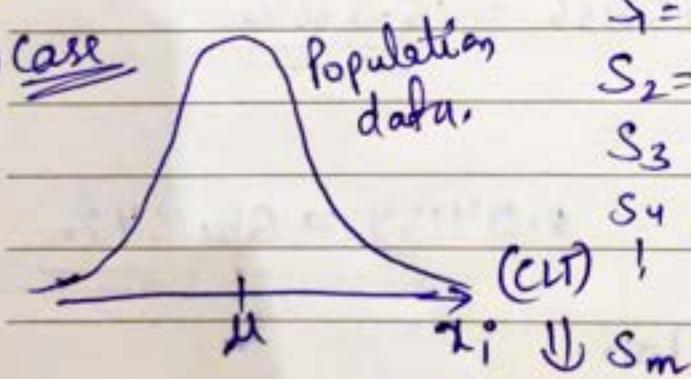
Date / /

→ Central Limit Theorem → It relies on the concept of a sampling distribution, which is the probability of a statistic for a large no. of samples taken from a population. The central limit theorem says that the sampling distribution of the mean will always be normally distributed, as long as the sample size is large enough. Regardless of whether the popⁿ has a normal, Poisson, binomial or any other distribution, the sampling distribution of the mean will be normal.

$$X \sim N(\mu, \sigma)$$

$n \rightarrow$ sample size $\xrightarrow{\text{any value}}$

① Case



$$S_1 = \{x_1, x_2, x_3, \dots, x_n\} = \bar{x}_1$$

$$S_2 = \{x_2, x_3, x_4, \dots, x_n\} = \bar{x}_2$$

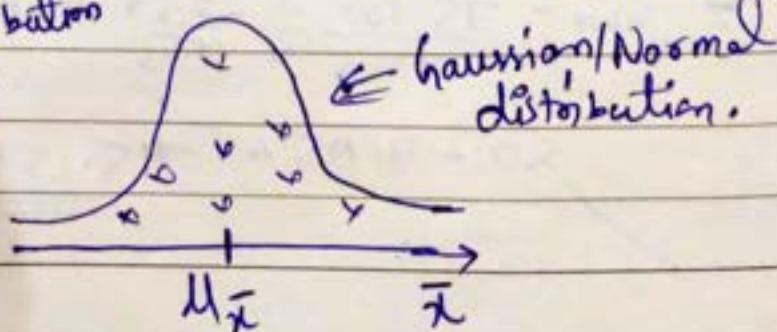
$$S_3 = \bar{x}_3$$

$$S_4 = \bar{x}_4$$

$$\vdots$$

$$\bar{x}_n$$

Sampling distribution
of the mean.



② case $X \not\sim N(\mu, \sigma)$

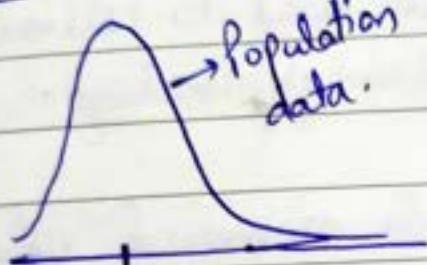
$[n \geq 30] \rightarrow$ sample size

$s_1 \rightarrow \bar{x}_1$

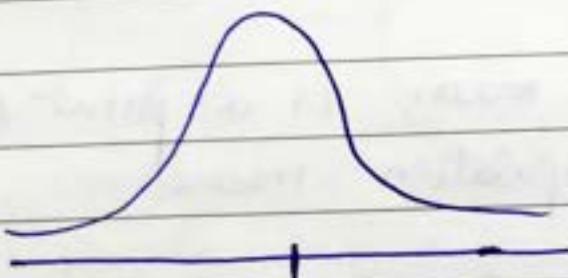
$s_2 \rightarrow \bar{x}_2$

:

$\frac{1}{\bar{x}_m}$

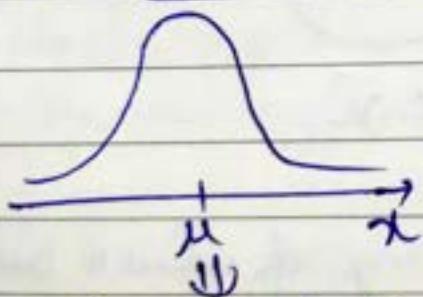


$\cup_{(CLT)}$



Normal Distribution

$$X \sim N(\mu, \sigma)$$



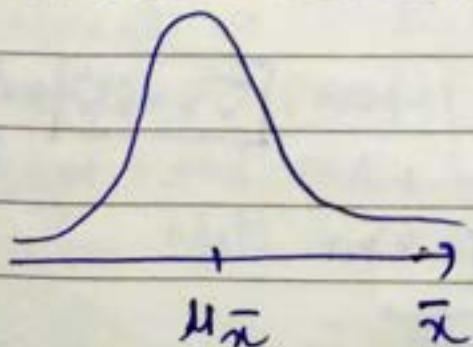
$\sigma \rightarrow$ population std. devia^{tion}

Sampling Distribution of mean

$\mu \rightarrow$ pop'n mean

$n \rightarrow$ sample size

$$X \sim N\left(\mu, \frac{\sigma}{\sqrt{n}}\right)$$



Inferential Statistics



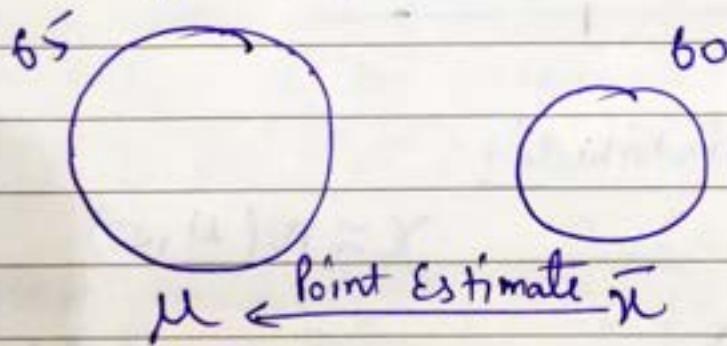
Mo Tu We Th Fr Sa Su

Date / /

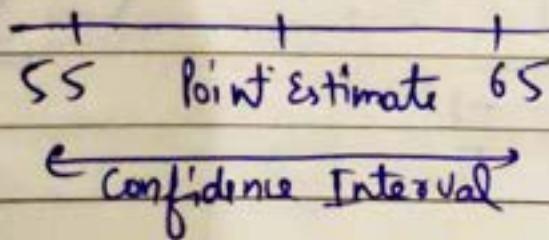
⇒ Estimate → It is a specified observed numerical value used to estimate an unknown population parameter.

Types :-

- ① Point Estimate → Single numerical value used to estimate an unknown population parameter.
eg Sample mean is a point estimate of a population mean



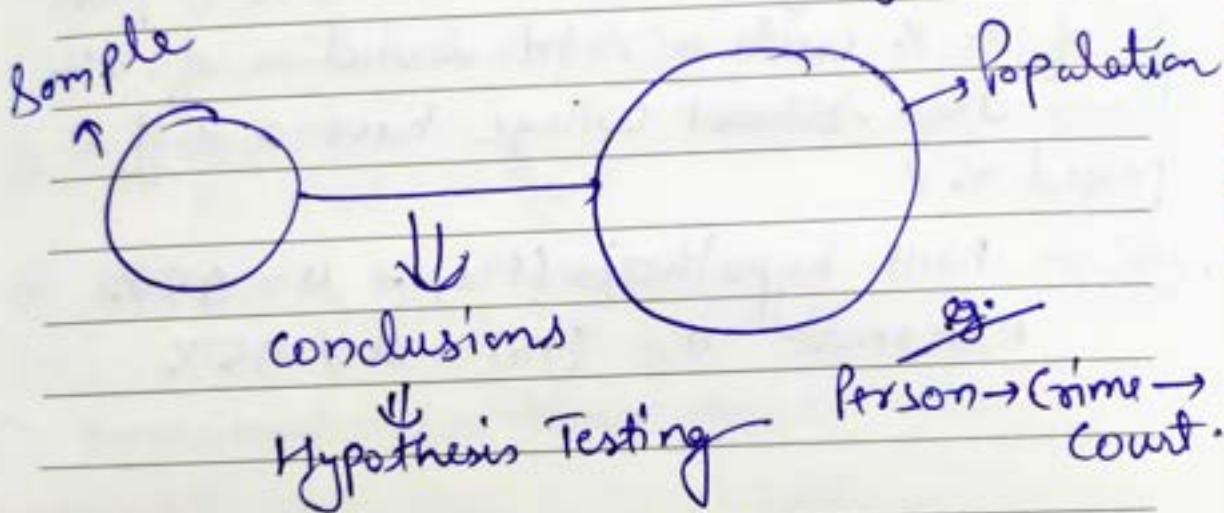
- ② Interval Estimate → Range of values used to estimate the unknown population parameter



$[55-65]$ → Sample mean

Hypothesis And Hypothesis Testing Mechanism

Inferential statistics → Conclusion or Inference.



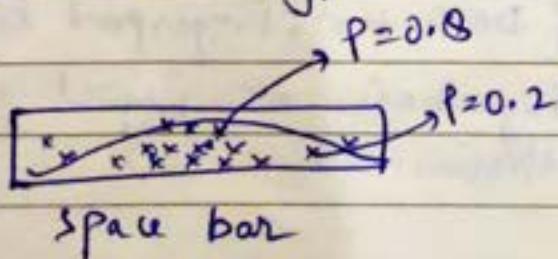
Hypothesis Testing Mechanism :-

- ① Null Hypothesis (H_0) - Person is not guilty until proved
- The assumption you are beginning with
- ② Alternative Hypothesis (H_1) - The person is guilty
- Opposite of null hypothesis
- ③ Experiments → Statistical Analysis.
→ Collect proofs (DNA test, fingerprint test)
- ④ Accept the Null hypothesis or Reject the Null hypothesis.

Eg College at district A states its avg passed % of students are 85%. A new college opened in the district and it was found that sample of student too have a pass % of 90% with a std. deviation of 4%. Does this ~~school~~ college have a diff. passed %?

Sol Null hypothesis (H_0) = $\mu = 85\%$.
 Alternate " (H_1) = $\mu \neq 85\%$.

→ P Value - It is a no., calculated from a statistical test, that describes how likely you are to have found a particular set of observations if the null hypothesis was true. P values are used in hypothesis testing to help decide whether to reject the null hypothesis.



out of 100 touches, we touch around 20 times in right side region and 80 times in middle region.



Mo Tu We Th Fr Sa Su

Date _____

Hypothesis Testing

Eg coin is fair or Not [100 times tossed]

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

$$P(H) = 0.6 \quad P(T) = 0.4$$

$$P(H) = 0.7 \quad P(T) = 0.3$$

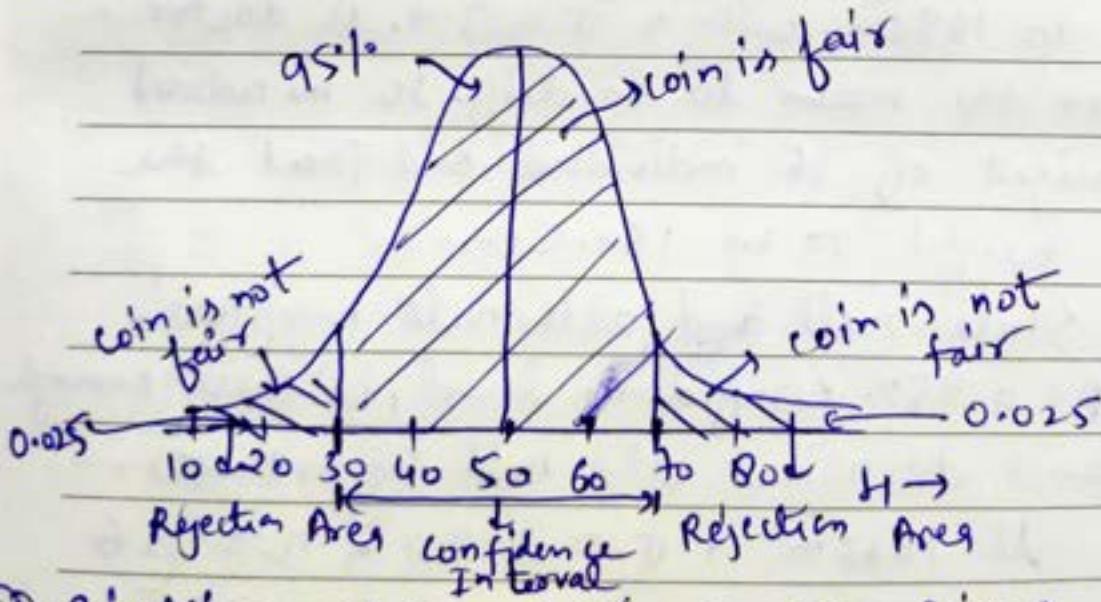
① Null hypothesis

$H_0 \rightarrow$ coin is fair

② Alternate hypothesis

$H_1 \rightarrow$ coin is not fair

③ Experiment \rightarrow 100 times tossed.

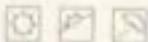


④ Significance Value : $\alpha = 0.05$ (given)

$P <$ Significance Value

\rightarrow Reject the Null hypothesis

else fail to reject the null hypothesis



Mo Tu We Th Fr Sa Su

Date / /

Hypothesis Testing and Statistical Analysis

- ① z test $\bar{y} \rightarrow$ Average \Rightarrow z table \rightarrow z score
- ② t test $\bar{y} \rightarrow$ Average \Rightarrow t table & (p value)
- ③ CHI SQUARE \Rightarrow Categorical data
- ④ ANNOVA \Rightarrow Variance

z test \rightarrow i) Population St.D. ii) $n \geq 30$
known.

eg) The avg heights of all residents in a city is 168cm. with a $\sigma = 3.9$. A doctor believes the mean to be diff. He measured the height of 36 individual and found the avg height to be 169.5cm.

- a) State null and alternate hypothesis
- b) At a 95% confidence level, is there enough evidence to reject the null hypothesis.

Ans $M = 168 \text{ cm}$, $\sigma = 3.9$, $n = 36$

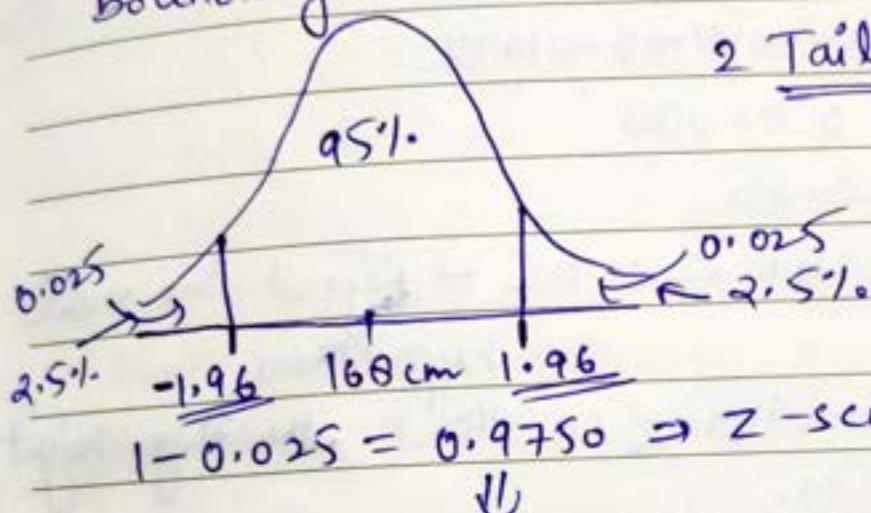
$$\bar{x} = 169.5 \text{ cm}, CI = 0.95$$

$\alpha = 1 - CI = 1 - 0.95 = 0.05$
(significance value)

a) Null hypothesis $H_0 = \mu = 160\text{ cm}$
 Alternative " $H_1 = \mu \neq 160\text{ cm.}$

b) Based on CI we will draw Decision boundary.

2 Tail test.



$$\text{Area} = +1.96.$$

if z is less than -1.96 or greater than $+1.96$, Reject the Null hypothesis.

* z -test

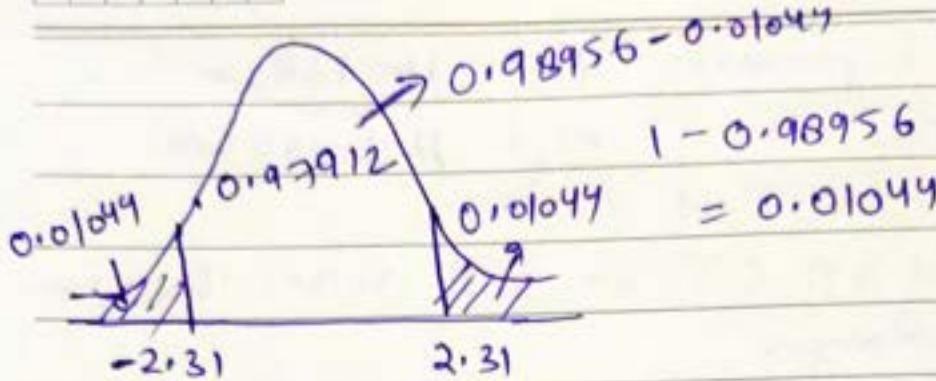
$$\begin{aligned} z_d &= \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{169.5 - 168}{3.9/\sqrt{36}} \\ &= \frac{1.5}{0.65} = 2.31 \end{aligned}$$

$2.31 > 1.96$, Reject the null hypothesis.
 $P < 0.05$



Mo Tu We Th Fr Sa Su

Date / /



$$\begin{aligned} P \text{ value} &= 0.01044 + 0.01044 \\ &= 0.02088 \end{aligned}$$

if $P < 0.05$

$0.02088 < 0.05 \rightarrow$ Reject the null hypotheses

* final conclusion is that the Avg. height $\neq 160 \text{ cm.}$

- Based on doctor sample data the avg height of seems to increase.

- Q A factory manufactures bulbs with a avg warranty of 5 years with std. deviation of 0.50. A worker believes that the bulb will malfunction in less than 5 year. He test a sample of 40 bulbs and find the avg. time to be 4.8 years.
- a) State null and alternate hypothesis

v) At a 2% significance level, is there enough evidence to support the idea that warranty should be revised?

~~10P~~ $\mu = 5 \text{ yrs}, \sigma = 0.50, n = 40$
 $\bar{x} = 4.8 \text{ yrs}, S.E = 0.02$

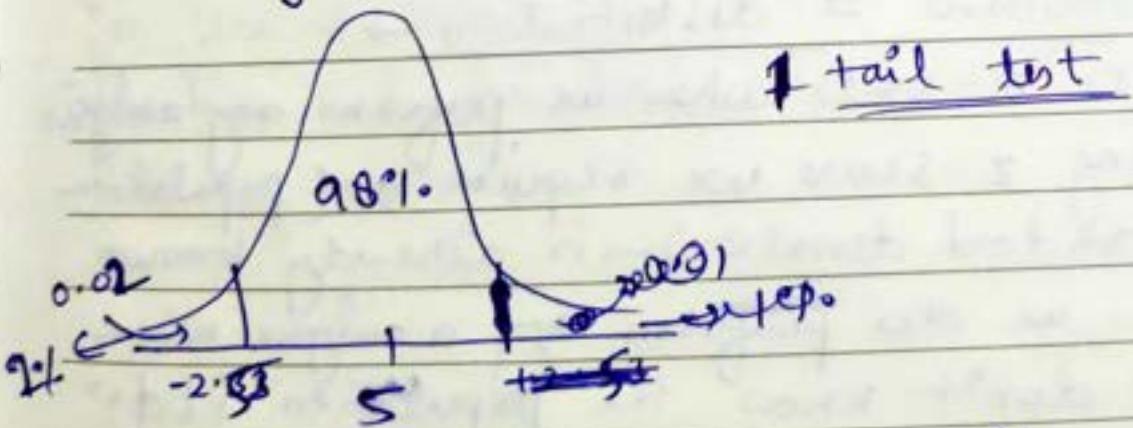
$$\alpha = 1 - CI$$

$$0.02 = 1 - CI \Rightarrow CI = 1 - 0.02$$

$$CI = 0.98$$

a) Null hypothesis (H_0) = $S_y \approx \mu$
 Alternate " (H_1) = $\mu \neq S_y$ (1-tail test)

b) Based on CI we should draw decision boundary.



$$1 - 0.01 = 0.99 \Rightarrow z\text{-score (from table)}$$

$$\Rightarrow A_{req} = +2.33$$



Mo Tu We Th Fr Sa Su

Date / /

if $-2.33 > Z_{\alpha=0.05}$, Reject the null hypothesis

* Z test

$$Z_d = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{4.8 - 5}{0.50/\sqrt{40}} = -0.2 \\ = -2.5300$$

Area under the curve with z score

$$-2.5300 \Rightarrow 0.0570$$

$$p\text{-value} = 0.0570$$

Compare p value with significance value (α)

$$0.0570 < 0.02 \Rightarrow \text{false}$$

* Conclusion → We accept the null hypothesis
or we failed to reject the null hypothesis.

→ Student t distribution →

In z stats when we perform any analysis using z-score we require σ (population standard deviation) → it already known.
Now we do perform any analysis when we don't know the population stdv



Student's t distribution

t stat

Date

$$z = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}}$$

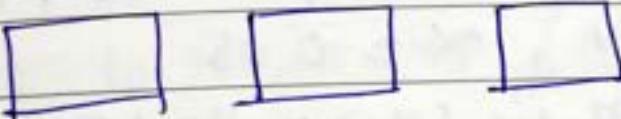
$$t = \frac{\bar{x} - \mu}{s/\sqrt{n}}$$

's → Sample std.
deviation

Degree of freedom

$$dof = n - 1$$

Ex There are 3 people in a room. and
there are 3 chairs in the room to sit.



if first person going to choose a chair
he will have 3 choice to him after getting
the chair only 2 chairs are empty. then
for second one there will be 2 choices
for him and for last one he will not
have any choice.



Mo Tu We Th Fr Sa Su

Date / /

Q In the population the avg IQ is 100. A team of researchers want to test a new medication to see if it has either a +ve or -ve effect on intelligence, or no effect at all. A sample of 30 participants who have taken the medication has a mean of 140 with a std deviation of 20. Did the medication affect intelligence? CI = 95%, $\alpha = 0.05$

Ans $H_0: \mu = 100, n = 30, \bar{x} = 140, s = 20$
 $CI = 95\%, \alpha = 0.05$

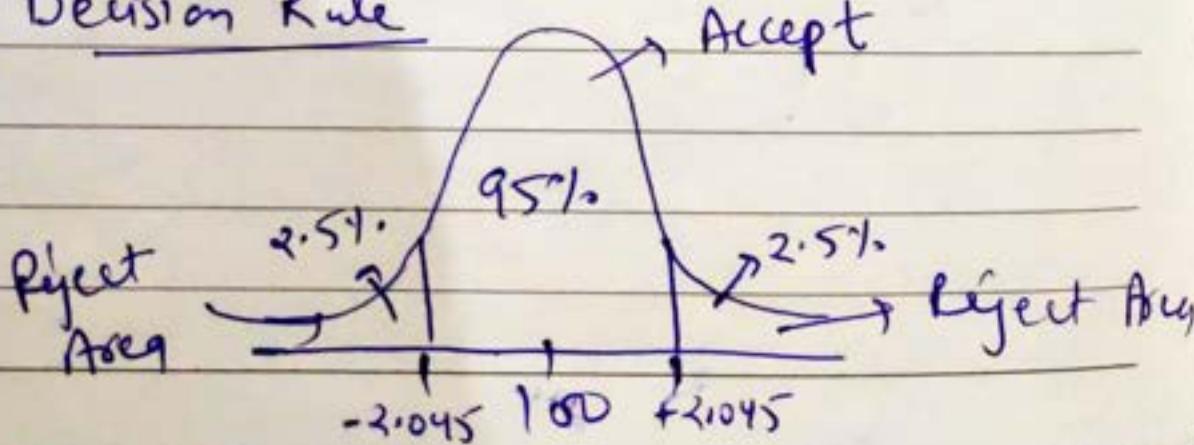
① Null hypothesis (H_0) = $\mu = 100$
Alternate " (H_1) = $\mu \neq 100$ (2 tail test)

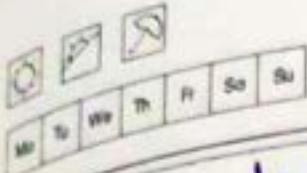
② $\alpha = 0.05$

③ Degree of freedom

$$dof = n - 1 = 30 - 1 \\ = 29$$

④ Decision Rule





Date / /

if t test is less than -2.045 or greater than 2.045, reject the null hypothesis

⑤ Calculate Test Statistics

$$t = \frac{\bar{x} - \mu}{\sigma/\sqrt{n}} = \frac{140 - 100}{20/\sqrt{30}} = \frac{40}{3.65} = 10.96$$

$$t = 10.96$$

Since, $t = 10.96 > 2.045$ {Reject the null hypothesis}

Conclusion:- Medication used has affect the intelligence

→ Medication has increased the Intelligence.

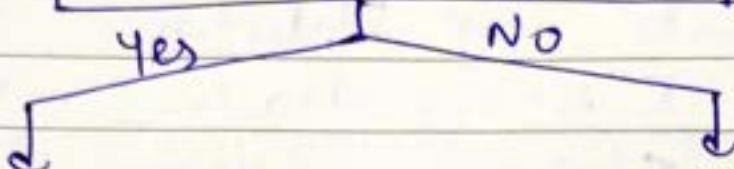


Mo Tu We Th Fr Sa Su

Date / /

→ When to use T-test vs Z-test

Do you know the
Population std. Deviation



Is the sample size above 30 ($n > 30$)

Yes | No

Use the z-test

Use the T-test

→ Type 1 & Type 2 Errors:-

Reality → 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
→ Good.

Outcome 2: we reject the Null hypothesis
when in reality it is True
→ Type 1 Error.

Outcome 3: we accept the Null hypothesis,
when in reality it is false.
→ Type 2 Error

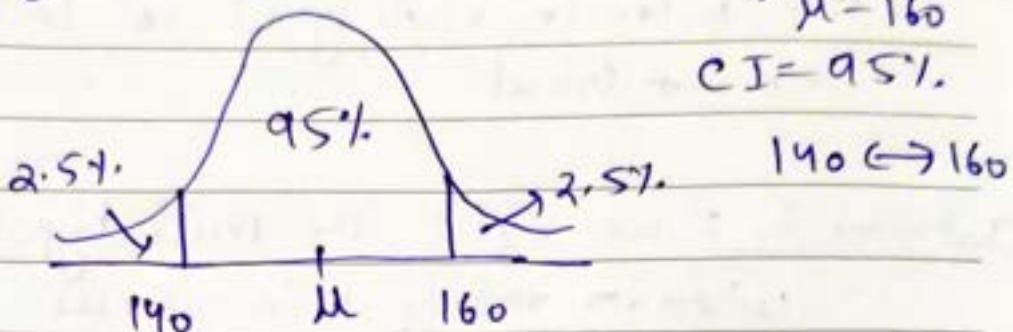
Outcome 4: we accept the Null hypothesis,
when in reality it is True
→ Good.



Mo Tu We Th Fr Sa Su

Date / /

→ Confidence Intervals and Margin of Error



Point Estimate → A value of any statistical that estimates the value of an unknown population parameter is called Point Estimate.

$$\bar{x} \longrightarrow \mu$$
$$\bar{x} = 2.95 \qquad n = 3.$$

Confidence Interval → We construct a confidence interval to help estimate what the actual value of the unknown population mean.

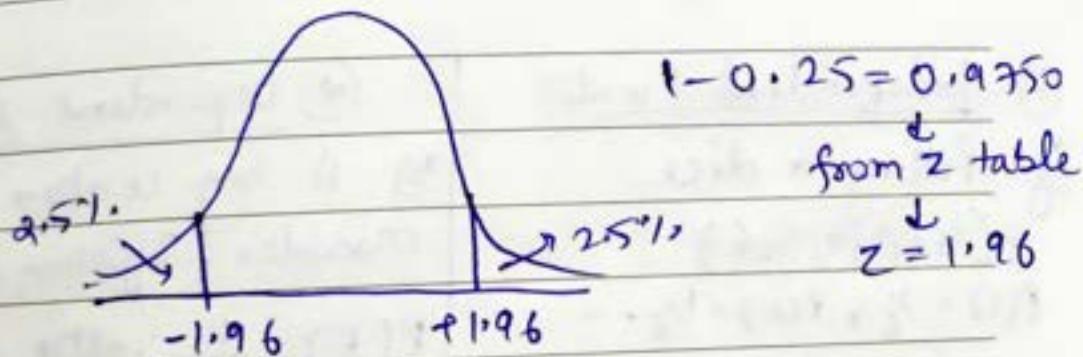
Point Estimate \pm Margin of Error.

z test
$$\boxed{\bar{x} \pm Z_{\alpha/2} \frac{\sigma}{\sqrt{n}}} \qquad \alpha = 0.05$$

$$Z_{0.05/2} \Rightarrow Z_{0.025}.$$

Q On the verbal section of CAT exam, the standard deviation is known to be 100. A sample of 25 test-takers has a mean of 520. Construct a 95% CI about the mean?

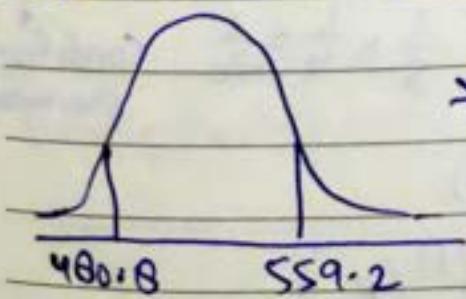
Ans $\bar{x} = 520, \sigma = 100, n = 25, CI = 0.95$
 $\alpha = 0.05$



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

$$\text{Lower CI} = 520 - (1.96) \times \frac{100}{\sqrt{25}} = 480.8$$

$$\text{Upper CI} = 520 + 1.96 \times \frac{100}{\sqrt{25}} = 559.2$$



* concluding → I am 95% confident that the mean CAT exam lies b/w 480.8 & 559.2.

⇒ Baye's Theorem → Bayesian statistics is an approach to data analysis and parameter estimation based on Baye's theorem.

Baye's theorem

Possibility → Independent Events
 ↳ Dependent Events.

① Independent Events

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

$$P(1) = \frac{1}{6}, P(2) = \frac{1}{6} \dots$$

Eg. Tossing a coin

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

→ In this all the events are independent of others.

② Dependent Events

Eg. A bag contains 5 marbles (3 yellow 2 Red)

$P(R) = \frac{2}{5}$. after getting out a red marble now

$P(Y) = \frac{3}{4}$ as a red marble is taken out.

$$P(R \text{ and } Y) = P(R) \times P(Y|R)$$

$$= \frac{2}{5} \times \frac{3}{4} = \frac{6}{20}$$

↓
Conditional probability.

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

$$P(A) \times P(B/A) = P(B) \times P(A/B)$$

$$\boxed{P\left(\frac{B}{A}\right) = \frac{P(B) \times P(A|B)}{P(A)}} \quad \text{Baye's theorem}$$

$$\boxed{P(A|B) = \frac{P(A) \times P(B|A)}{P(B)}} \quad \text{Baye's theorem}$$

$A, B \rightarrow$ events.

$P(A|B) \rightarrow$ Probability of A given B is true.

$P(B|A) \rightarrow$ Probability of B given A is true.

$P(A), P(B) \rightarrow$ Independent probabilities of A and B.

eg DATAS ET	Independent feature	O/P/Dependent feature
Size of House	No. of Rooms	Location
x_1	x_2	x_3

y

$$\boxed{P(y/x_1, x_2, x_3) = \frac{P(y) \times P(x_1, x_2, x_3|y)}{P(x_1, x_2, x_3)}} \quad \text{Baye's theorem}$$

\Downarrow
Baye's theorem

→ CHI SQUARE TEST

The chi square test for Goodness of fit test claims about population proportions.

- It is a non-parametric test that is performed on categorical data. [ordinal and nominal data].

- Q There is a population of male who likes diff. colour of bikes.

Category	Theory	Sample	* Goodness of fit test
Yellow Bike	1/3	22	
Red Bike	1/3	17	
Orange Bike	1/3	59	

↓ observed categorical distribution

Theory categorical distribution

Goodness of fit test

- In a science class of 75 students, 11 are left handed. Does this class fit the theory that 12% of people are left handed.

Category	Obs.	Exp.	$\frac{(O-E)^2}{E}$
Left Hand	11	9	$\frac{(11-9)^2}{9} = 0.44$
Right Hand	64	66	$\frac{(64-66)^2}{66} = 0.03$

- CHI-SQUARE for Goodness of fit →

- Q) In 2010 census of the city, the weight of the individuals in a small city were found to be the following

$< 50 \text{ kg}$	$50-75$	> 75
20%	30%	50%

In 2020, weight of $n=500$ individuals were sampled. Below are the results

< 50	$50-75$	> 75
140	160	200

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

Any 2020	< 50	$50-75$	> 75
Expected	0.2×500	0.3×500	0.5×500

- ① Null hypothesis H_0 : the data meets the expectation.
- ② Alternative " (H_1) : the data doesn't meet expectation.
- ③ $\alpha = 0.05$, CI = 95%.



Mo Tu We Th Fr Sa Su

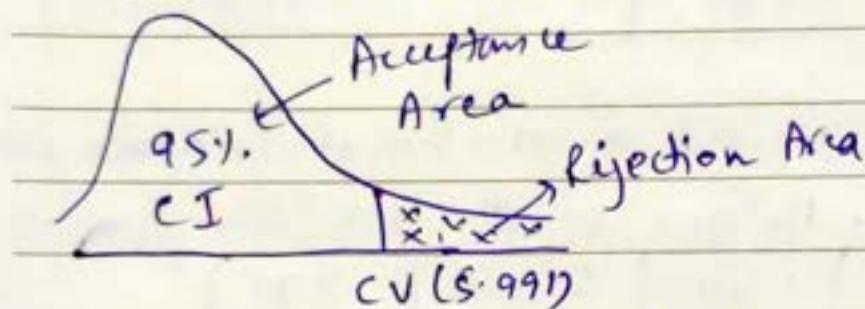
Date / /

③ Degree of freedom

$$df = k - 1 = 3 - 1 = 2$$

$k \rightarrow$ no of categories ($< 50, 50-75, > 75$)

④ Decision Boundary



using chi square table CV will be
5.99

If χ^2 is greater than 5.99, Reject H₀
else we fail to reject the null hypothesis

⑤ Calculate Chi square test statistics, 2020 sample

$$\chi^2 = \frac{\sum (O - E)^2}{E}$$

< 50	$50-75$	> 75
140	160	200

$$= \frac{(140 - 100)^2}{100} + \frac{(160 - 150)^2}{150} + \frac{(200 - 250)^2}{250} \quad \text{2020 expected}$$

$$= \frac{1600}{100} + \frac{100}{150} + \frac{2500}{250}$$

< 50	$50-75$	> 75
100	150	250

Date / /

$$= 16 + 0.66 + 10$$

$$= 26.66$$

$$X^2 = 26.66$$

~~If X^2 Henu~~, $26.66 > 5.99$, Reject H_0 .

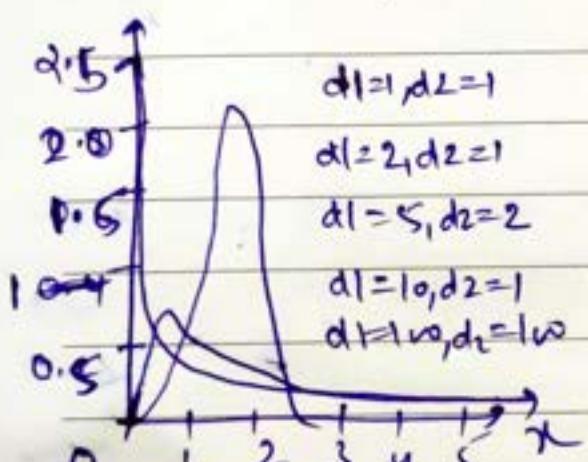
*Conclusion → The weight of 2020 population
are diff. than those expected
in the 2020 population.



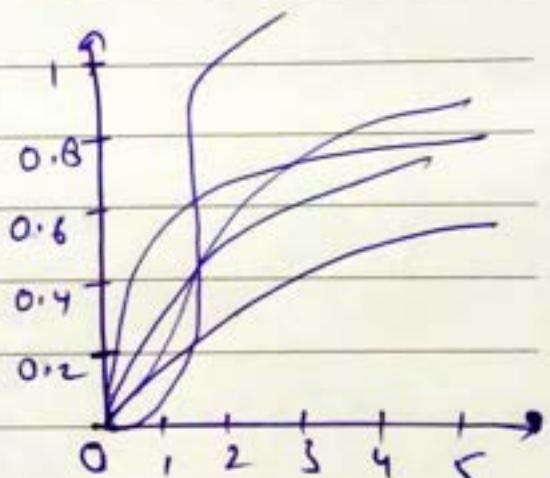
Mo Tu We Th Fr Sa Su

Date / /

→ F distribution :- In prob. or statistics, the F-distribution or F-ratio, also k/a Snedecor's F distribution or the Fisher-Snedecor distribution (after Ronald Fisher and George W. Snedecor) is a continuous prob. distribution that arises frequently as the null distribution of a test statistic, most notably in the analysis of variance (ANOVA) and other f-tests.



Probability density fun"



Cumulative distribution funct.

→ Parameters $\rightarrow d_1, d_2 > 0 \Rightarrow$ degree of freedom
support $\chi^2 \in (0, +\infty)$

$$\text{pdf} = f(x; d_1, d_2) = \frac{\Gamma(d_1 + d_2)}{\Gamma(d_1) \Gamma(d_2)} x^{d_1 - 1} (1 - x)^{d_2 - 1} \cdot \underbrace{\Gamma\left(\frac{d_1}{2}, \frac{d_2}{2}\right)}_{\text{Beta function}}$$

$$\Gamma(m, n) = \frac{(m-1)! (n-1)!}{(m+n-1)!} = \frac{m+n}{mn} \binom{m+n}{m}$$

→ F distribution with d_1 and d_2 degree of freedom in the distribution of

$$X = \frac{S_1/d_1}{S_2/d_2}, S_1 \& S_2 \rightarrow \text{Independent random variables with Chi square distribution}$$

F-test [Variance Ratio Test] $d_1, d_2 \rightarrow$ Respective degree of freedom.



Mo Tu We Th Fr Sa Su

Date / /

F test [Variance Ratio Test] :-

Eg The following data shows the no. of bulbs produced daily for some days by 2 workers A and B.

A	B	→ Can we consider based on the data worker B is more stable and efficient
40	39	
30	38	
38	41	$\alpha = 0.05$
41	33	
38	32	
35	39	
40		
34		

Ans ① Null hypothesis $H_0 \rightarrow \sigma_1^2 = \sigma_2^2$
Alternate $H_1 \rightarrow \sigma_1^2 \neq \sigma_2^2$

② Calculation of Variance

$$S_1^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$



Mo Tu We Th Fr Sa Su

Date / /

Worker A

x_1	\bar{x}_1	$(x_1 - \bar{x}_1)^2$
40	37	9
30	37	49
38	37	1
41	37	16
38	37	1
35	37	4
		$\sum (x_1 - \bar{x}_1)^2$
		80

$$\bar{x}_1 = 37$$

$$= 80$$

Worker B

x_2	\bar{x}_2	$(x_2 - \bar{x}_2)^2$
39	37	4
38	37	1
41	37	16
33	37	16
32	37	25
39	37	4
40	37	9
		$\sum (x_2 - \bar{x}_2)^2$
		84

$$\bar{x}_2 = 37$$

$$S_1^2 = \frac{80}{n-1} = \frac{80}{6-1} = \frac{80}{5} = 13$$

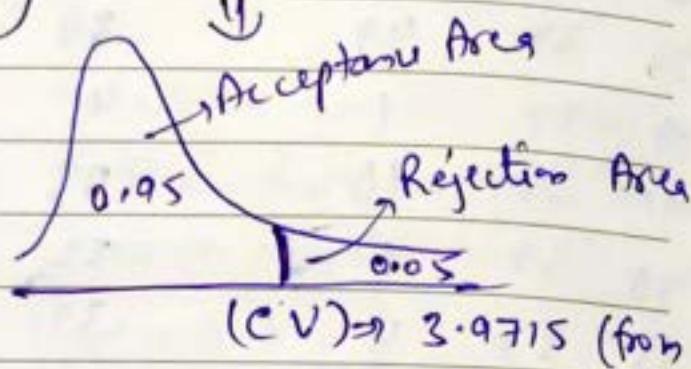
$$S_2^2 = \frac{84}{m-1} = \frac{84}{7-1} = \frac{84}{6} = 12$$

→ Calculation of Variance Ratio Test (F test)

$$F = \frac{S_1^2}{S_2^2} = \frac{13}{12} = \underline{\underline{1.33}}$$

* Decision Rule

$$\begin{aligned} df_1 &= 6 - 1 = 5 \\ df_2 &= 8 - 1 = 7 \\ \alpha &= 0.05 \end{aligned} \quad \left. \begin{array}{l} \text{using F table} \\ \Downarrow \end{array} \right.$$



Decision :-

F test is greater than 3.9715, Reject the null hypothesis

$1.33 < 3.9715$, we fail to Reject the null hypothesis.

Conclusion :-

→ Worker B is not more stable or efficient when compared to worker A.



Mo Tu We Th Fr Sa Su

Date / /

Analysis of Variance (ANOVA)

Defn → ANOVA is a statistical method used to compare the means of 2 or more groups.

ANOVA

① factors (variable)

② Levels

e.g. Medicine (factor)

[Dosage] 5 mg 10 mg 15 mg → levels

→ Mode of Payments (factor)

C Pay Phone IMPS NEFT

Levels.

⇒ Analysis of Variance (ANOVA) →

→ Assumption of ANOVA :-

① Normality of Sampling Distribution of Mean
- The distribution of sample mean is normally distributed.

② Absence of Outliers

- Outline score need to be removed from



Mo Tu We Th Fr Sa Su

Date / /

the datasets

③ Homogeneity of Variance

→ Population variance in different levels of each independent variable are equal

$$[\sigma_1^2 = \sigma_2^2 = \sigma_3^2]$$

④ Samples are independent and random

→ Types of ANOVA :-

(3 types)

① One way ANOVA - One factor with at least 2 levels, these levels are independent

e.g Doctor wants to test a new medication to decrease headache. They split the ~~population~~ participants in 3 conditions

[10 mg, 20 mg, 30 mg]. Doctors ask the participants to rate the headache [1-10]

Meditation → factor

10 mg	20 mg	30 mg
5	7	2
3	4	6
-	-	-
-	-	-

② Repeated Measures ANOVA → One factor with at least 2 levels, levels are dependent.

eg Running → factor

Day 1	Day 2	Day 3 → levels
8 km	5 km	4 km
7 km	4 km	9 km

③ Factorial ANOVA - Two or more factors, each of which with at least 2 levels, their levels can be independent or dependent.

Running → factor

Level	Day 1	Day 2	Day 3 → levels
Male	8	5	4
Gender	9	4	3
Female	2	4	6
Factor	7	8	3



Mo Tu We Th Fr Sa Su

Date / /

→ Hypothesis Testing in ANOVA (Partitioning of Variance in the ANOVA)

Null hypothesis $H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_k$.

Alternative $H_1:$ At least one of the sample mean is not equal.

$$\boxed{\mu_1 \neq \mu_2 \neq \mu_3 = \dots \neq \mu_k}$$

not applicable

Test Statistics

$$F \text{ test} = \frac{\text{Variance b/w sample}}{\text{Variance within sample}}$$

Variance b/w Sample			$H_0 \rightarrow \bar{X}_1 = \bar{X}_2 = \bar{X}_3$
$\sum x^2$	x_1	x_2	x_3
Variance within sample	1	6	5
	2	7	6
	4	3	3
	5	2	2
	3	1	4
	$\bar{x}_1 = 15/5 = 3$	$\bar{x}_2 = 19/5$	$\bar{x}_3 = 4$

One Way ANOVA

one factor with atleast 2 levels, levels are independent.

- d) Doctor wants to test a new medicine which reduces headache. They split the participants into 3 cond'n [15mg, 30mg, 45mg]. Later on the doctor ask the patient to rate the headache b/w [1 - 10]. Are there any differences b/w the 3 conditions using $\alpha = 0.05$?

Any	15 mg	30 mg	45 mg
9	7	4	
8	6	3	
7	6	2	
8	7	3	
8	8	4	
9	7	3	
8	6	2	

- ① Define Null and Alternate Hypotheses :-
- $H_0: \mu_{15} = \mu_{30} = \mu_{45}$
- $H_1: \text{not all } \mu \text{ are equal.}$

② Significance $\alpha = 0.05$, CI = 0.95

③ Calculate degree of freedom.

$$N = 21, a = 3, n = 7$$

N → total no. of participants ($9 \times 3 = 27$)

n → participants in a single level (e.g. 15mg)

a → factors (20mg, 15mg, 45mg) = 3.

$$df_{between} = a - 1 = 3 - 1 = 2 \quad \gamma \quad df_1, df_2 \\ (2, 18)$$

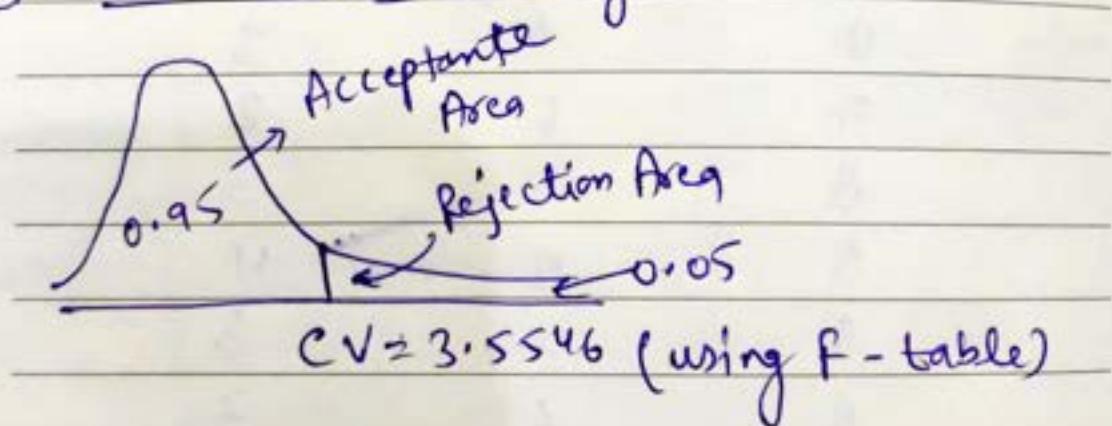
$$df_{within} = N - a = 21 - 3 = 18 \quad \cup$$

$$df_{total} = N - 1 = 20 \quad f\text{-table}$$

$$\alpha = 0.05$$

Critical Value

④ Decision Boundary



Decision Rule

If F is greater than 3.5546, Reject the Null hypothesis

⑤ Calculate F test statistic

	SS	df	MS	F
Between	98.67	2	49.34	
Within	10.29	10	0.54	
Total	108.96	20		

i) $SS_{\text{between}} = \frac{\sum (\sum a_i)^2}{n} - \frac{T^2}{N}$

$$15 \text{mg} := 9 + 8 + 7 + 8 + 8 + 9 + 8 = 57$$

$$30 \text{mg} := 7 + 6 + 6 + 7 + 8 + 7 + 6 = 47$$

$$45 \text{mg} := 4 + 3 + 2 + 3 + 4 + 3 + 2 = 21$$

$$SS_{\text{b/w}} = \frac{57^2 + 47^2 + 21^2}{7} - \frac{[57^2 + 47^2 + 21^2]}{21}$$

$$= 98.67$$

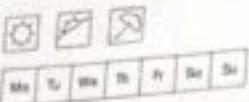
ii) $SS_{\text{within}} = \sum y^2 - \frac{\sum (\sum a_i)^2}{n}$

$$\sum y^2 = 9^2 + 8^2 + 7^2 + 8^2 + \dots$$

$$= 853$$

$$SS_{\text{within}} = 853 - \left[\frac{57^2 + 47^2 + 21^2}{7} \right]$$

$$= 10.29$$



Date

$$f \text{ test} = \frac{MS_{\text{between}}}{MS_{\text{within}}} = \frac{98.67/2}{10.29/2} \approx \frac{49.34}{0.54}$$

$| F = 86.56$

If f is greater than 3.5546 , Reject
the null hypothesis

$86.56 > 3.5546$, Reject the H_0