

18/10/21.

## UNIT-3.

Central tendency:- (measures).

- ① Arithmetic mean (value at the middle),
- ② Median (middle value),
- ③ Mode (
- ④ Geometric mean
- ⑤ Harmonic mean

Averages and are statistical constants which enable us to comprehend in a single effort the ~~simp~~ significance of the entire data (UR) whole. (this gives us an idea about the concentration of values in the central part of the data)

Requisites for an ideal measure :-

- (i) It has to be define rigidly. (To evaluate it need a exact formula to calculate the mean)
- (ii) It should be readily comprehensible and easy to calculate.
- (iii) It is based on all the observations.

Classification of data:

Example:-

- (i) Individual observations  $\rightarrow 30, 35, 37, 30, 43, 46.$
- (ii) Discrete observations  $\rightarrow \begin{cases} x: & 30 & 35 & 37 & 43 & 46 \\ f: & 2 & 1 & 1 & 1 & 1 \end{cases}$
- (iii) Continuous observations  $\rightarrow 30-35 - 2$  (including 35)

$$30-35 - 2$$

$$35-40 - 2$$

$$40-45 - 1$$

$$45-50 - 1$$

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

$$35-40 - 1$$

$$40-45 - 1$$

$$45-50 - \frac{1}{6}$$

## Definition of Average:

Arithmetic

(i) Individual observations:  $\bar{x}$  ( $\bar{x}$ -mean)

$$\bar{x} = \frac{\sum x_i}{n}$$

Add all the values divided by no. of observation that would be our average.

(ii) Discrete Observations:

$$\bar{x} = \frac{\sum f x_i}{\sum f}$$

frequency. Multiplication of frequent divided by the frequencies.

(iii) continuous observations: (frequency distribution)

$$\bar{x} = A + \frac{\sum f d}{\sum f} \times i$$

$\rightarrow$  length of the class interval.

Assumed mean (middle value)  $d = \frac{m - A}{i}$   $i \rightarrow$  length of the CI

19/10/21

## Properties of Arithmetic Mean (AM):

① Algebraic sum of the deviations of a set of values from their AM is zero!

$$\text{T.P.: } \sum_i f_i (\underbrace{x_i - \bar{x}}_{\text{deviations from AM}}) = 0 \quad \{ \text{Discrete case} \}$$

[If  $x_i$  and  $f_i$  is provided we can find the AM].

(if the  $f_i$  is '1' then it is individual observations)

Proof:

LHS:

Consider  $\sum_i f_i (x_i - \bar{x})$

$$= \sum_i f_i x_i - \sum_i f_i \bar{x}$$

$$\boxed{\bar{x} = \frac{\sum f_i x_i}{\sum f_i}}$$

$$= \bar{x} \sum_i f_i - \bar{x} \sum_i f_i = 0$$

② The sum of the squares of the deviations of the set of values is the minimum when taken about mean.

Proof:-

$$x: x_1 \ x_2 \dots x_n$$

$$f: f_1 \ f_2 \dots f_n$$

$$Z = \sum_{i=1}^n f_i (x_i - A)^2$$

T.P:-

The minimum value is obtained when

$$A = \bar{x}$$

Dif. w.r.t. A;

$$\frac{dZ}{dA} = \sum_{i=1}^n 2 f_i (x_i - A) (-1)$$

$$\frac{d^2 Z}{dA^2} = \sum_{i=1}^n 2 f_i (-1) (-1)$$

$$\Rightarrow \sum 2 f_i > 0$$

$\Rightarrow$  Minimum value is obtained when

$$-2 \sum f_i (x_i - A) = 0 //$$

(-2 cannot be zero)

Dispersion

$$\Rightarrow \sum f_i x_i - A \sum f_i = 0$$

S.D.

$$A = \frac{\sum f_i x_i}{\sum f_i} = \bar{x}$$

all these are  
about mean and  
will always be  
minimum.

③ Combined AM

If  $\bar{x}_i$  ( $i=1, 2, \dots, k$ ) are the means of  $k$

component series of size  $n_i$  ( $i=1, 2, \dots, k$ )

respectively, then the mean  $\bar{x}$  of the

composites series obtained by combining the component series is given by ;

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2 + \dots + n_k \bar{x}_k}{n_1 + n_2 + \dots + n_k} \quad \left. \begin{array}{l} \text{combined} \\ \text{mean} \end{array} \right\}$$

$\frac{\bar{x}_1}{n_1}$  Amp: 1       $\frac{\bar{x}_2}{n_2}$  Amp: 2       $\frac{\bar{x}_k}{n_k}$  Amp K -

Proof:

$$\bar{x}_1 = \frac{(x_{11} + x_{12} + \dots + x_{1n_1})}{n_1}$$

$$\bar{x}_{1k} = \frac{(x_{k1} + x_{k2} + \dots + x_{kn_k})}{n_k}$$

$$\bar{x} = \frac{(x_{11} + x_{12} + x_{13} + \dots + x_{1n_1}) + \dots + (x_{k1} + x_{k2} + \dots + x_{kn_k})}{n_1 + n_2 + \dots + n_k}$$

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2 + \dots + n_k \bar{x}_k}{n_1 + n_2 + \dots + n_k}$$

(n<sub>1</sub> + n<sub>2</sub> + ... + n<sub>k</sub>) no. of observations

- Find the mean of 32, 45, 28, 13, 0, 16, 32, 20, 27.

Soln:

Sum of observations

No. of observations

$$\bar{x} = \frac{\sum x_i}{n}$$

Formula.

$$0, 13, 16, 32, 28, 32, 20$$

$$\frac{20+21}{2} = 20.5$$

$$\bar{x} = \frac{32+45+28+13+0+16+32+27}{8}$$

$$= \frac{193}{8}$$

$$\bar{x} = 24.125.$$

2. The algebraic sum of the deviations of 20 observations measure from 30 is 2. Find the mean?

Soln:  $\sum_{i=1}^{20} (x_i - A) = 2$

$$\sum_{i=1}^{20} (x_i - 30) = 2.$$

$$\sum_{i=1}^{20} x_i - 20 \times 30 = 2$$

$$\sum_{i=1}^{20} x_i = 602.$$

$$\bar{x} = \frac{\sum_{i=1}^{20} x_i}{20} \rightarrow \boxed{\bar{x} = \frac{\sum x_i}{n}}$$

$$= 30.11$$

3. calculate the mean for the following.

$x: 0 \quad 10 \quad 20 \quad 30 \quad 40 \quad 50 \quad 60$

$f: 8 \quad 10 \quad 11 \quad 16 \quad 20 \quad 25 \quad 15$ .

x	f	fx
0	8	0
10	10	100
20	11	220
30	16	480
40	20	800
50	25	1250
60	1	60

$$\bar{x} = \frac{\sum f x}{\sum f}$$

Sum of the frequencies = 105.

$$\bar{x} = \frac{3750}{105}$$

$$\bar{x} = 34.35 \text{. } \cancel{74} \quad (\text{should be in 3 decimal places})$$

A. calculate the mean for the following frequency distribution:

CI -	0-8	8-16	16-24	24-32	32-40	40-48
f -	4	12	20	28	36	44

Soln:

$$\bar{x} = A + \frac{\sum f d}{\sum f} \times i \rightarrow \begin{array}{l} \text{Assumed mean (middle value)} \\ \text{length of class CI} \\ \text{midvalue of CI} \end{array}$$

$$\text{where } d = \frac{m - A}{i}$$

$$\text{midpt} = \frac{(I) \rightarrow \text{Add Cs}}{2}$$

CI	midpt.	f	$d = \frac{m - A}{i}$	fd	$\sum f m$	
0-8	4	4	$d = \frac{4 - 28}{8} = -3$	$4 \times -3 = -12$	16	$\sum f = 144$
8-16	12	12	-2	$12 \times -2 = -24$	144	$\sum f d = 68$
16-24	20	20	-1	-20	400	$\sum f m = 4576$
24-32	28	28	0	0	784	
32-40	36	36	1	36	1296	
40-48	44	44	2	88	1936	

$$\bar{x} = A + \frac{\sum fd}{\sum f} \times i$$

144  
27

$$= 28 + \frac{68}{144} \times 8$$

$$= \frac{4032 + 68 \times 8}{144} = \frac{4576}{144}$$

(26 - 13)

$$= 31.777 \dots / 31.778 \quad -0.5 \\ -1.5 \quad 1$$

(or)

$$\bar{x} = \frac{\sum fm}{\sum f}$$

$$= \frac{4576}{144} = 31.77 \dots / 31.778$$

$$\begin{array}{r} 1.8 - 8 \\ \hline 0.8 \end{array} \quad \frac{7}{2} \quad 3.0 \quad \frac{0.17}{2} \\ \hline \frac{7}{2} \end{array}$$

$$= 3.8$$

5. Find the mean of the following:

$$(I) : \underline{0-7.5} \quad \underline{8-14} \quad \underline{15-22} \quad \underline{23-30} \quad \underline{31-38} \quad \underline{39-46}$$

$$f: \quad \underline{8} \quad \underline{7} \quad 16 \quad 24 \quad 15 \quad 4$$

Soln:

Since it unequal values so that  
it is not calculated

CI	f	midpt.	d = m - A	fd	$\frac{\sum fd}{\sum f}$
-0.5 - 7.5	8	3.5	-23	-184	$7.5 + 14.5$
7.5 - 14.5	7	11	-15.5	-108.5	$\frac{22}{2}$
14.5 - 22.5	16	18.5	-8	-128	$7.5 + 14.5$
22.5 - 30.5	24	26.5	0	0	$\frac{26.5 + 14.5}{2}$
30.5 - 38.5	15	34.5	+8	120	$26.5 + 14.5$
38.5 - 46.5	7	42.5	+16	112	$\frac{26.5 + 14.5}{2}$
	<u>77</u>			<u>-188.5</u>	

$$\sum f = 77$$

$$\sum fd = -188.5$$

$$\bar{x} = A + \frac{\sum fd}{\sum f}$$

$$\bar{x} = 26.5 + \frac{(-188.5)}{77}$$

$$= 24.05$$

20/10/21

Recall:

→ AM: It's a measure to represent the entire data.

→ Individual observation:  $\bar{x} = \frac{\sum x_i}{n}$   
n → no. of observations

→ Discrete observation:  $\bar{x} = \frac{\sum f_i x_i}{\sum f_i}$  frequency

→ Assumed mean

→ continuous data:  $\bar{x} = A + \frac{\sum f_i d_i}{\sum f_i} x_i$

→ mid-pt of the CI

$$d = \frac{m - A}{i}$$

i → length of the CI.

→ Combined AM:-

$$\left( \frac{n_1}{\bar{x}_1} \right) \cdot \left( \frac{n_2}{\bar{x}_2} \right) \cdots \left( \frac{n_k}{\bar{x}_k} \right)$$

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2 + \cdots + n_k \bar{x}_k}{n_1 + n_2 + \cdots + n_k}$$

- (1) The average salary of male employees in a firm was Rs. 5200 and that of females was Rs. 4200. The mean salary of all employees was Rs. 5000. Find the percentage of male and female employees.

Soln:

$$\bar{x}_1 = 5200$$

Average salary of male employee

male	female
5200	4200

Combined AM  $\Rightarrow 5000$ .

$\bar{x}_2 = 4200$  (Average salary for female employee)

$$\bar{x} = 5000.$$

$n_1 \rightarrow$  is no. of male employees

$n_2 \rightarrow$  no. of female employees.

$$\left( \bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2} \right) \rightarrow \text{Formula}$$

$$5000 = \frac{n_1 (5200) + n_2 (4200)}{n_1 + n_2}$$

$\therefore$  of male  
per

$$5000n_1 + 5000n_2 = 5200n_1 + 4200n_2.$$

$$5000n_2 - 4200n_2 = 5200n_1 - 5000n_1$$

$$800n_2 = 200n_1$$

$$\frac{n_1}{n_2} = \frac{800}{200} = \frac{4}{1}$$

$$\% \text{ of male employees} = \frac{4}{5} \times 100 = 80$$

$$\% \text{ of female employees} = \frac{1}{5} \times 100 = 20.$$

80-1 are males and 20-1 are females.

### Disadvantages of AM:

1. It is affected by the extreme values.
2. It can't be determined by inspection nor it can be located graphically.
3. AM can't be calculated if any observation is missing.
4. In extremely asymmetrical distribution AM is not a suitable measure of location.

Median:

↳ (The middle value of the data.) NO. of observations

Individual observations: - size of  $\left(\frac{N+1}{2}\right)^{\text{th}}$  item.

Discrete observations: - size of  $\left(\frac{N+1}{2}\right)^{\text{th}}$  item.

Continuous observations: Median =  $L + \frac{\frac{N}{2} - cf}{f} \times i$

lower limit of the  
median class

f

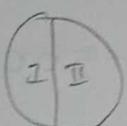
frequency  
of median class

preceding  
no. of median  
class

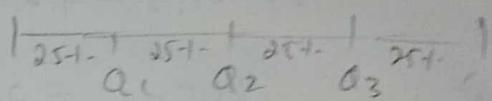
length of the  
class

Median is nothing but cutting the two parts

equally.



quartiles



$$Q_1 = L + \frac{\frac{N}{4} - cf}{f} \times i \quad (\text{1st quartile})$$



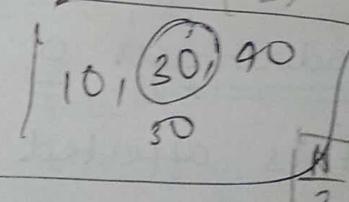
$$Q_3 = L + \frac{\frac{3N}{4} - cf}{f} \times i \quad (\text{3rd quartile})$$

What about for Indi. & Discrete observations we have written  $\left(\frac{N+1}{2}\right)^{\text{th}}$  and for continuous observation there is  $\left(\frac{N}{2}\right) \rightarrow ?$

Eg:

(i) Find the median of 20?

There is no median.



(ii) Find the median of 10, 20, 30, 40.

In this median is 25.

$$\frac{20+30}{2}$$

$$\frac{50}{2}$$

Advantage:

\* It is unaffected by the extreme values.  
(outliers)

1. Find the median for the data: (Individual observations)  
 $40, 25, 60, 35, 44, 50$

Soln:

\* Arrange in ascending / descending order

$25, 35, \boxed{40}, \boxed{44}, 50, 60.$

No. of observation ( $N$ ) = 6.

Median :- size of  $(\frac{N+1}{2})^{\text{th}}$  item. = size of  $3.5^{\text{th}}$  item

= Avg. of  $3^{\text{rd}}$  &  $4^{\text{th}}$  position

$$= \frac{40+44}{2}$$

$$= 42.$$

The median is  $42\frac{1}{2}$ .

~~20/44  
44  
2~~

2. Find the median of  $30, 60, 20, 15, 70.$

Soln:

\* Arrange in ascending order

$15, 20, \boxed{30}, 60, 70.$

$$N = 5.$$

Median := size of  $(\frac{N+1}{2})^{\text{th}}$  item = size of  $3^{\text{rd}}$  item

= Avg

$$\text{Median} = 30\frac{1}{2}$$

3. Obtain median for the following (Discrete observation)

$$x: 10 \quad 20 \quad 30 \quad 40 \quad 50 \quad 60 \quad 70 \quad 80$$

$$f: 3 \quad 15 \quad 20 \quad 6 \quad 10 \quad 8 \quad 7 \quad 11.$$

Soln:

\* Form the table as follows;

$n$     $f$     $C$

10	3	3
20	15	18
30	20	38
40	6	44
50	10	54
60	8	62
70	7	69
80	11	80
	86	

$$\text{Median} = \text{size } \left( \frac{N+1}{2} \right)^{\text{th}} \text{ item}$$

$$= \text{size } \left( \frac{80+1}{2} \right)^{\text{th}} \text{ item}$$

$$N = \sum f = 80.$$

$$= \text{size } (40.5)^{\text{th}} \text{ item.}$$

In cf column we should see the no's that is ( $\geq$ ) to the size of the item

$$\begin{cases} 10 & (10.5) \\ \text{In this soln it is } 44 \end{cases}$$

$$\begin{matrix} 1 \\ 50 & (10.5) & 21 \\ \geq & \geq & \geq \end{matrix}$$

$$\text{Median} = 40 //$$

Suppose if we want to calculate

$$Q_1 = \text{size of } \left( \frac{N+1}{4} \right)^{\text{th}} \text{ item}$$

$$Q_1 = \frac{N+1}{4}$$

$$= \text{size of } \left( \frac{8}{4} \right)^{\text{th}} \text{ item}$$

$$= \frac{18}{4}$$

$$= 20.25 //$$

$$Q_1 \leq 30, \text{ Median} = 40.$$

4. Find the median. (continuous observation).

CI: 10 - 20 20 - 30 30 - 40 40 - 50 50 - 60

$f:$  20 90 150 100 40.

Soln:

Form the table::

P.T.O.

(I)	f	cf
10-20	20	20
20-30	90	110 cf
30-40	150	260
40-50	100	360
50-60	70	430
	130	

$$\frac{N}{2} = \frac{430}{2} = 215.$$

$\frac{430}{2} = 215$

$\geq$

$$\begin{aligned} \text{Median} &= L + \frac{\frac{N}{2} - f}{f} \times i \\ &= 30 + \frac{215 - 110}{150} \times 10 \\ &= 30 + \frac{105}{15} \Rightarrow 37 \end{aligned}$$

5. Given median is  $46$ . Determine the missing frequencies

(I): 10-20 20-30 30-40 40-50 50-60 60-70 70-80

f: 12 30 ? 65 ? 25 18

sum of all frequencies is 229.

Soln:

Form a table.

(I)	f	cf
10-20	12	12
20-30	30	42
30-40	x	$(42+x) cf$
40-50	$65 f$	$107+x$
50-60	y	$107+x+y$
60-70	25	$132+x+y$
70-80	18	$150+x+y$
	<u>229</u>	

frequency column  
is called cumulative frequency  
and value of the  
Given:

median =  $46$ .

Median class = 40-50.

$$150 + x + y = 229$$

$$x + y = 229 - 150$$

$$\boxed{x + y = 79}.$$

$$L = 40, \frac{N}{2} = \frac{229}{2} = \boxed{114.5.}$$

$$cf = 42 + x + f = 65$$

$$i = 10$$

$$\frac{229}{2}$$

$$\text{Median} = L + \frac{\frac{N}{2} - f}{f} \times i$$

$$46 = 40 + \frac{114.5 - 42 - x}{65} \times 10.$$

$$46 = 40 + \dots$$

$$6 = \frac{76.5 - x}{65} \times 10$$

$$6 \times 65 = 76.5 - x \times 10.$$

$$\frac{6 \times 65}{10} = 76.5 - x$$

$$x = 76.5 - \frac{390}{10}$$

$$x = 76.5 - 39$$

$$\boxed{x = 37.5}$$

$$x + y = 79$$

$$37.5 + y = 79$$

$$y = 79 - 37.5$$

$$\boxed{y = 41.5}$$

21/10/21

Recall:

Median:  $\rightarrow$  Middle value of the data.

Formulas:

Individual observations - size of  $(\frac{N+1}{2})^{\text{th}}$  item  $\xrightarrow{\text{No. of observations}}$

Discrete observations - size of  $(\frac{N+1}{2})^{\text{th}}$  item,  $N = \sum f$

continuous observations -  $L + \frac{N - f}{2} \times i$   $\xrightarrow{\text{lower limit of median class}}$   $\xrightarrow{\text{frequency of median class}}$

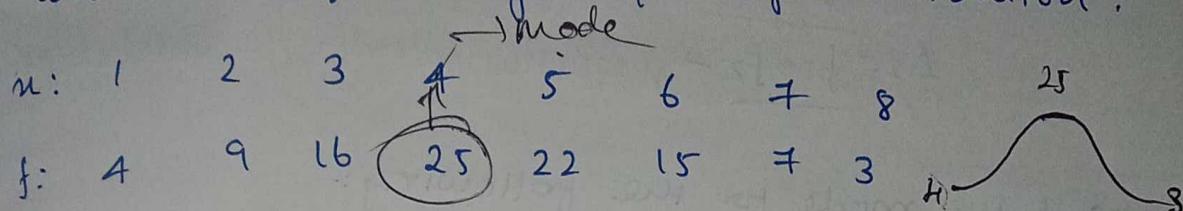
$\xrightarrow{\text{cumulative freq. preceding the median class}}$   $\xrightarrow{\text{length of class}}$

## Mode:

Mode is the value which occurs most frequently in a set of observations and around which the other items of the set cluster density.

### FQ: 1-

consider the following frequency distribution.



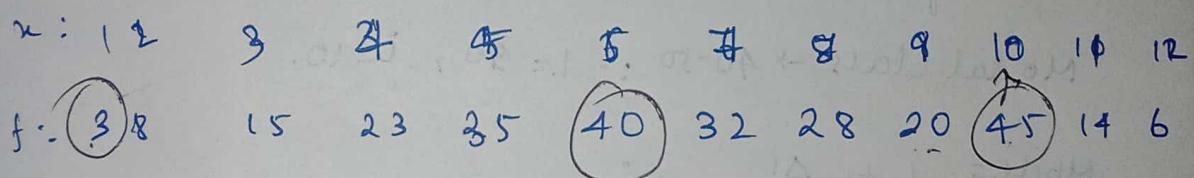
Find Mode = ?

(How to find mode)?

Mode = 4 // . 25 Which is maximum frequency that should be considered as the mode.

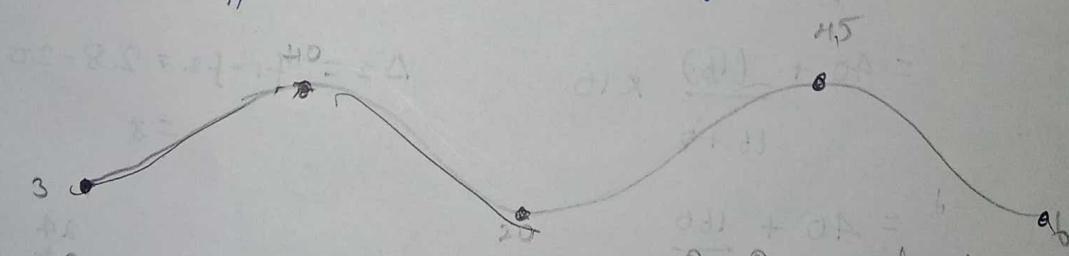
### FQ: 2

consider the following frequency distribution:



Find mode = ?

Mode = 10 // (Highest frequency).



(when we have like this data it is called as irregular data)

Relation b/w mean, median, mode.

$$\text{Mode} = 3\text{median} - 2\text{mean} \rightarrow \text{Formula}$$

## Continuous Observation:

$$\text{Mode} = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times i \rightarrow \text{length of the CI.}$$

lower limit of  $\Delta_1 + \Delta_2$

$$f_1 - f_0$$

modal class

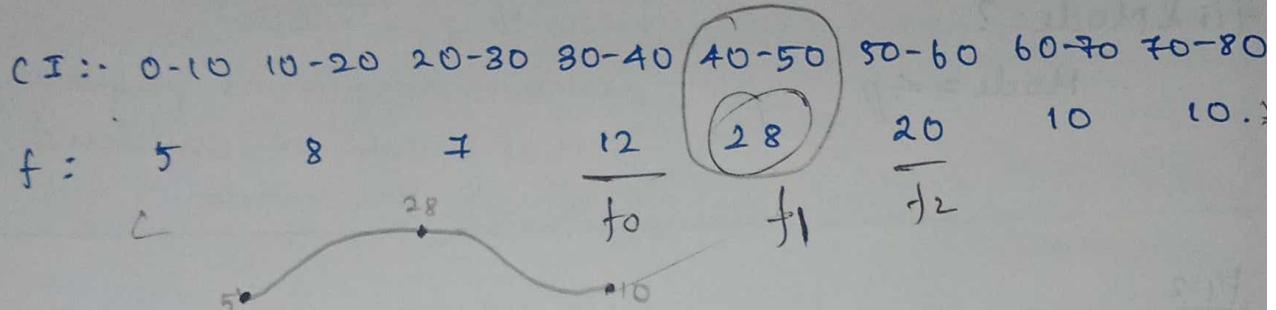
(CI which have highest frequency)

Where  $\Delta_1 = f_1 - f_0 \rightarrow$  frequency of modal class (highest frequency)

$\Delta_2 = f_1 - f_2 \rightarrow$  frequency preceding the modal class.

$\Delta_2 = f_1 - f_2 \rightarrow$  frequency succeeding the modal class.

• Find the mode for the following:-



Soln:-

Highest frequency  $f_1 = 28$ ,  $f_0 = 12$ ,  $f_2 = 20$ .

Modal class  $\rightarrow 40-50 \therefore L = 40$ ,  $i = 10$ .

$$\text{Mode} = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times i.$$

$$\Delta_1 = f_1 - f_0 = 28 - 12$$

$$= 16$$

$$= 16$$

$$= 40 + \frac{(16)}{16+8} \times 10$$

$$\Delta_2 = f_1 - f_2 = 28 - 20$$

$$= 8$$

$$= 40 + \frac{160}{24}$$

$$\begin{array}{r} 24 \\ 4 \\ \hline 960 \\ 1160 \\ \hline 1120 \end{array}$$

$$= \frac{960 + 160}{24}$$

$$= \frac{1120}{24}$$

$$= 46.666\ldots \text{ (approx)}$$

1922 Cheshire

$x$	$f_1$	columns. column 3		$f_4$	$f_5$	$f_6$
		$f_2$	$f_3$			
1	3	11				
2	8		23		46	
3	15	38			78	
4	23		58			
5	35	75		98		
6	40		72	107		
7	32	60			100	
8	28		48	80		
9	20	65		93		
10	45		59		79	
11	14	20		65		
12	6					

this is called grouping method Find the maximum value in column 2 i.e., 75,  $\rightarrow$  adding 1<sup>st</sup> & 2<sup>nd</sup> frequencies.

Find the maximum value in column 3 i.e., 72,  $\rightarrow$  2<sup>nd</sup> & 3<sup>rd</sup> frequencies.

In column 4  $\rightarrow$  we should add 1<sup>st</sup>, 3<sup>rd</sup> frequencies.

$\hookrightarrow$  Find the maximum value in C4 i.e., 98,

In C5 leave the 1<sup>st</sup> frequency and from 2<sup>nd</sup> add those frequencies i.e., 8 + 15 + 23

(b) Find the maximum value in C5 i.e., 107.

$\otimes$  n = 5, 6, 4, has got a role to play calculation of the mode.

In C6 omit the 1<sup>st</sup> & 2<sup>nd</sup> frequencies and sum of 3<sup>rd</sup> frequencies  $\rightarrow$  add

(c) Find the maximum value in C5 i.e., 100.

C1  $\rightarrow$  frequency

(2)  $\rightarrow$  sum of two frequency (consecutive)

(3)  $\rightarrow$  omit the 1<sup>st</sup> frequency and sum the two consecutive freq.)

(4)  $\rightarrow$  the sum of the 1<sup>st</sup> & 3<sup>rd</sup> frequency

(5) → omit the 1<sup>st</sup> frequency and sum the '3' consecutive

(6) → omit the 1<sup>st</sup> '2' frequency and sum the '3' " " "

{Identify the maximum in each column}

Max in C<sub>1</sub> → 45 i.e;  $\textcircled{10}$ .

Max in C<sub>2</sub> → 75 i.e; 5, 6.

Max in C<sub>3</sub> → 72 i.e; 6  $\textcircled{7}$

Max in C<sub>4</sub> → 98 i.e;  $\textcircled{4}, \textcircled{5}, \textcircled{6}$

Max in C<sub>5</sub> → 107 i.e; 5, 6, 7

Max in C<sub>6</sub> → 100 i.e; 6, 7  $\textcircled{8}$ .

(When we have

irregular data  
like this we want  
to do.)

$$x = 4 - f(1)$$

that like  
we want  
to calculate.

x	f
4	1
5	3
6	5 $\textcircled{5}$
7	3
8	1
10	1

→ Highest frequency  $\rightarrow 5$ .

So, the mode is 6 //

25/10/21

Recall: (last week)

central tendency.

concentrated in middle part of the data.

1) Mean, 2) Median, 3) Mode, 4) Geometric mean

5) Harmonic mean.

Classification of data:

1) Individual observation, 2) Discrete observation

3) continuous observation.

Mean:  $\bar{x} = \frac{\sum x}{n} \rightarrow$  Individual observations,

$\bar{x} = \frac{\sum f x}{\sum f} \rightarrow$  Discrete observations.

$\bar{x} = A + \frac{\sum f d}{\sum f} \times i \rightarrow$  Continuous observation  
Assumed mean where  $d = \frac{m - A}{i}$   
 $\rightarrow$  mid pt.  
 $i \rightarrow$  length of class

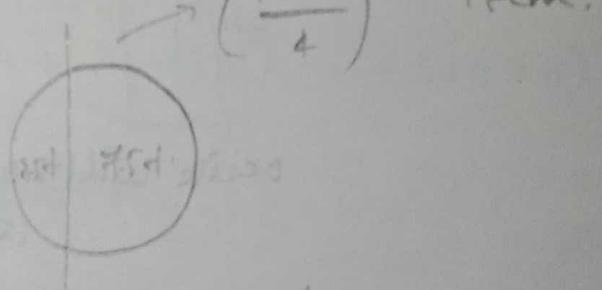
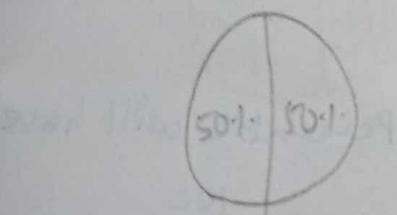
Median:  $\bar{x} =$  size of  $(\frac{N+1}{2})^{\text{th}}$  item  $\rightarrow$  Indi. obs

$$\rightarrow N = \sum f$$

size of  $(\frac{N+1}{2})^{\text{th}}$  item  $\rightarrow$  Discrete obs

$L + \frac{N}{2} - 4$   $\rightarrow$  cumulative freq. preceding the median class  
 $\downarrow$   $\frac{N}{2}$   $\times i \rightarrow$  continuous obs  
lower limit  $f \rightarrow$  frequency of the median class

Partining the values::



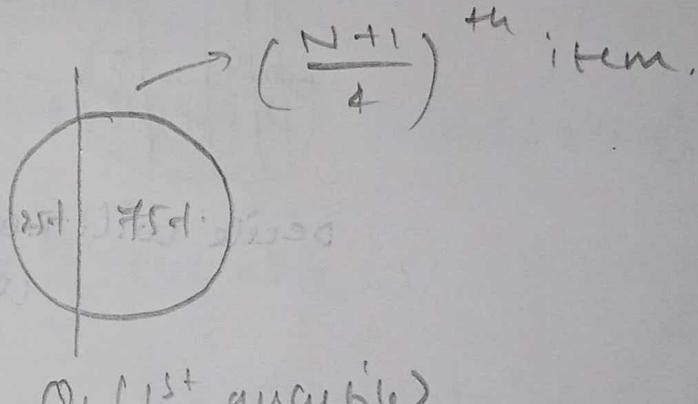
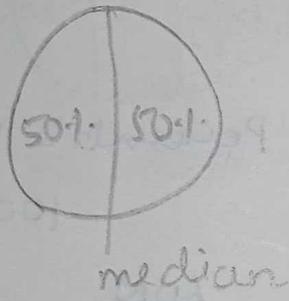
$Q_1 = \text{size of } (\frac{N+1}{4})^{\text{th}}$  item.

$$Q_1 = L + \frac{\frac{N}{4} - 4}{f} \times i$$

\* Median:  $\bar{x}_z$  size of  $(\frac{N+1}{2})^{\text{th}}$  item  $\rightarrow$  Indi. obs  
 $N = \sum f$   
 size of  $(\frac{N+1}{2})^{\text{th}}$  item  $\rightarrow$  Discrete obs

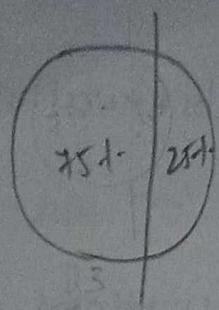
$L + \frac{N}{2} - \frac{f}{4}$   $\rightarrow$  cumulative freq. preceding the median class  
 $\downarrow$   
 $\frac{f}{4} \times i$   $\rightarrow$  continuous obser  
 lower limit  $f$   $\rightarrow$  frequency of the median class  
 of the median class

Partitioning the values:-



$Q_1$  is size of  $(\frac{N+1}{4})^{\text{th}}$  item.

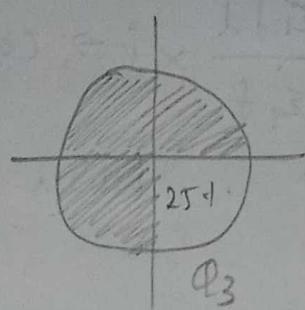
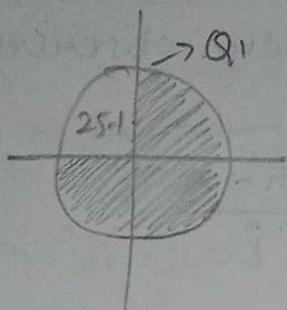
$$Q_1 = L + \frac{\frac{N}{4} - \frac{f}{4}}{f} \times i$$



$Q_3 = \text{size of } \left( \frac{N+1}{4} \right)^{\text{th}} \text{ item.}$

$$Q_3 = L + \frac{\frac{3N}{4} - 4}{f} \times i$$

↓  
lower limit of the 3<sup>rd</sup> quartile

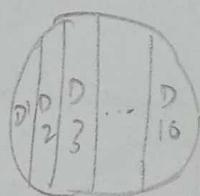


### Percentile, Decile

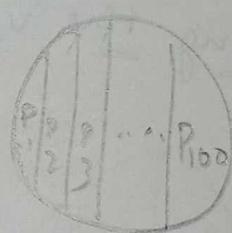
$$\text{Decile} = L + \frac{\frac{N}{10} - 4}{f} \times i$$

$$k^{\text{th}} \text{ decile} = L + \frac{\frac{kN}{10} - 4}{f} \times i$$

$$30 \rightarrow k = 3$$



Decile will have  
10



Percentile will have  
100

$$\text{If we want } 60^{\text{th}} \text{ percentile} = L + \frac{\frac{60N}{100} - 4}{f} \times i$$

↓  
6<sup>th</sup> decile.

$$\text{If we want } 50^{\text{th}} \text{ percentile} = L + \frac{\frac{50N}{100} - 4}{f} \times i$$

↓ 5<sup>th</sup> decile.  
Median  
(half half so it is median)

[Our 5<sup>th</sup> decile and 50<sup>th</sup> percentile are same as median]

Mode :- Most repeated value

$$\text{Mode} = L + \frac{\Delta_1}{\Delta_1 + \Delta_2} \times i$$

Modal class

Where,  $\Delta_1 = f_1 - f_0$   $\rightarrow$  frequency preceding the modal class,

$\Delta_2 = f_1 - f_2$   $\rightarrow$  frequency of the modal class

$\rightarrow$  frequency succeeding the modal class.

$$\boxed{\text{Mode} = 3 \text{Median} - 2 \text{Mean}}$$

$\hookrightarrow$  Relationship b/w mean, median, mode.

Geometric mean:

Geometric mean of a set of  $n$  observations is the  $n^{\text{th}}$  root of their product.

Individual observations:  $x_1, x_2, \dots, x_n$  are the given values.

$$G.M = G = (x_1 x_2 \dots x_n)^{1/n}$$

$$\log G = \frac{1}{n} (\log x_1 + \log x_2 + \dots + \log x_n)$$

$$G = \text{Antilog} \left\{ \frac{1}{n} \sum_{i=1}^n \log x_i \right\},$$

Formula for frequency (Discrete data):-

$$G = \text{Antilog} \left\{ \frac{1}{n} \sum_{i=1}^n f_i \log x_i \right\} \quad \& \quad N = \sum_{i=1}^n f_i.$$

$$G.M = \left( x_1^{f_1}, x_2^{f_2}, \dots, x_n^{f_n} \right)^{1/n} \rightarrow \text{where } n = \sum f$$

Continuous data:

$$AM = \text{Antilog}$$

$$x: 5 \quad 15 \quad 25 \quad 35$$

This is the discrete data

$$f: 3 \quad 4 \quad 7 \quad 6$$

Convert into continuous data.

Soln.

$$(I): 0-10 \quad 10-20 \quad 20-30 \quad 30-40$$

$$f: 3 \quad 4 \quad 7 \quad 6$$

If this is given we want to discrete data  
then we should find mid-pt

Disadvantages of Geometric mean (GM):

- (1) If one value is (-)ve then GM becomes difficult.
- (2) If any one value is zero, then  $GM=0$ .

Relationship b/w AM & GM.

$$AM \geq GM$$

Proof:

$a, b$  are the two values.

$$AM = \frac{a+b}{2}$$

$$GM = \sqrt{ab}$$

To prove:

$AM \geq GM$ . Equivalently, to prove  $AM - GM$

≥ 0.

$$\text{consider } AM - GM = \frac{a+b}{2} - \sqrt{ab}$$

$$= \frac{a+b - 2\sqrt{ab}}{2}$$

$$= \frac{(\sqrt{a}-\sqrt{b})^2}{2} \geq 0.$$

$\therefore AM \geq GM$ .

Hence proved //

1. Find the GM of  $-3, -5$ .

$$GM = \sqrt{-3 \times -5}$$

$$\sqrt{(-1)^2} = 1.$$

$$GM = \sqrt{15}.$$

$$\sqrt{-3}.$$

$$i = -\sqrt{1}$$

$$-1 = i^2 = \sqrt{i^4} = \sqrt{i^2 \cdot i^2}$$

$$= \sqrt{(-1)(-1)}$$

$$\sqrt{(-a)(-b)} = \pm \sqrt{ab}$$

$$= \sqrt{1}$$

$$= 1 //$$

$$\Rightarrow -1 = 1$$

Geometric mean of combined grp:

$$\log G = \frac{n_1 \log a_1 + n_2 \log a_2 + \dots + n_k \log a_k}{n_1 + n_2 + \dots + n_k}$$

$$a_i = (x_{i1}, x_{i2}, \dots, x_{in_i})^{1/n_i}$$

$$\log a_i = \frac{1}{n_i} \left\{ \sum_{j=1}^{n_i} \log x_{ij} \right\}$$

$$n_1 \log a_1 = \sum_{i=1}^{n_1} \log x_{1,i}$$

$$a_k = (x_{k1}, x_{k2}, \dots, x_{kn_k})^{1/n_k}$$

$$n_k \log a_k = \sum_{i=1}^{n_k} \log x_{ki}$$

combined grp:

$$G = (x_{11}, x_{12}, \dots, x_{1n_1}, x_{21}, x_{22}, \dots, x_{2n_2}, \dots, x_{k1}, \\ x_{k2}, \dots, x_{kn_k})^*$$

$$* \cdot \frac{1}{n_1 + n_2 + \dots + n_k}$$

$$\log G = \frac{n_1 \log a_1 + n_2 \log a_2 + \dots + n_k \log a_k}{n_1 + n_2 + \dots + n_k}$$

26/10/21: Recall:-

geometric mean:

$$a_M = (x_1 x_2 \dots x_n)^{1/m}$$

For calculation purpose:

(individual observation)

$$a = \text{Antilog} \left\{ \frac{1}{n} \sum_i \log x_i \right\}$$

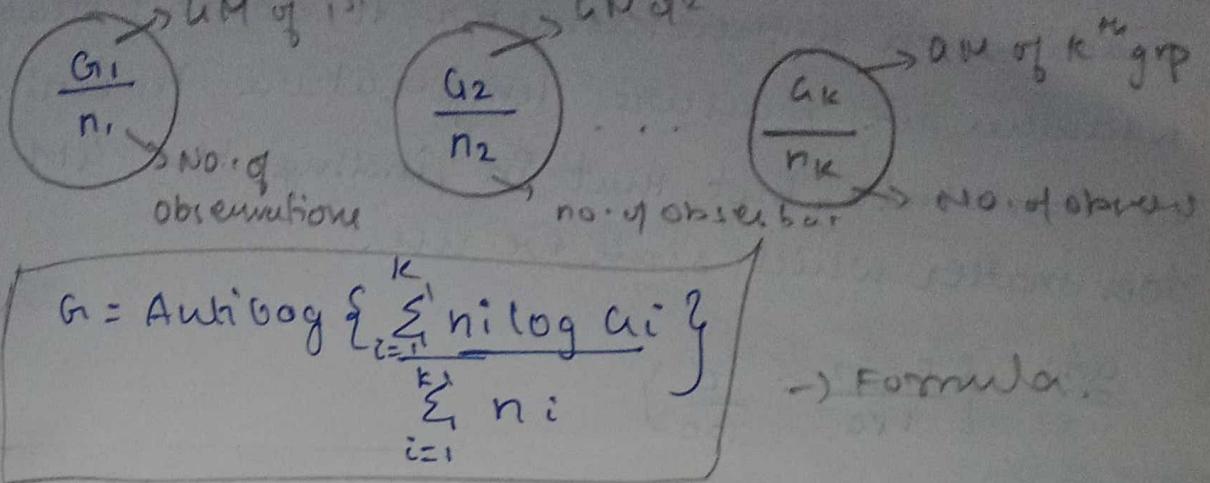
(discrete distributions)  $\rightarrow$

$$a = \text{Antilog} \left\{ \frac{1}{N} \sum_i f_i \log x_i \right\}, N = \sum_i f_i$$

↓ for both sum

(continuous distribution)

combined data for AM:



• the AM of 10 observations on a certain variable was calculated as 16.2. It was later discovered that one of the variable observation was wrongly recorded as 12.9 instead of 21.9. calculate the correct AM?

Soln:

Let  $x_1$  be 12.9. → was wrongly entered

$x_1'$  be 21.9. ,  $n=10$ .

$$a = (x_1 x_2 \dots x_n)^{1/10} = 16.2.$$

$$= (12.9 \times x_2 \times x_3 \dots x_{10})^{1/10} = 16.2$$

$$= (x_2 x_3 \dots x_{10})^{1/10} = \frac{16.2}{(12.9)}^{1/10}.$$

$$\text{Corrected AM} = (x_1' x_2 x_3 \dots x_{10})^{1/10} = \frac{(x_1')^{1/10}}{(x_2 x_3 \dots x_{10})^{1/10}}$$

$$= \frac{(21.9)^{1/10} \times 16.2}{(12.9)^{1/10}}.$$

(or)

$$\frac{(16.2 \times 21.9)^{1/10}}{(12.9)^{1/10}}$$

2. show that in finding the AM of a set of readings on a thermometer it doesn't matter whether we measure temperature in centigrade or Fahrenheit, but that in finding the GM it does matter which scale we use.

$$\text{AM: } \left\{ \frac{F-32}{180} = \frac{C}{100} \right\}$$

Soln:

$$F - 32 = \frac{\frac{9}{5}}{100} \times C$$

$$F = 32 + \frac{9}{5} C$$

suppose  $c_1, c_2, \dots, c_n$  are in centigrade

$$\text{AM} \Rightarrow \bar{C} = \frac{1}{n} \{ c_1 + c_2 + \dots + c_n \}$$

$$\text{GM} \Rightarrow G_C = (c_1 c_2 \dots c_n)^{1/n}$$

The observations corresponding to Fahrenheit would be;

$$32 + \frac{9}{5} c_1, 32 + \frac{9}{5} c_2, \dots, 32 + \frac{9}{5} c_n$$

$$\begin{aligned} \bar{f} &= \frac{1}{n} \left\{ \left( 32 + \frac{9}{5} c_1 \right) + \left( 32 + \frac{9}{5} c_2 \right) + \dots + \left( 32 + \frac{9}{5} c_n \right) \right\} \\ &= 32 + \frac{9}{5} \left\{ \frac{c_1 + c_2 + \dots + c_n}{n} \right\} = 32 + \frac{9}{5} \bar{C}. \end{aligned}$$

$$\begin{aligned} \text{G}_{cf} &= \frac{1}{n} \left\{ \left( 32 + \frac{9}{5} c_1 \right) \left( 32 + \frac{9}{5} c_2 \right) \dots \left( 32 + \frac{9}{5} c_n \right) \right\} \\ &\neq \frac{9}{5} \left( c_1 c_2 \dots c_n \right)^{1/n} + 32, \end{aligned}$$

3. In a frequency table the upperbound of each class interval has a constant ratio to the lower bound. show that GM is given by  $\log a$

$$\log a = \pi_0 + \frac{c}{N} \sum_i f_i (i-1)$$

where;  $\pi_0$  - logarithm of midvalue of 1<sup>st</sup> CI.

$c$  - " " " constant ratio.

ratio of upperbound and lowerbound

Soln:

let the  $i^{\text{th}}$  CI is denoted as;

$I_i - I_{i+1}$  & the frequency is  $f_i$

$$\text{Given } \frac{I_2}{I_1} = \frac{I_3}{I_2} = \dots = \frac{I_i}{I_{i-1}} \dots = \lambda \quad \begin{cases} i = 1, 2, \dots, n \\ I_1 - I_2 f_1 \\ I_2 - I_3 f_2 \\ I_3 - I_4 f_3 \end{cases}$$

$$I_i = \lambda I_{i-1}$$

$$= \lambda (\lambda I_{i-2}) = \dots = \lambda^{i-1} I_1.$$

let  $x_i$  be the mid-pt. of  $i^{\text{th}}$  class

$$x_i = \frac{1}{2} (I_1 + I_2) = \frac{1}{2} (I_1 + \lambda I_1)$$

$$= \frac{I_1}{2} (1 + \lambda) \quad I_2 = \lambda I_1$$

$$x_i = \frac{1}{2} (I_i + I_{i+1}) = \frac{1}{2} \left\{ \lambda^{i-1} I_1 + \lambda^i I_1 \right\}$$

$$= \frac{1}{2} \lambda^{i-1} I_1 (1 + \lambda).$$

$$x_1 = \frac{(1+\lambda)}{2} I_1 \quad \& \quad x_i = \frac{1}{2} \lambda^{i-1} I_1 (1 + \lambda)$$

$$x = \lambda^{i-1} x_1$$

let  $a$  be the GM;

$$\log a = \frac{\sum f_i \log x_i}{\sum f_i N}$$

$$\sum f_i N$$

$$= \frac{1}{N} \sum f_i \log (\lambda^{i-1} x_i)$$

$$\log a = \frac{1}{N} \sum f_i \log (x_i \lambda^{i-1})$$

$$= \frac{1}{N} \left\{ \sum f_i (\log x_i + \log \lambda^{i-1}) \right\}$$

$$= \frac{1}{N} \sum f_i \log x_i + \frac{1}{N} \sum f_i \log \lambda^{i-1}$$

$$= \frac{1}{N} \sum f_i \log x_i + \frac{1}{N} \sum f_i (i-1) \log \lambda.$$

$$= \frac{\log x_i}{N} \sum f_i + \frac{\log \lambda}{N} \sum f_i (i-1).$$

one and one  
 same

$$= x_0 + c \sum f_i (i-1),$$

### Harmonic Mean:- (HM)

HM of a no. of observations, none are zero ie; (nothing should be zero or in the given data) is the reciprocal of AM of the reciprocal of the given values.

i.e.;  $x_1, x_2, \dots, x_n$  is the given data then the harmonic mean  $\bar{x}$ ,

$$HM = \frac{1}{\frac{1}{n} \sum_{i=1}^n \left( \frac{1}{x_i} \right)}$$

(Individual  
observation)

$$HM = \frac{1}{\frac{1}{N} \sum_{i=1}^n \left( \frac{f_i}{x_i} \right)}, N = \sum f_i$$

(Both for discrete &  
continuous)

In case of continuous observation take  
 $x_i$  at the mid pt, and in discrete case as  
such as i.e.,  $f_i/x_i$ .

1. A cyclist spreads from his home to his college  
at a speed of 10 km/hr and back from college  
to his house at 15 km/hr. Find the average  
speed.

Solns:

$$\text{cyl} \rightarrow H \quad H \rightarrow \text{cyl}$$

To find the average  
speed use HM

Let ' $x$ ' be the distance b/w home and cyl

- (I) Home to college  $\rightarrow \frac{x}{10}$  km/hr (distance travelled)  
(II) College to home  $\rightarrow \frac{x}{15}$  hr (distance travelled).

The total distance  $2x$  is covered in  $\frac{x}{10} + \frac{x}{15}$  hrs.

Average speed =  $\frac{\text{Total distance}}{\text{Total time taken}}$  if  $\frac{\frac{x}{10} + \frac{x}{15}}{2}$

$$= \frac{2x}{x\left(\frac{1}{10} + \frac{1}{15}\right)}$$

$$= \frac{2}{\left(\frac{1}{10} + \frac{1}{15}\right)}$$

$$= \frac{2 \times 10 \times 15}{35} \\ = 12 \text{ hrs/1}$$

27/10/21

## Dispersion:-

Dispersion is scatteredness.

Why to study this?

\* To have an idea about the homogeneity and heterogeneity of the gn. data (obs. distribution,

Defn:-

The degree to which numerical data tend to spread about an average value is called dispersion / variation of data.

## Measures of dispersion which are widely used:-

- Range  $\rightarrow$  co-efficient of range
- Quartile deviation  $\rightarrow$  co-efficient of quartile dev.
- Mean deviation  $\rightarrow$  co-efficient of mean dev.
- Standard deviation  $\rightarrow$  co-efficient of variation.

↳ moments.

All these are independent  
of units

Squarus      Kurtosis  
 ↓                ↓  
 Size of the data      (height of the gn-data)  
                     Flatness of the data.

Suppose if we wish to compare dispersion of 2 samples.

We use co-efficient of dispersion.

(It is independent of units) Dispersion.

Formulas:-

$$* \text{Range} = L - S$$

$L \rightarrow$  smallest value  
 S  $\rightarrow$  largest value

$$\text{Co-efficient of range} = \frac{L-S}{L+S}$$

(The same formula for all 3 observations)

i.e; Individual, Discrete, Continuous

\* Quartile deviation =  $\frac{Q_3 - Q_1}{2}$   $\rightarrow$  3rd quartile

$$\text{Co-eff. of QD} = \frac{Q_3 - Q_1}{Q_3 + Q_1}$$

Mean deviation:

$$MD \text{ from } \alpha = \frac{1}{N} \sum_i f_i |x_i - \alpha|$$

$\alpha \Rightarrow$  mean  $\rightarrow$  mean deviation abt mean

$\alpha \Rightarrow$  median  $\rightarrow$  mean deviation abt median

$\alpha \Rightarrow$  mode  $\rightarrow$  mean deviation abt mode

\* Co-eff. of MP =  $\frac{MD}{\text{mean}}$

\* Co-eff. of median =  $\frac{MD}{\text{median}}$

\* Co-eff. of mode =  $\frac{MD}{\text{mode}}$

Standard deviation: ( $\sigma$  sigma)

$$\sigma = \sqrt{\frac{1}{N} \sum_i f_i (x_i - \bar{x})^2} \text{ Formula.}$$

$$\sum f_i = N \text{ (Total).}$$

\* called as Root mean square value.

\*  $\sigma^2$  = variance of the distribution.

Equivalent formula for SD :: (Alternate formula)

$$\sigma = \sqrt{\frac{\sum f d^2}{\sum f} - \left( \frac{\sum f d}{\sum f} \right)^2}$$

where  $d = \frac{m - A}{i}$

(Co-eff of variation) =

$$CV(x) = \frac{\sigma}{\bar{x}} \times 100$$

SD considered to be

(\*) one bcz it's  
used to compare

$$\sigma \rightarrow \sigma^2$$

1. For a group of 200 students the mean and standard deviation (SD) of scores are found to be 40 and 15 respectively. Later it was discovered that the scores 43 and 35 were misread as 34 and 53 respectively. Find the correct mean and SD?

Soln:

An:  $n = 200$ ,  $\sigma = 15$

$$\bar{x} = 40$$

To calculate corrected mean & SD.

$$\bar{x} = \frac{\sum x}{n} = \frac{\bar{x} \cdot n}{40 \cdot 40}$$

Corrected value  
34 → 43  
53 → 35  
misread 35

$$\sum x = n \cdot \bar{x} = 200(40) = 8000, 53 \text{ " } 35$$

$$\begin{aligned} \text{corrected } \sum x &= \sum x - 34 - 53 + 43 + 35 \\ &= 8000 - 9. \end{aligned}$$

$$\text{Corrected mean} = \frac{7791}{200}$$

$$\sigma^2 = \frac{\sum x^2}{n} - \left(\frac{\sum x}{n}\right)^2 \quad (\text{Individual observation})$$

$$\text{Corrected } \sum x^2 = 40$$

$$15^2 = \frac{\sum x^2}{200} - (40^2)$$

$$\bar{x} = \frac{\sum f}{n}$$

$$225 = \frac{\sum x^2}{200} - (1600)$$

$$\frac{\sum x^2}{200} = 1825$$

$$\sum x^2 = 1825 \times 200$$

$$\sum x^2 = 365000$$

$$\begin{aligned} \text{Corrected } \sum x^2 &= \sum x^2 - \underbrace{34^2}_{1156} - \underbrace{53^2}_{2809} + \underbrace{43^2}_{1849} + \underbrace{35^2}_{1225} \\ &= 365000 - 1156 - 2809 + 1849 \\ &\quad + 1225 \end{aligned}$$

$$= 364109$$

$$\text{Corrected } \sigma^2 = \frac{364109}{200} - \left(\frac{7991}{200}\right)^2$$

$$= 224.1429$$

$$\sigma = \sqrt{224.1429}$$

$$\sigma = 14.97$$

2. The 1<sup>st</sup> of the two samples has  $\frac{100}{n}$  items with mean 15 and SD 3. If the whole group

has 250 items with mean 15.6 and SD  $\sqrt{13.44}$

Find the SD of the 2<sup>nd</sup> grp.

Soln:

{ Combined SD is given by

$$\sigma = \sqrt{\frac{\sum d_i^2}{n_1 + n_2 + \dots + n_k}}$$

$$\frac{n_1(\sigma_1^2 + d_1^2) + n_2(\sigma_2^2 + d_2^2) + \dots + n_k(\sigma_k^2 + d_k^2)}{n_1 + n_2 + \dots + n_k}$$

Where  $d_i = \bar{x}_i - \bar{x}$  ( $i=1, 2, \dots, k$ )

$\bar{x}$  = combined mean

Soln:

$$\bar{x}_1 = 15, n_1 = 100, n_2 = 150$$

20

$\bar{x}$  = 15.6 (combined mean)

$$\sigma = 3, \sigma^2 = 13.44 \text{ (combined variance)}$$

$$\bar{x} = \frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2} \Rightarrow 15.6 = \frac{100(15^2 + 150(\bar{x}_2))}{250}$$

$$\bar{x}_2 = \frac{15.6 \times 250 - 1500}{150}$$

$$\bar{x}_2 = \frac{2400}{150} = 16$$

m

$$\sigma^2 = \frac{n_1(\sigma_1^2 + d_1^2) + n_2(\sigma_2^2 + d_2^2)}{n_1 + n_2}$$

$$d_1 = \bar{x}_1 - \bar{x}$$

$$13.44 = \frac{100(9 + 0.36) + 150(\sigma_2^2 + 0.16)}{250} = 15 - 15.6 \\ = -0.6$$

$$d_2 = \bar{x}_2 - \bar{x}$$

$$13.44 \times 250 = 936 + 150(\sigma_2^2 + 0.16) = 16 - 15.6 \\ = 0.4$$

$$3360 - 936 = 150 (\sigma_2^2 + 0.16)$$

$$2424 = 150 (\sigma_2^2 + 0.16)$$

$$\sigma_2^2 = 16 //$$

$$\sigma_2 = 4 //$$

$$\bar{x}_2 = 16$$

$$\sigma_2 / \sqrt{10} = 4$$

28/10/21

Recall:

\* Measures of dispersion:

⇒ Range → L-S

⇒ Quartile deviation →  $\frac{Q_3 - Q_1}{2}$  divide by

⇒ Mean deviation (MD) →  $\frac{\sum f |m - \bar{x}|}{\sum f}$

⇒ Standard deviation (SD) →  $\sigma = \sqrt{\frac{1}{n} \sum x^2 - (\frac{\sum x}{n})^2}$

⇒ moments.

$\sigma^2$  = variance (when we square the SD).

Co-eff. of dispersion:

⇒ Co-eff. of range →  $\frac{L-S}{L+S}$

(Independent  
of units)

⇒ Co-eff. of Q.D →  $\frac{Q_3 - Q_1}{Q_3 + Q_1}$

⇒ Co-eff. of MD →  $\frac{MD}{\text{mean}}$

⇒ Co-eff. of SD variation →  $\frac{\sigma}{\bar{x}} \times 100$ .

Combined SD:

$$\frac{n_1(\sigma_1^2 + d_1^2) + n_2(\sigma_2^2 + d_2^2) + \dots + n_k(\sigma_k^2 + d_k^2)}{n_1 + n_2 + \dots + n_k}$$

$$n_1 + n_2 + \dots + n_k$$

where  $d_i = \bar{x}_i - \bar{x}$

( $\rightarrow$ ) combined mean

$$\frac{n_1\bar{x}_1 + n_2\bar{x}_2 + n_3\bar{x}_3 + \dots + n_k\bar{x}_k}{n_1 + n_2 + n_3 + \dots + n_k}$$

1. calculate  $Q_D$  and co-eff. of  $Q_D$  for the foll.

C.I.: 20-30 30-40 40-50 50-60 60-70 70-80 80-90

f: 3 61 132 153 140 51 2

soln:

C.I.	f	c.f
20-30	3	3
30-40	61	64
40-50	132	196
50-60	153	349
60-70	140	489
70-80	51	540
80-90	2	542
	<u>542</u>	

$$Q_D = \frac{Q_3 - Q_1}{2}$$

$$Q_3 = L + \frac{\frac{3N}{4} - c}{f} \times i$$

$$\frac{3N}{4} = \frac{3(542)}{4}$$

$$= 406.5$$

$\Sigma$  the comming no.

$$i.e., 489$$

$$Q_3 = 60 + \frac{406.5 - 349}{140} \times 10.$$

$$Q_3 = 64.107 //$$

$$Q_1 = L + \frac{\frac{N}{4} - c}{f} \times i$$

$$\frac{N}{4} = \frac{542}{4} = 135.5$$

$$Q_1 = 40 + \frac{135.5 - 64}{132} \times 10$$

$$40 + 0.5416 \times 10$$

$$Q_1 = 45.416 //$$

$$Q_1 = 45.416, Q_3 = 64.107$$

$$\begin{aligned} QD &= \frac{Q_3 - Q_1}{2} \\ &= \frac{64.107 - 45.416}{2} \\ &= 9.345 \end{aligned}$$

$$\text{Co-eff of } QD = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{9.345}{109.523} = 0.085 // 0.17068$$

2. Find the SD for the following-

Mean: 0-4 4-8 8-12 12-16 16-20 20-24.

No. of st. 10 12 18 4 5 3.  
Students

Soln:

whichever we can take the assumed mean  
but for calculation it's better to take  
big no.

C.I	f	m	$d = \frac{m - A}{4}$	$\sum fd$	$\sum fd^2$
0-4	10	2	-3	-30	90
4-8	12	6	-2	-24	48
8-12	18	10	-1	-18	18
12-16	4	14	0	0	0
16-20	5	18	1	5	5
20-24	3	22	2	6	12
	55			-61	173

$$\begin{aligned} \sigma &= \sqrt{\frac{\sum fd^2}{\sum f} - \left( \frac{\sum fd}{\sum f} \right)^2} \times 4 \\ &= \sqrt{\frac{173}{55} - \left( \frac{-61}{55} \right)^2} \times 4 \end{aligned}$$

$$= \sqrt{3.145 - 123} \times 4$$

$$= 1.3838 \times 4$$

$$= 5.535$$

$$\sigma = \sqrt{\frac{\sum fm^2}{\sum f} - \left( \frac{\sum fm}{\sum f} \right)^2}$$

(or this formula)

3. Find the mean deviation about mean and co-eff of M.D for the foll.

CI: 20-25 25-30 30-35 35-40 40-45 45-50 50-55 55-60

f: 35 45 70 105 90 74 51 30

Soln:

CI	f	m	$d = \frac{m-A}{5}$	fd
20-25	35	22.5	-3	-105
25-30	45	27.5	-2	-90
30-35	70	32.5	-1	-70
35-40	105	37.5	0	0
40-45	90	42.5	1	90
45-50	74	47.5	2	148
50-55	51	52.5	3	158
55-60	30	57.5	4	120
				246
				500

$$\bar{x} = A + \frac{\sum fd}{\sum f} \times i$$

$$= 37.5 + \left( \frac{246}{500} \right) \times 5$$

$$= 37.5 + 2.46$$

$$\bar{x} = 39.96$$

$$MD = \frac{\sum f |m - \bar{x}|}{\sum f}$$

$$= \frac{3904.6}{500}$$

$$MD = 7.809$$

	$ m - \bar{x} $	$f m - \bar{x} $
17.46	611.1	
12.46	560.7	
7.46	522.2	
2.46	258.3	
2.54	228.6	
7.54	557.96	
12.54	639.54	
17.54	526.2	3904.6

$$\text{Co-eff. MD} = \frac{7.809}{39.96} \times \frac{\text{MD}}{\text{mean.}}$$

$$= 0.1954.$$

$$\boxed{\text{MD} = \frac{\sum f |m - \bar{x}|}{\sum f}}$$

4. An analysis of monthly wages paid to the workers of two firms A and B belonging to the same industry gives the foll.

	firm A	firm B
No. of workers n	500 $\rightarrow n_1$	600 $\rightarrow n_2$
Avg. daily wage $\bar{x}_1$	₹ 186 $\rightarrow x_1$	₹ 175 $\rightarrow x_2$
Variance of wage $\sigma^2$	$\sqrt{81} \rightarrow \sigma_1^2$	$\sqrt{100} \rightarrow \sigma_2^2$

$\sigma^2 = \sigma_1^2 + \sigma_2^2 - 2 \cdot \text{Combined mean} \cdot \text{SD}$ .  
we want to find in this

- (i) Which firm A or B has large wage variability?
- (ii) Which firm A or B is there greater variability in individual wages?  $\rightarrow$  Comb. variance.
- (iii) calculate the combined mean avg. daily wage and variance of the distribution of the wages of all the workers in firm A and firm B together?

Soln:

$$(i) CV(A) = \frac{\sigma_1}{\bar{x}_1} \times 100$$

$$\boxed{CV(x) = \frac{\sigma}{\bar{x}} \times 100.}$$

$$= \frac{9}{186} \times 100 = \frac{900}{186}$$

$$= 4.838.$$

$$CV(x^2) = \frac{\sigma^2}{\bar{x}^2} \times 100.$$

$$\sqrt{100} = 10.$$

$$CV(B) = \frac{\sigma_2}{\bar{x}_2} \times 100$$

$$= \frac{10}{175} \times 100 = \frac{1000}{175}$$

$$= 5.714.$$

$$CV(B) > CV(A)$$

$\Rightarrow$  B has got greater variability when compared with A.

(II)  $n_1 = 500, n_2 = 600$

$$\bar{x}_1 = 186, \bar{x}_2 = 175$$

Total wages in Firm A =  $n_1 \bar{x}_1$

$$= 500 \times 186$$

$$= 93000$$

$$\begin{array}{r} 186 \\ \underline{- 53} \\ 93000 \end{array}$$

Total wages in Firm B =  $n_2 \bar{x}_2$

$$= 600 \times 175$$

$$= 105000$$

More wage is paid by firm B.

(III) Combined mean =  $\frac{n_1 \bar{x}_1 + n_2 \bar{x}_2}{n_1 + n_2}$

$$= \frac{93000 + 105000}{500 + 600}$$

$$\begin{array}{r} 105000 \\ - 93000 \\ \hline 12000 \end{array}$$

$$= \frac{198000}{1100} = 180$$

Combined SD =  $\sqrt{\frac{n_1 (\sigma_1^2 + d_1^2) + n_2 (\sigma_2^2 + d_2^2)}{n_1 + n_2}}$

$$= \frac{500 \cdot (81 + 36) + 600 \cdot (125)}{500 + 600}$$

$$= \frac{500(117) + 600(125)}{1100}$$

$$= \frac{500(117) + 600(125)}{1100} = \frac{133500}{1100}$$

$$= 121.36$$

01/11/21.

Recall:

Measures of dispersion:

① Range = L - S

② Quartile deviation =  $\frac{Q_3 - Q_1}{2}$

③ Mean deviation  
about mean  $\bar{x} = \frac{\sum f |x - \bar{x}|}{\sum f}$

④ Standard deviation  $\Rightarrow \sigma = \sqrt{\frac{\sum f d^2}{\sum f} - \left( \frac{\sum f d}{\sum f} \right)^2}$   
where  $d = \frac{m - A}{i}$

Variance =  $(\sigma^2)$ .

Co-eff of dispersion:

① Co-eff of range =  $\frac{L - S}{L + S}$

② Co-eff of mean deviation =  $\frac{M.D \text{ of mean}}{\text{mean}}$

③ Co-eff of QD =  $\frac{Q_3 - Q_1}{Q_3 + Q_1}$

④ Co-eff of variation =  $\frac{s}{\bar{x}} \times 100$

(Variability of the data)

## Moments:

Generalization of the SD.

The  $r^{\text{th}}$  moment of a variable  $x$  about any pt.  $x = A$ , denoted by  $M_r'$  is given by

$$(\text{raw moment}) M_r' = \frac{1}{\sum f} \sum_i f_i (x_i - A)^r$$

The  $r^{\text{th}}$  moment of a variable  $x$  about the mean  $\bar{x}$ , denoted by  $M_r$  is defined as  $M_r = \frac{1}{\sum f} \sum_i f_i (x_i - \bar{x})^r$

$$M_r = \frac{1}{\sum f} \sum_i f_i (x_i - \bar{x})^r$$

### Note:

$$M_r = \frac{\sum_i f_i (x_i - \bar{x})^r}{\sum f_i}$$

(i) when  $r=0$ ;

$$M_0 = \frac{\sum f_i}{\sum f_i} = 1.$$

(ii)  $r=1$ ;

$$M_1 = \frac{\sum f_i (x_i - \bar{x})}{\sum f_i} = \frac{\sum f_i x_i - \bar{x} \sum f_i}{\sum f_i}$$

$$= \bar{x} - \bar{x} \cdot 1 = 0.$$

(iii)  $r=2$ ;

$$M_2 = \frac{\sum f_i (x_i - \bar{x})^2}{\sum f_i} = \sigma^2.$$

# Relationship b/w $M_r$ and $M_1'$

$$\text{Let } d_i = x_i - A \quad \& \quad \bar{x} = A + \frac{\sum f_i d_i}{\sum f_i}$$

$$x_i - A = M_1'$$

$$= A + M_1'$$

$$\boxed{\Rightarrow \bar{x} = A + M_1'}$$

$$N = \sum f_i$$

$$x_i - A = M_1'$$

$$M_r' = \frac{1}{N} \sum f_i (x_i - \bar{x})^r \cdot \boxed{x_i - A = M_1'} \\ = \frac{1}{N} \sum f_i \{x_i - A + A - \bar{x}\}^r \\ = \frac{1}{N} \sum f_i \{d_i - M_1'\}^r$$

$$M_r = \frac{1}{N} \sum f_i (d_i - M_1')^r$$

$$= \frac{1}{N} \sum_i f_i \{d_i^r - r_{c_1} M_1'^{r-1} d_i^{r-1} + r_{c_2} M_1'^{r-2} d_i^{r-2} (M_1')^2 + \dots + (-1)^r (M_1')^r\}$$

$$= \frac{1}{N} \sum f_i d_i^r - r_{c_1} M_1' \underbrace{\frac{\sum f_i d_i^{r-1}}{N}}_{M_{r-1}} + r_{c_2} (M_1')^2 \underbrace{\frac{\sum f_i d_i^{r-2}}{N}}_{M_{r-2}} + \dots + (-1)^r (M_1')^r.$$

$$M_r = M_r' - r_{c_1} M_{r-1} M_1' + r_{c_2} M_{r-2} (M_1')^2 + \dots + (-1)^r (M_1')^r$$

$$M_r = M_r' - r_{c_1} M_{r-1} M_1' + r_{c_2} M_{r-2} (M_1')^2 + \dots + (-1)^r (M_1')^r$$

when  $r=2$ ;

$$M_2 = M_2' - 2c_1 M_1' M_2' + (-1)^2 (M_1')^2 = M_2' - (M_1')^2$$

(Covariance)

when  $r=3$ :

$$M_3 = M_3' - 3C_1 M_2' M_1' + 3C_2 M_1' (M_1')^2 - \cancel{3C_3} (M_1')^3$$
$$M_3 \Rightarrow M_3' - 3M_2' M_1' + 2(M_1')^3.$$

when  $r=4$ :

$$M_4 = M_4' - 4C_1 M_3' M_1' + 4C_2 M_2' (M_1')^2 - (4C_3) M_1' (M_1')^3$$
$$+ (-1)^4 (M_1')^4$$
$$= M_4' - 4M_3' M_1' + 6M_2' M_1' - 3(M_1')^4.$$

$$M_1 = \underbrace{\frac{1}{4} \int f(x-\bar{x}) dx}_{\text{if}} = \underbrace{\frac{1}{4} \int f(x) dx}_{\text{if}} - \bar{x} = 0$$

$$\boxed{M_1 = M_1' - M_1' = 0} \rightarrow \textcircled{1}$$

Skewness  $\beta_1 = \frac{M_3^2}{M_2^3}$  {bulkness of the data}.

Kurtosis  $\beta_2 = \frac{M_4}{M_2^2}$  {flatness of the curve}

For a symmetrical distribution:

$$\boxed{\beta_2 = 3}$$

① Raw moments

② central moments

③ Relationship b/w raw and central moments

④ Particular cases in central moments.

$$M_2 = \text{variance} = M_2' - (M_1')^2$$

$$M_3 \& M_4.$$

⑤  $\beta_1 = \frac{M_3^2}{M_2^3}$

① The first four moments about the value  $x=4$   
 are  $-1.5, 17, -30, 108$ . Find skewness & kurtosis.  
 soln:-  $M_1' M_2' M_3' M_4'$

$$M_1' = -1.5, M_2' = 17, M_3' = -30, M_4' = 108.$$

$$\text{To find } \beta_1 = \frac{M_3^2}{M_2^3} \text{ & } \beta_2 = \frac{M_4^2}{M_2^2}$$

$$M_2 = M_2' - (M_1')^2$$

$$= 17 - (-1.5)^2$$

$$= 17 - 2.25$$

$$= 14.75$$

$$M_3 = M_3' - 3M_2'(M_1')^2 + 2(M_1')^3$$

$$\beta_1 = \frac{M_3^2}{M_2^3}$$

$$\beta_2 = \frac{M_4^2}{M_2^2}$$

$$= -30 - 3(17)(-1.5)^2 + 2(-1.5)^3$$

$$= -30 + 51(1.5) - 2(1.5)^3$$

$$= -30 + 76.5 - 67.5$$

$$= 39.75$$

$$\beta_1 = \frac{M_3^2}{M_2^3} = \frac{39.75^2}{(14.75)^3} \quad \text{Skewness.}$$

$$M_4 = M_4' - 4M_3'M_1' + 6M_2'(M_1')^2 - 3(M_1')^4$$

$$= 108 - 4(-30)(-1.5) + 6(17)(-1.5)^2 - 3(-1.5)^4$$

$$= 108 - 120(1.5) + 6(17)(2.25) - 3(2.25)^4$$

$$= 142.3125$$

$$\beta_2 = \frac{M_4^2}{M_2^2} = \frac{142.3125}{(14.75)^2} \quad \text{Kurtosis}$$

$$= 0.6541$$

Note:-

$$\bar{x} = A + M_1' = 4 + (-1.5) = 2.5$$

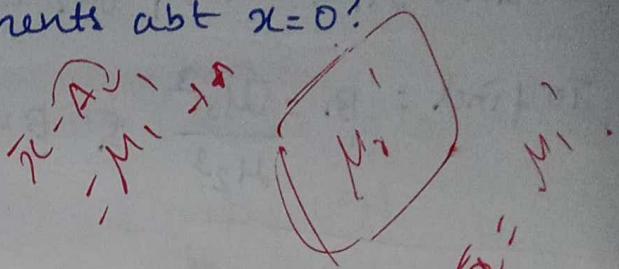
$$\boxed{\bar{x} = 2.5}$$

How will you find moments abt  $x=0$ ?

$$\boxed{\bar{x} = A + M_1'}$$

Here assume  $A=0$ ,

$$\Rightarrow \bar{x} = M_1' = 2.5$$



$$M_2 = M_2' - (M_1')^2$$

$$\Rightarrow M_2' = M_2 + (M_1')^2 = 14 - 75 + (2.5)^2$$

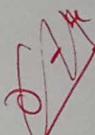
$$= 14 - 75 + 6.25$$

$$= 21\frac{1}{4}$$

2. Find the 1st & 4 moments for the following.

$x: 0 \ 1 \ 2 \ 3 \ 4 \ 5 \ 6 \ 7 \ 8$

$f: 1 \ 8 \ 28 \ 56 \ 70 \ 56 \ 28 \ 8 \ 1$



Soln:-

$x$	$f$	$d = x - A$	$fd$	$fd^3$	$fd^2$	$fd^4$	$\sum fd = 0$
0	1	-4	-4	-64	16	256	
1	8	-3	-24	-216	72	648	$\sum fd^2 = 512$
2	28	-2	-56	-224	112	560	$\sum fd^3 = 0$
3	56	-1	-56	-56	56	51	
4	70	0	0	0	0	0	$\sum fd^4 = 2816$
5	56	1	56	56	56	56	
6	28	2	56	224	112	448	
7	8	3	24	216	72	648	
8	1	4	4	64	16	256	

$$\bar{x} = A + \frac{\sum fd}{\sum f} = 4 + \frac{0}{4} = 4.$$

$$M_1 = 0$$

4  $\times$  &  $+$

$$M_2 = \frac{\sum fd^2}{\sum f} = \frac{512}{256} = 2$$

$\begin{cases} 2 \\ 2 \\ 2 \\ 2 \end{cases}$

$$M_3 = \frac{\sum fd^3}{\sum f} = \frac{0}{256} = 0$$

$$M_4 = \frac{\sum fd^4}{\sum f} = \frac{2816}{256} = 11$$

2/11/21

correlation.

Amount of relationship b/w two variables.

linear correlation!

A relationship b/w two variables

x	↑	↓	↑	↓
y	↑	↓	↓	↑

multiple correlation

partial correlation

$r_{12 \cdot 3} \rightarrow$  partial

$R_{12 \cdot 3} \rightarrow$  multiple

Carl Pearson's co-efficient of correlation

r or  $r$  (rho).

$$r = \frac{N \sum xy - \sum x \sum y}{\sqrt{N \sum x^2 - (\sum x)^2} \sqrt{N \sum y^2 - (\sum y)^2}}$$

Individual  
(correlation  
formula)

$$r = \frac{\text{covariance}(x, y)}{\sigma_x \sigma_y} \quad i.e., r = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

$\text{SD of } x \text{ and } y \leftarrow \sigma_x \sigma_y$

If x and y are independent then it is zero.

if;  $r=0$  x and y are independent.

Properties:

### Properties of correlation co-eff.

1.  $r$  lies b/w  $-1 \leq r \leq 1$

2. correlation is unaffected by the change of scale and origin. (Reduce big to smaller no's)

### Individual observations:

$$r = \frac{N \sum xy - \bar{x} \bar{y}}{\sqrt{N \sum x^2 - (\sum x)^2} \sqrt{N \sum y^2 - (\sum y)^2}}$$

### Discrete / continuous:

$$r = \frac{N \sum f dx dy - \sum f dx \sum f dy}{\sqrt{N \sum f dx^2 - (\sum f dx)^2} \sqrt{N \sum f dy^2 - (\sum f dy)^2}}$$

$\downarrow$

$$f = \sum f, dx = \frac{x-A}{n}, dy = \frac{y-B}{K}$$

i. calculate correlation co-eff.

$x: 65 \ 66 \ 67 \ 68 \ 69 \ 70 \ 72$ .

~~$y: 67 \ 68 \ 65 \ 68 \ 72 \ 72 \ 69 \ 71$~~

Soln:-

$$r = \frac{N \sum xy - \bar{x} \bar{y}}{\sqrt{N \sum x^2 - (\sum x)^2} \sqrt{N \sum y^2 - (\sum y)^2}}$$

(or)

$$r = \frac{N \sum d x d y - \sum d x \sum d y}{\sqrt{N \sum d x^2 - (\sum d x)^2} \sqrt{N \sum d y^2 - (\sum d y)^2}}$$

$\rightarrow dx = \frac{x-A}{n}$   
 $\rightarrow dy = \frac{y-B}{K}$

individual

$x$	$y$	$dx = x - 68$	$dy = y - 69$	$dx^2$	$dy^2$	$dx dy$
65	67	-3	-2	9	4	6
66	68	-2	-1	4	1	2
67	65	-1	-4	1	16	4
67	68	-1	-1	1	1	1
68	71	0	3	0	9	0
69	72	1	3	1	9	3
70	69	2	0	4	0	0
72	71	4	2	16	4	8
		<u>0</u>	<u>0</u>	<u>36</u>	<u>14</u>	<u>24</u>

$$N = 8$$

$$\sum dx = 0 = \sum dy$$

$$\sum dy \sum dx^2 = 36$$

$$\sum dy^2 = 44$$

$$\sum dx dy = 24$$

$$r = \frac{N \sum dx dy - \sum dx \sum dy}{\sqrt{N \sum dx^2 - (\sum dx)^2} \sqrt{N \sum dy^2 - (\sum dy)^2}}$$

$$= \frac{8(24) - 0 \times 0}{\sqrt{8(36) - 0^2} \sqrt{8(44) - 0^2}}$$

$$= \frac{192}{6\sqrt{0 \times 352}}$$

$$= \frac{192}{6\sqrt{2816}}^{32}$$

$$= \frac{32}{\sqrt{2816}} = 0.603 //$$

2. The following table gives the number of blind per lakh of population in different age groups. Find out the correlation b/w age and blindness.

Age: 0-10 10-20 20-30 30-40 40-50 50-60 60-70 70-80

No. of blind 55 67 100 111 150 200 300 500  
blindness.

Soln:

$\Sigma dy$	$m$	$y$	$\sum dx = m - 45$	$\sum dx^2$	$\sum dy = y - 150$	$\sum dy^2$	$\sum dx dy$
0-10	5	55	-4	16	-95	9025	380
10-20	15	67	-3	9	-93	6889	249
20-30	25	100	-2	4	-50	2500	100
30-40	35	111	-1	1	-49	2401	49
40-50	45	116	0	0	0	0	0
50-60	55	200	1	1	50	2000	50
60-70	65	300	2	4	150	92000	300
70-80	75	300	3	9	350	122500	1050
			<u>-4</u>	<u>44</u>	<u>273</u>	<u>168315</u>	<u>2178</u>

$$N = 8, \sum dx = -4, \sum dx^2 = 44, \sum dy = 273, \sum dy^2 = 168315$$

$$\frac{1}{N} \sum dx dy = 2178$$

$$r = \frac{\sum N \sum dx dy - \sum dx \sum dy}{\sqrt{N \sum dx^2 - (\sum dx)^2} \sqrt{N \sum dy^2 - (\sum dy)^2}}$$

$$= \frac{8(2178) - (-4)(273)}{\sqrt{8(-4) - (-4)^2} \sqrt{8(168315) - (273)^2}}$$

$$= \frac{18516}{\sqrt{336} \sqrt{1271991}}$$

$$= \frac{18516}{18.333 \times 112782}$$

$$= 0.8956 //$$

09/11/21

The following table gives the marks obtained by 100 students. Find the co-eff. of correlation.

$$\begin{matrix} x & y & n & y \\ \uparrow & \downarrow & \uparrow & \downarrow \\ \{ & \{ & \{ & \} \end{matrix} (E)$$

Age Mental	18	19	20	21	Total.
10-20	4	2	2	-	8
20-30	5	4	6	4	19
30-40	6	8	10	11	35
40-50	4	4	6	8	22
50-60	2	2	4	4	10
60-70	-	2	3	1	6
Total	19	22	31	28	100.

$$r = \frac{N \sum xy - \sum x \sum y}{\sqrt{N \sum x^2 - (\sum x)^2} \sqrt{N \sum y^2 - (\sum y)^2}}$$

↓                      ↓  
 $s.d(x)$        $s.d(y)$

Properties:

$$(i) -1 \leq r \leq 1$$

(ii) Co-eff of correlation is unaffected change of scale and origin.

$$r = \frac{N \sum f dx dy - \sum f dx \sum f dy}{\sqrt{N \sum f dx^2 - (\sum f dx)^2} \sqrt{N \sum f dy^2 - (\sum f dy)^2}}$$

$$r = \frac{N \sum f dx dy - \sum f dx \sum f dy}{\sqrt{N \sum f dx^2 - (\sum f dx)^2} \sqrt{N \sum f dy^2 - (\sum f dy)^2}}$$

where, No. of observation

$$N = \sum f$$

$$dx = \frac{x - A}{h}$$

$$dy = \frac{y - B}{k}$$

where;

$$N = \sum f$$

$$dx = \frac{y - A}{h}$$

$$dy = \frac{y - B}{k}$$

$x$	18	19	20	21	Total	
$y$	-2	-1	0	1		$dx = x - 20$
$dy$	$\frac{y-20}{x-20}$	$\frac{y-20}{x-20}$	$\frac{y-20}{x-20}$	$\frac{y-20}{x-20}$		$dy = \frac{y-15}{10}$
10-20	-3	4	2	2	0	8
20-30	-2	5	4	6	4	19
30-40	-1	6	8	10	11	35
40-50	0	4	4	6	8	22
50-60	1	0	2	4	4	10
60-70	2	0	2	3	1	6
f	-	19	22	31	28	120

$$(P(x) - P(20)) / (x - 20) = f(x)$$

$$(P(x) - P(20)) / (x - 20) = f(x)$$

$$f(x) = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

express between 21 and 22 is  $f(20+0)(21)$

approximate value also

$$\frac{f(21) - f(20)}{21 - 20} = \frac{f(21) - f(20)}{1} = f(21)$$

$$(P(x) - P(20)) / (x - 20) = f(x)$$

$$x = 21$$

$$f(21) = \frac{f(21) - f(20)}{21 - 20} \cdot 21 + f(20)$$

$$f(21) = 120$$

$x$	18	19	20	21	f	$\sum f dx$	$\sum f dy$	$\sum f dy^2$	$\sum f dx dy$
$dx$	-2	-1	0	1	-				
$dy$	4	2	2	0	8	-24			
10-20	-3	4	2	0	8	-24	72	30	
20-30	-2	5	4	6	4	19	-38	76	20
30-40	-1	6	8	10	11	35	-35	35	9
40-50	0	4	4	6	8	22	10	0	0
50-60	1	0	2	4	4	10	10	10	2
60-70	2	0	2	3	1	6	12	24	-2
$f$	-	19	22	31	28	(100)	*	*	59
$\sum f dx$	-32	-22	0	28					
$\sum f dx^2$	76	72	0	28					
$\sum f dy$	56	16	0	-13	59				

$$dx = x - 20$$

$$dy = \frac{y - 45}{10}$$

$$\sum f dx dy = 59$$

$$\sum f dx = -32$$

$$\sum f dy = -75$$

$$\sum f dx^2 = 126$$

$$\sum f dy^2 = 217$$

$$r = \frac{\sum f dx dy - \sum f dx \sum f dy}{\sqrt{N \sum f dx^2 - (\sum f dx)^2} \sqrt{N \sum f dy^2 - (\sum f dy)^2}}$$

$$= \frac{100(59) - (-32)(-75)}{\sqrt{100(126) - (-32)^2} \sqrt{100(217) - (-75)^2}}$$

$$= 0.2566 //$$

21. A computer while calculating co-eff of correlation b/w  $x$  and  $y$  from 25 pairs of observations obtained the following.

$$\sum xy = 508, n = 25; \sum x = 125, \sum x^2 = 650, \sum y = 100, \sum y^2 = 460$$

\* It was later discovered that at the time of checking he had copied two pairs wrongly and while all correct values

$x$	$y$
6	14
9	6

all  $\begin{array}{|c|c|} \hline x & 8 & 6 \\ \hline y & 6 & 8 \\ \hline \end{array}$  obtain correct correlation co-eff

Soln:

$x$	$y$
8	6
6	8

$$r = \frac{N \sum xy - \sum x \sum y}{\sqrt{N \sum x^2 - (\sum x)^2} \sqrt{N \sum y^2 - (\sum y)^2}}$$

$$\text{corrected } \sum x = 125 - 6 - 9 + 8 + 6 = 124.$$

$$\text{corrected } \sum y = 100 - 14 - 6 + 6 + 8 = 94.$$

$$\text{corrected } \sum x^2 = 650 - 36 - 81 + 64 + 36 = 633.$$

$$\text{corrected } \sum y^2 = 460 - 196 - 36 + 36 + 64 = 328$$

$$\text{corrected } \sum xy = 508 - 6 \times 14 - 9 \times 6 + 6 \times 8 + 8 \times 6$$

$$= 466.$$

$$\text{corrected } r = \frac{25(466) - (124)(94)}{\sqrt{25(633) - (124)^2} \sqrt{25(328) - (94)^2}}$$

Rank correlation or (Spearman's correlation)

(co-eff).

Data is given according to two different characteristics if we require co-eff of correlation, we use rank correlation.

$x$ : → Data arranged according to height of the student

$y$ : → Data arranged according to marks.

We 1<sup>st</sup> assign ranks and find the correlation

$$\rho = \text{rho} = 1 - \left[ \frac{6 \sum D^2}{N(N^2 - 1)} \right]$$

Rank<sup>x</sup> rank<sup>y</sup>

$$D = R_x - R_y$$

$N \rightarrow$  No. of observations

For repeated ranks;

$$\rho = 1 - \frac{6 \sum D^2 + \frac{1}{12} (m^3 - m) + \dots}{N(N^2 - 1)}$$

→ If rank '3' is repeated  
then take 'm' as '3'.

$m = \text{no. of times}$

→ If rank '5' is repeated.  
then take 'm' as '5'.

a particular rank  
gets repeated

Example: Find the rank correlation co-eff  
of the following

$X: 65 \ 67 \ 63 \ 62 \ 65$

$Y: 19 \ 18 \ 23 \ 18 \ 17$

rank<sup>x</sup> rank<sup>y</sup>  
5/1 3/2 4/3 3/2 5/1  
2/2 1/1 3/3 2/2 2/2

Soln:

$X$	$Y$	$R_x$	$R_y$	$D = R_x - R_y$	$D^2$
65	19	2.5	2	-0.5	0.25
67	18	1	3.5	-2.5	6.25
63	23	4	3	1	9.00
62	18	5	3.5	1.5	2.25
65	17	2.5	5	-2.5	6.25

$$\rho = 1 - \frac{6 \left\{ \sum D^2 + \frac{1}{12} (m^3 - m) \right.}{N(N^2 - 1)} + \frac{\left. \frac{1}{12} (m^3 - m) \right\}}$$

$$= 1 - \frac{6 \left\{ 24 + \frac{1}{12} (8-2) + \frac{1}{12} (8-2) \right\}}{5(25-1)}$$

$$= 1 - \frac{6 \times 24}{5 \times 24}$$

$$= 1 - \frac{6 \left\{ 24 + \frac{1}{2} + \frac{1}{2} \right\} 5}{5 \times 24}$$

$$= 1 - \frac{\sum D^2}{2N} = -\frac{1}{4}$$

$$= -0.25 //$$

Q11121

Rank correlation:

$$\rho = 1 - \frac{6 \left\{ \sum D^2 + \frac{1}{12} (m^3 - m) + \frac{1}{12} (m^3 - m) + \dots \right\}}{N(N^2 - 1)}$$

$$D = R_x - R_y$$

$m \rightarrow$  no. of times a particular rank gets repeated

quality of two observations:-

$x \rightarrow$  Height of the students

$y \rightarrow$  Marks scored by the students.

10. Ten competitions in a beauty contest were ranked by 3 judges.

R<sub>x</sub>: 4 2 8 6 1 5 3 9 10 7

R<sub>y</sub>: 2 5 9 3 6 7 1 10 8 4

R<sub>z</sub>: 4 3 6 9 2 8 7 5 1 10.

which pair of judges has the nearest approach to common taste?

<u>soln.</u>	R <sub>x</sub>	R <sub>y</sub>	R <sub>z</sub>	D <sub>1</sub> = R <sub>x</sub> - R <sub>y</sub>	D <sub>1</sub> <sup>2</sup>	D <sub>2</sub> = $\frac{D_2^2}{R_x - R_z}$	D <sub>2</sub> <sup>2</sup>	D <sub>3</sub>	D <sub>3</sub> <sup>2</sup>
	4	2	4	+2	4	0	0	-2	4
	2	5	3	-3	9	-1	1	2	4
	8	9	6	-1	1	2	4	3	9
	6	3	9	3	9	-3	9	-6	36
	1	6	2	-5	25	-1	1	4	16
	5	7	9	-2	4	-3	9	-6	36
	3	1	7	2	4	-4	16	5	25
	9	10	5	-1	1	4	16	7	49
	10	8	1	2	4	9	81	-6	36
	7	1	10	-3	9	-3	9	-6	36

$$\sum D_1^2 = 70$$

$$\sum D_2^2 = 146$$

$$\sum D_3^2 = 216$$

$$P_{XY} = 1 - \frac{6 \sum D_1^2}{N(N^2-1)}$$

$$= 1 - \frac{6(70)}{10(99)} = 0.542$$

$$P_{XZ} = 1 - \frac{6 \sum D_2^2}{N(N^2-1)}$$

$$= 1 - \frac{6(146)}{10(99)} = 0.115$$

$$P_{YZ} = 1 - \frac{6 \sum D_3^2}{N(N^2-1)}$$

$$= 1 - \frac{6(216)}{10(99)} = -0.309$$

$P_{XY} = 0.542 \rightarrow$  Judges X & Y have common taste

$$P_{YZ} = -0.309$$

$$P_{XZ} = 0.115$$

2. Find the rank correlation co-eff, given

(X)

Rank in math: 48 57 82 65 39 78 26 47 50.85 48 57

(Y)

Rank in stats: 56 43 80 45 79 50 46 72 65 71 32 59

R <sub>X</sub>	R <sub>Y</sub>	D = R <sub>X</sub> - R <sub>Y</sub>	D <sup>2</sup>	$\sum D^2 = 259$
8.5	7	1.5	2.25	
15.5	11	5.5	30.25	
2	1	1	1.00	
4	10	-6	36	
11	2	9	81	
3	8	-5	25	
12	10.9	3	9	
10	3	7	49	
7	5	2	4	
1	4	-3	9	
9.5	12	-3.5	12.25	
6	0.5	0.5	0.25	

$$t = 1 - \frac{6 \left\{ \sum D^2 + \frac{1}{12} (m^3 - m) + \dots \right\}}{N(N^2-1)}$$

$$= 1 - \frac{6 [259 + 1/12 (2^3 - 2) + \frac{1}{2} (2^3 - 2)]}{130(12(143))}$$

$$= 1 - \frac{6 (259 + 1)}{12 (143)}$$

$$= 1 - \frac{130}{143}$$

$$= \frac{13}{143}$$

## Regression lines ::

Approx value for the gn. variable -

\* Regression line of  $x$  on  $y$

( $x$  depends on  $y$ )

$$(x - \bar{x}) = b_{xy} (y - \bar{y})$$

where;  $b_{xy} = \frac{N \sum xy - \sum x \sum y}{N \sum y^2 - (\sum y)^2}$

$$\Rightarrow r = \frac{\sigma_x}{\sigma_y}$$

Standard deviation  
SD

$b_{xy}$  also called as the regression of co-eff  $x$  on  $y$ .

\* Regression line of  $y$  on  $x$ :

( $y$  depends on  $x$ )

$$(y - \bar{y}) = b_{yx} (x - \bar{x})$$

where;  $b_{yx} = \frac{N \sum xy - \sum x \sum y}{N \sum x^2 - (\sum x)^2}$

$$\Rightarrow r = \frac{\sigma_y}{\sigma_x}$$

Property :: (Regression co-eff)

$$(1) r^2 = b_{xy} b_{yx}$$

i.e, when we multiply the regression co-eff then automatically square of the correlation co-eff

$$\frac{x-y}{\downarrow} \in z \rightarrow r_{12.3}$$

Partial correlation  
(keep one variable constant)

Multiple correlation

(effect of one variable over other variables)

It is denoted by  $R_{123}$

NOTE:

If both  $b_{xy}$  and  $b_{yx}$  are (-)ve then  
 $r = -ve.$

Eq:  $b_{xy} = -0.5$  &  $b_{yx} = -\frac{1}{0.5}$

$$r^2 = (-0.5) \left(-\frac{1}{0.5}\right) = 1$$

$$\Rightarrow r = -1$$

Eg: Calculate regression lines for the following.

X: 1 2 3 4 5

Y: 10 12 7 3 8

Soln:-

X	Y	$x^2$	$y^2$	$xy$
1	10	1	100	10
2	12	4	144	24
3	7	9	49	21
4	3	16	9	12
5	8	25	64	40
Total	<u>15</u>	<u>40</u>	<u>366</u>	<u>107</u>

$$b_{xy} = \frac{N \sum xy - \sum x \sum y}{N \sum x^2 - (\sum x)^2}$$

$$= \frac{5(107) - (15)(40)}{5(55) - (15)^2}$$
$$= -1.3$$

$$b_{yx} = \frac{N \sum xy - \sum x \sum y}{N \sum y^2 - (\sum y)^2}$$

$$= \frac{5(107) - (15)(40)}{5(366) - (40)^2}$$

$$b_{xy} = -1.3$$

$$b_{yx} = -0.283.$$

$$= -0.283.$$

$$\bar{x} = \frac{\sum x}{N} = \frac{15}{5} = 3; \quad \frac{\sum y}{N} = \bar{y} = \frac{40}{5} = 8.$$

$\Rightarrow x \text{ on } y$  ( $x$  depends on  $y$ )

$$x - \bar{x} = 5 \times y (\bar{y} - \bar{y})$$

$$x - 3 = -0.283 (\bar{y} - 8)$$

$$x = 3 - 0.283 y + 8(0.283)$$

$$x = 5.264 - 0.283 \cancel{y}$$

$$\boxed{x = 5.264 - 0.283 \cancel{y}}$$

$\Rightarrow y \text{ on } x$  ( $y$  depends on  $x$ )

$$y - \bar{y} = b y x (x - \bar{x})$$

$$y - 8 = -1.3 (x - 3)$$

$$y = 8 - 1.3 x + 3 - 9$$

$$\boxed{y = 11.9 - 1.3 x}$$

$$x \text{ on } y \Rightarrow x = 5.264 - 0.283 y.$$

$$y \text{ on } x \Rightarrow y = 11.9 - 1.3 x.$$

Given  $x$ , to find  $y$ , use  $y \text{ on } x$ .

Given  $y$ , to find  $x$ , use  $x \text{ on } y$ .

Note

E.g. • find  $y$ , given  $x = 6$

$$y \text{ on } x \text{ is } y = 11.9 - 1.3 x$$

$$\Rightarrow y(6) = 11.9 - 1.3(6)$$

$$= 11.9 - 7.8 = 4.1$$

• find  $x$ , given  $y = 6$ .

$$x \text{ on } y \text{ is } x = 5.264 - 0.283 y$$

$$\Rightarrow x(6) = 5.264 - 0.283(6)$$

$$= 3.566$$

$y \text{ on } x$  is  $11.9 - 1.3 x$ .

$$\boxed{x=1} \quad y = 11.9 - 1.3(1)$$

$X=1$

$$= 11.9 - 1.3$$

$$= 10.6$$

$$X=1, Y=10.6.$$

$$\begin{array}{r} 11.9 \\ 1.3 \\ \hline 10.6 \end{array}$$

111121

i) For 2 variables  $x$  and  $y$  the equations of the regression lines are  $9y - x - 288 = 0$ , and  $x - 4y + 38 = 0$ . Find

- (i) mean values of  $x$  &  $y$ .
- (ii) co-eff of correlation regression.
- (iii) The ratio of the SD. of  $y$  to that of  $x$ .
- (iv) Most probable value of  $y$  when  $x = 145$ .
- (v) Most probable value of  $x$  when  $y = 35$ .

→ (i)  $x - \bar{x} = b_{xy}(y - \bar{y})$  } means

$y - \bar{y} = b_{yx}(x - \bar{x})$  } P. book of x series

Soln:-

(i)  $x - \bar{x} = b_{xy}(y - \bar{y})$

$y - \bar{y} = b_{yx}(x - \bar{x})$

~~$-x + 9y - 288 = 0$~~

~~$x - 4y + 38 = 0$~~

~~$\frac{5y - 250 = 0}{+}$~~

$-x + 9y - 288 = 0$

$-x + 9(50) - 288 = 0$

$-x + 450 - 288 = 0$

$x = -288 + 9y$

$5y = 250$

$y = \frac{250}{5}$

$\boxed{Y = 50}$

$\boxed{\bar{x} = 172}$

(ii) Let  $9Y - X - 288 = 0$  represent regression line of  $X$  on  $Y$ .

$$X = 9Y - 288$$

$$X = 9 \left( Y - \frac{288}{9} \right)$$

$$X = 9Y$$

$$b_{XY} = 9.$$

The line  $X - 4Y + 38 = 0$  is regression line of  $Y$  on  $X$ .

$$4Y = X + 38$$

$$Y = \frac{1}{4} (X + 38)$$

$$Y = \frac{1}{4}$$

$$b_{YX} = \frac{1}{4}$$

$$\boxed{r^2 = b_{XY} b_{YX}}$$
$$= 9 \times \frac{1}{4}$$

$$r = \frac{3}{2} > 1. \text{ (impossible)}$$

∴  $9Y - X - 288 = 0$  is regression line of  $X$  on  $Y$

$X - 4Y + 38 = 0$  is regression line of  $Y$  on  $X$

$$X = 9Y - 288 = 9 \left( Y - \frac{288}{9} \right) \quad b_{XY} = 9$$

$$+ 4Y = X + 38 \Rightarrow Y = \frac{1}{4} (X + 38) \quad b_{YX} = 1$$

$$9Y = X + 288 \Rightarrow Y = \frac{1}{9} (X + 288) \Rightarrow b_{YX} = \frac{1}{9}$$

$$X = 4Y - 38 \Rightarrow X = 4 \left( Y - \frac{38}{4} \right) \Rightarrow b_{XY} = 4$$

$$r^2 = b_{XY} b_{YX} = \frac{4}{9} \Rightarrow r = \sqrt{\frac{4}{9}} = \frac{2}{3}$$

$$(iii) \frac{\sigma_y}{\sigma_x} = ?$$

by  $x = \frac{1}{9}$

$$\frac{\sigma_y}{\sigma_x} = \frac{1}{9}$$

$$\frac{\sigma_y}{\sigma_x} = \frac{1}{9} \gamma = \frac{1}{9 \times \frac{2}{3}} = \frac{1}{6}$$

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

(iv) Regression line of  $y$  on  $x$  is  $94 - x - 28 = 0$

when  $x = 145$

$$94 - 145 - 28 = 0$$

$$94 = 293.433$$

$$y = \frac{433}{9}$$

(v) Regression line of  $x$  on  $y$  is  $4x - 4y + 38 = 0$

when  $y = 35$

$$x = 4y - 38$$

$$x = 4(35) - 38$$

$$= 140 - 38$$

$$= 102$$

2. A panel of judges A and B graded seven detectors and independently awarded the marks.

Judge A: 40 34 28 30 44 38 31 36

Judge B: 32 39 26 30 38 34 27

Eighth debator was awarded 36 marks by judge A while judge B was not present. If judge B was also present, how many marks would be allot for the eighth debator.

Soln:-

Judge A =  $x$ , Judge B =  $y$ .

To find regression line of  $y$  on  $x$

then substitute  $x=36$  in the above to

find  $y$ .

$x$	$y$	$dx = x - \bar{x}$	$dy = y - \bar{y}$	$\sum dx dy$	$\sum dx^2$	$\sum dy^2$	
40	32	5	0	0	25	0	
34	39	-1	7	-7	1	49	$y = A + \frac{\sum dy}{n}$
28	26	-7	-6	42	49	36	
36	13	-5	-2	10	25	4	
44	38	9	6	54	81	36	$= 32 + \frac{3}{7}$
38	34	3	2	6	9	4	$= 32 + 0.4$
31	28	-4	-4	16	16	16	$= 32.4 //$
		<u>0</u>	<u>3</u>	<u>121</u>	<u>206</u>	<u>145</u>	

$$\bar{x} = 35$$

$$b_{yx} = \frac{N \sum dx dy - \sum dx \sum dy}{N \sum dx^2 - (\sum dx)^2}$$

$$= \frac{7(121)}{7(206)} = 0.5873 //$$

Regression line of  $y$  on  $x$  is

$$y - \bar{y} = b_{yx} (x - \bar{x})$$

$$Y - 32 \cdot 4 = 0.5873X - 0.05873(35) + 32 \cdot 4$$

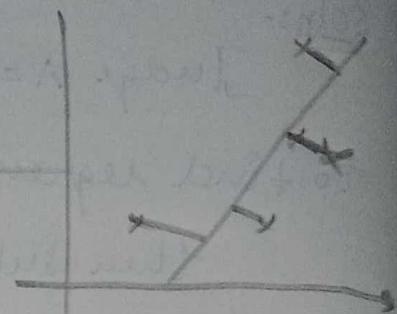
$$Y = 0.5873X + 32 \cdot 4 - 0.5873(35)$$

Given  $x = 36 \Rightarrow Y = 0.5873(36) + 32 \cdot 4 - 20.55$

$$= 53.546 - 20.55$$

$$= 32.991$$

Method of least squares:



\* Fit  $y = a + bx$

$$\begin{array}{l} \text{Normal eqns} \\ \left\{ \begin{array}{l} \sum y = na + b \sum x \\ \sum xy = a \sum x + b \sum x^2 \end{array} \right. \end{array} \rightarrow \text{(fitting a line)}$$

\* Fit  $y = a + bx + cx^2$ .

$$\begin{array}{l} \text{Normal eqn.} \\ \left\{ \begin{array}{l} \sum y = na + b \sum x + c \sum x^2 \\ \sum xy = a \sum x + b \sum x^2 + c \sum x^3 \\ \sum x^2y = a \sum x^2 + b \sum x^3 + c \sum x^4 \end{array} \right. \end{array} \rightarrow \text{(fitting a parabola)}$$

\* Fit  $y = ae^{bx}$  (exponential curve).

$$y = ae^{bx}$$

$$\log_e y = \log_e(ae^{bx})$$

$$\log_e y = \log_e a + \log_e e^{bx}$$

$$\log_e y = A + bx \quad (A = \log_e a)$$

$$x = A + bx$$

$$x = \log_e y.$$

$$\begin{array}{l} \text{Normal eqn} \\ \left\{ \begin{array}{l} \sum x = nA + b \sum x \\ \sum x^2 = A \sum x + b \sum x^2 \end{array} \right. \end{array}$$

$$\left\{ \begin{array}{l} \sum x^2 = A \sum x + b \sum x^2 \end{array} \right.$$