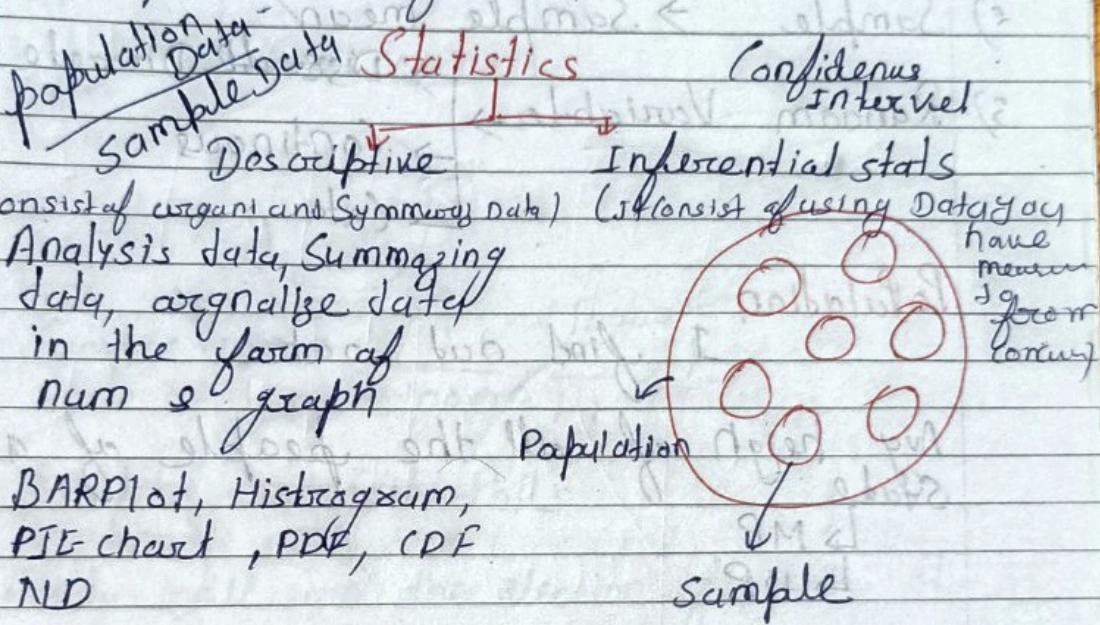


11 Statistics

Date []

Statistics is the discipline that concerns the collections, organization, analysis, interpretation, and presentation of data.



- 3) Measure of Center
Tendency
mean, median, mode

Party 1 → :
Party 2 → ..
Party 3 → x.

- 4) Measure of Variance
SD, Variance

take Sample
↳ Inferences and Conclusion

- 1) Z Test
2) T Test
3) Chi square

Population

A circled checkmark symbol, indicating a correct answer or a mark for grading.

Date

Population vs Sample

- 1) Population \rightarrow Population mean \rightarrow made
 - 2) Sample \rightarrow Sample mean
 - 3) Random Variable \rightarrow
 - \rightarrow Discrete
 - \rightarrow Continuous
 - \rightarrow Categorical

Population

I find out out

Avg height of all the people of a state

LMP
LUP

	$x x x x x x$	1 million
MP	$x x x x x x$	
Population	$x x x x x x$	
	$x x x x x x$	

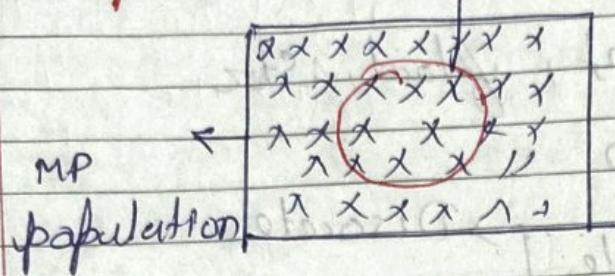
μ = mean = Avg height of all the people.

$$\text{population mean} = \frac{1}{1M} \sum_{i=1}^{1M} \text{Op}$$

(13)

Date: ~~10/10/19~~

Sample

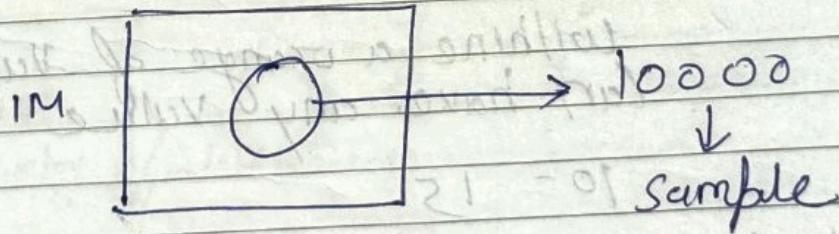
1000
Sample

add
xx
xx

$$\text{Sample mean} = \frac{1}{10000} \sum_{i=1}^{10000} \text{hrs}$$

eg Sample \rightarrow Exit Polls

- \rightarrow Party will win the election from the state
- \rightarrow News channel Exit Poll \rightarrow Sample Data



option which
Party is win

Population mean > Sample mean

① Simple Random Sampling



Random Variable:

Variable = int, float, str

$$x = 10$$

Random Variable

→ Discrete

→ Continuous

Discrete:-

whole no, not be floating no.

Eg.

① no of bank account

2, 3, 4, 5,

× 2.5

② Population of state

1 million

× 1.5 person

Continuous:-

Within a range of value we can have any value

10 - 15

10.1, 10.2, 10.3, 14.9, 15, 14.9

↳ ~~discrete~~

int num is also

Eg: 5.11

5.8

Measure of Central Tendency

- 1) Mean
- 2) Median
- 3) Mode

Center Tendency refers to
the measure used to determine
the "center" of the distribution
of data.

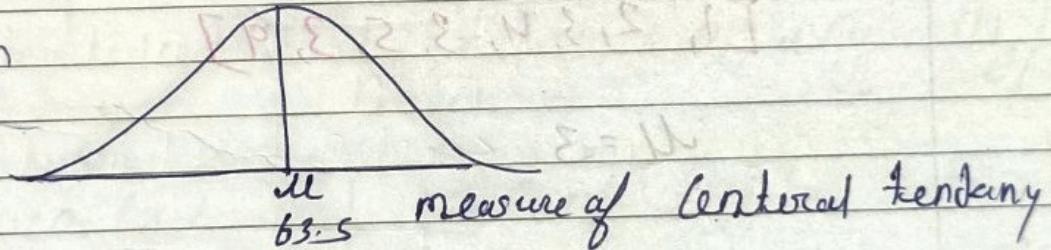
Mean - {Average}

Sample of Height - {168, 170, 150, 160, 182, 140, 175}

S, 27, 3, 4, 5)

$$\bar{M} = \frac{\sum_{i=1}^n x_i}{n} = \frac{\text{Sum}}{7} = 63.5$$

M_2 Population
 n Sample



Median:

Disadvantage of mean

$$\text{eg } [1, 2, 3, 4, 5] = \frac{15}{5} = 3$$

$$[1, 2, 3, 4, 5, 50] = \frac{65}{5} = 13$$

$$\bar{M} = 3$$

ON $\bar{M} = 3$

$$\bar{M} = 13$$

new mean = 13

Very Different

(X)

Medium:

the no of shaded circle

 $[1, 2, \boxed{3, 4}, 5, 50]$

1) find the shaded no

2) find the center element

$$\text{even no} = \frac{3+4}{2} = 3.5$$

 $[1, 2, \boxed{3}, 4, 5]$

odd no = 3

Mode:-

maximum we no is selection

 $[1, 2, 3, 4, 3, 5, 3, 9]$

$M = 3$

eg

Any

23

24

27

32

35

21

to fill missing value

mean

medium

mode

(17)

Measure of Dispersion (Covariance)-

Date _____

one of the very imp. topic Data
processing

I have two Random Variable
Quantify a relation

Size	price
1200 sqft	1000k
1500 sqft	3000k
2000 sqft	5000k

Size \leftrightarrow price

find the relation b/w Two
Variable.

ST	PT
SD	PD

Covariance-

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

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

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

$$\boxed{\text{Cov}(x, x) = \text{Var}(x)}$$

$$x \uparrow y \uparrow = +ve$$

$$x \downarrow y \downarrow = +ve$$

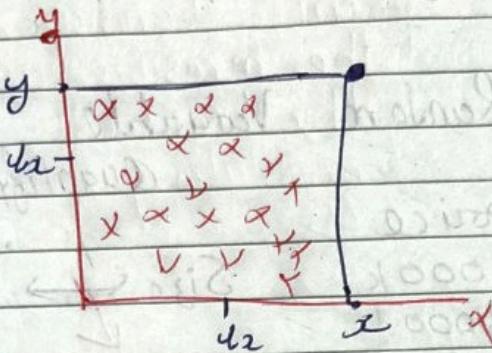
$$x \uparrow y \downarrow = -ve$$

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

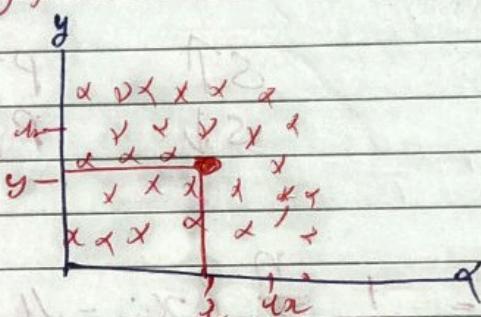
(18)

Date _____

$$\text{Cov}(x,y) = (+ve) \times (+ve) = +ve$$



$$\text{Cov}(x,y) = (-ve) \times (-ve) = +ve$$



$$\text{Cov}(xy) = (-ve) \times (ve) = -ve$$

$$\text{Cov}(yx) = (ve) \times (-ve) = -ve$$

$x \uparrow y \uparrow$ = +ve How much +ve

$x \uparrow y \downarrow$ = -ve How much -ve

Standard deviation = $\sqrt{\text{Var}}$

19

Date

Pearson Correlation

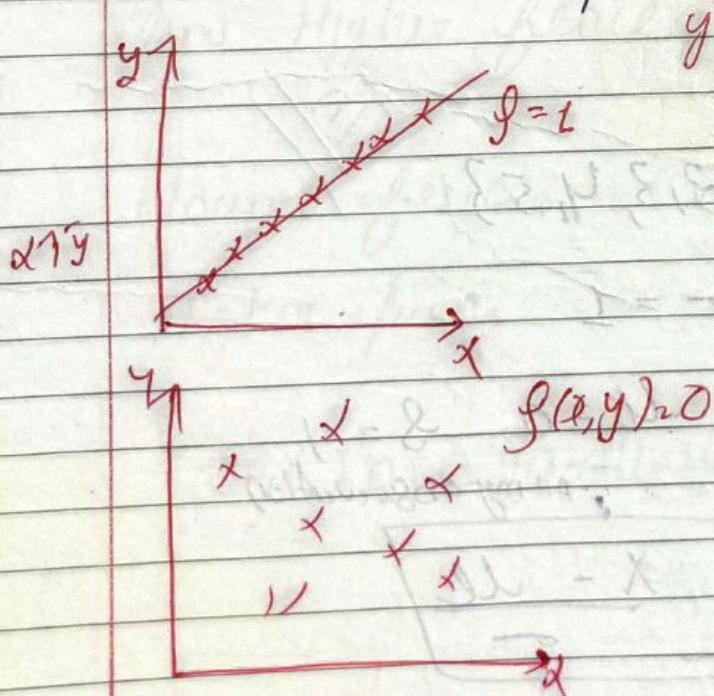
$$\text{Co-Variance} = \text{Cov}(x, y) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$x \uparrow y \uparrow = \text{pos}$
 $x \downarrow y \downarrow = \text{pos}$
 $x \uparrow y \downarrow = \text{neg}$
} find the direction of relationship

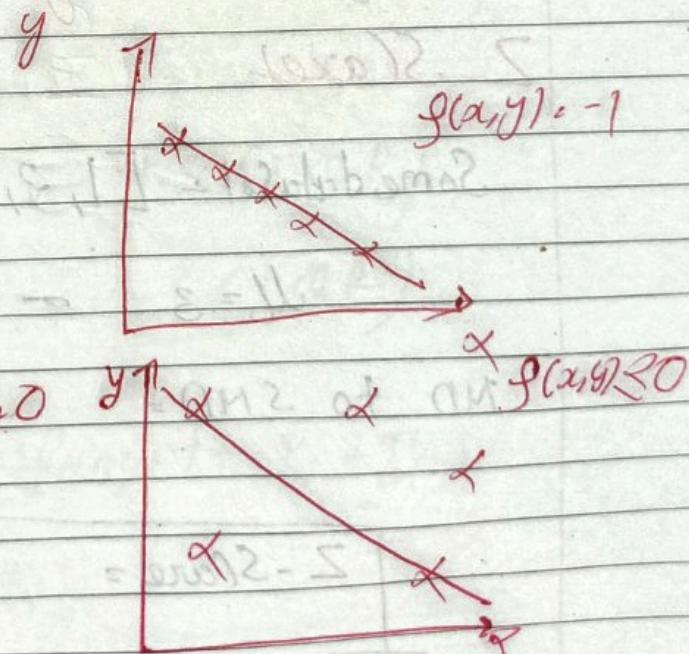
$$\text{Pearson Cov} = \rho(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y}$$

$$\rho(x, y) = -1 \leq \rho \leq 1$$

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



$$\rho(x, y) = 0$$



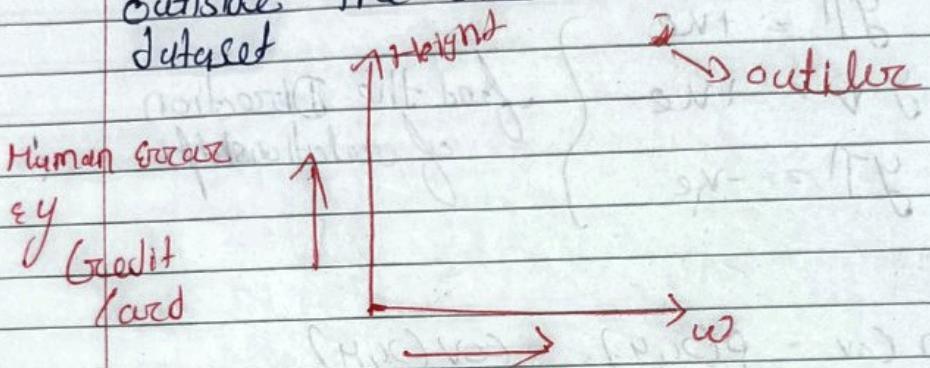
$$\rho(x, y) < 0$$

(20)

Date

Outlier

An outlier is a data point in a dataset that is distant from all other observations. A data point that lies outside the overall distribution of the dataset.



Identify outlier

- 1) Z-Score
- 2) Interquartile Range

Z-Score

Some dataset = $\{1, 2, 3, 4, 5\}$

$$\mu = 3 \quad \sigma = 1$$

$$ND \text{ to } SND = \mu = 0, \sigma = 1$$

all my observation

$$\boxed{Z\text{-Score} = \frac{x - \mu}{\sigma}}$$

(21)

Date

II Interquartile Range

Find Number Summary

- 1) Minimum
- 2) First Quartile (Q_1)
- 3) Median
- 4) Third Quartile (Q_3)
- 5) Max

Removing the outlier

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

find Lower fence

find Higher fence

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

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

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

$$IQR = Q_3 - Q_1$$

$$IQR = Q_3 - Q_1$$

(22)

Date

$$q_1 = 28 \times 1 - \frac{25}{100} \times 19 + 1$$

$$\frac{25}{100} \times 20 = 5 \text{ index}$$

$$q_1 = 5$$

$$q_3 = \frac{75}{100} \times 20 = 15 \text{ index}$$

$$q_3 = 7$$

$$IQR = q_3 - q_1 = 7 - 5 = 4$$

$$\text{Lower fence} = q_1 - 1.5(IQR)$$

$$5 - 1.5 \times 4$$

$$3 - 6 = -3$$

$$\text{Higher fence} = q_3 + 1.5(IQR)$$

$$7 + 1.5 \times 4$$

$$7 + 6 = 13$$

(23)

Date

[lower fence - Higher fence]

[-3 - 13]

what outlier = 27

Min = 1

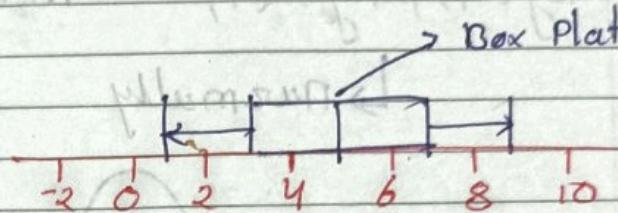
Box Plot

$Q_1 = 3$

median = 5

$Q_3 = 7$

$MgX = 9$



TYPE of Distributions

Gaussian Distribution / Normal Distribution

Eg >

Harm

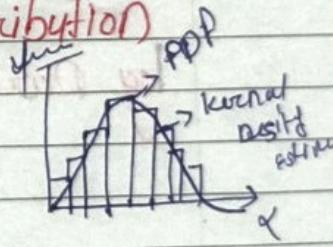
weigh

size

out

$$\text{DC} \sim \text{GSD}(\mu, \sigma)$$

mean SD

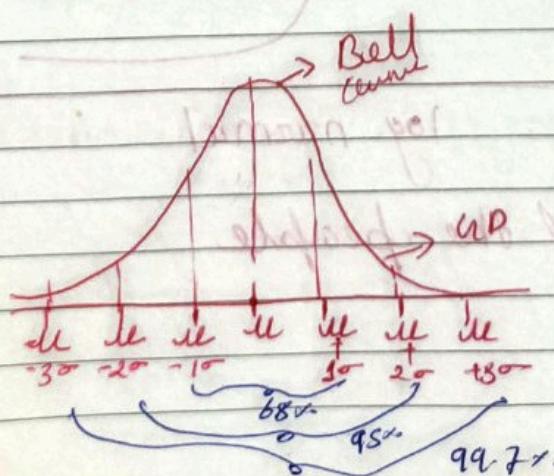


$$\mu = \frac{1}{n} \sum_{i=1}^n x_i$$

T-

gym

spare



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

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

- (1) Symmetric
- (2) 68 - 95 - 99.7% Empirical Rules

24 Central Limit Theorem \Rightarrow Sample distribution
NPR 30 Date []

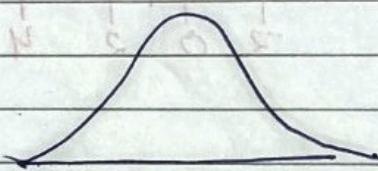
log Normal Distribution

x is log normal Distribution if $\ln(x)$ is normally distributed

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

$$\ln(x_1), \ln(x_2), \ln(x_3)$$

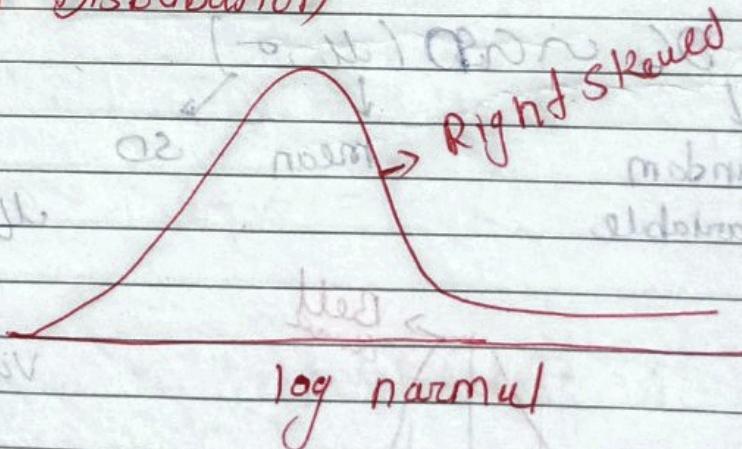
normally



x is log normal Distribution if

$$\ln(x) \sim N(\mu, \sigma^2)$$

log Normal Distribution

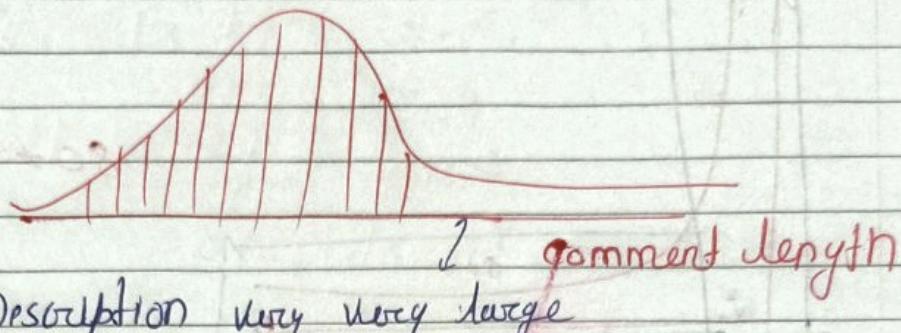


e.g. - Income of the people

2

Date

Eg: product comments (Amazon products)



larger Description very very large

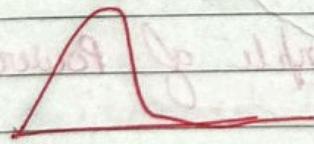
Eg of ND

1) Height

2) weight

Eg of log ND

①

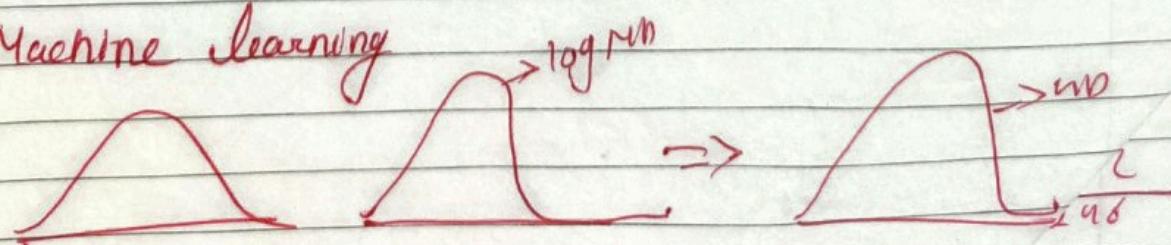


①

wealth Distribution

②

Machine learning

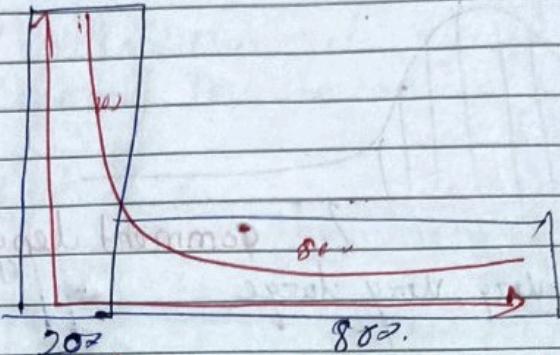


Data Transformation

2

Date

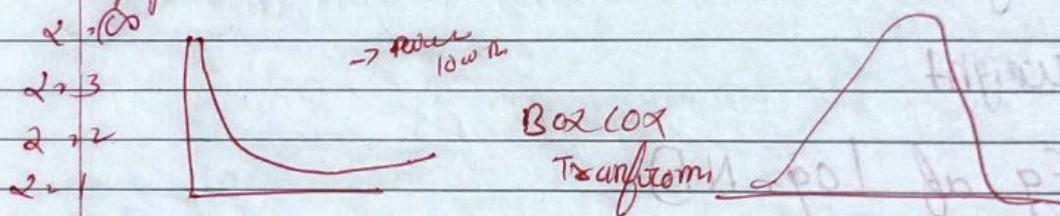
Power Law Distribution



80+20% rule

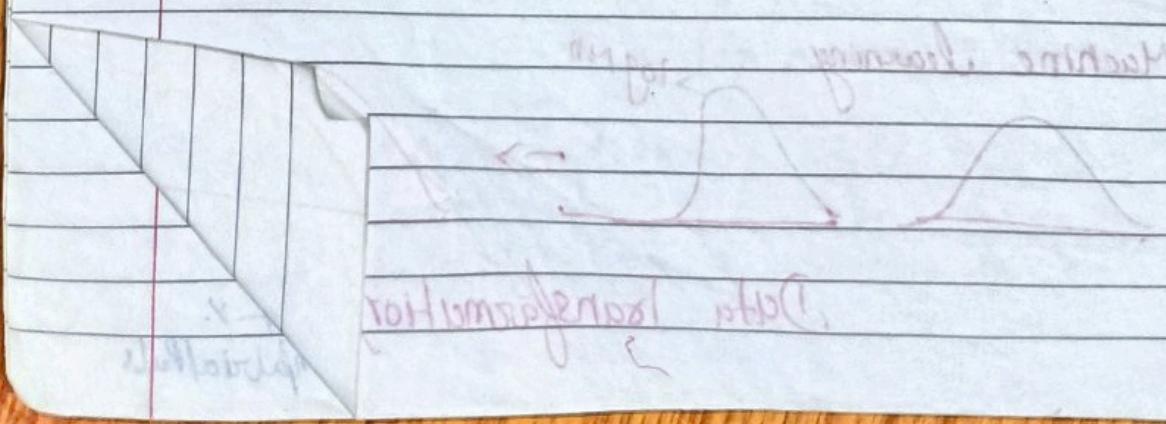
Eg 20% of Team is responsible for winning 80% of math

→ 80% of world's distributed 20% of the total
population



Pareto Distribution

- ii This is a example of Power
- iii Distribution
- iv it is not normal Distribution



(27)

Date _____

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

x	\bar{x}	$(x - \bar{x})^2$	$(x - \bar{x})^2$
1	2.83	-1.83	3.34
2	2.83	-0.83	0.6889
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
			<u>10.84</u>
	<u>2.83</u>		

$$\bar{x} = 2.83$$

$$\text{Var} = s^2 = \frac{10.84}{5} = 2.168$$

Sample Var

$$SD = \sqrt{\text{Var}} = \sqrt{2.168}$$

$$\text{Standard Deviation} = 1.472$$

$$\text{Dataset} \rightarrow \bar{x} = 2.83 \quad SD = 1.472$$

$$\begin{array}{r} 2.83 \\ 1.472 \\ \hline 1.358 \end{array}$$

