

Basic Statistics

For S1 =

P1+P2+P3+F1+F2

Measure of central tendency (Average)

1) Arithmetic Mean

2) Median

3) Mode

4) Geometric Mean

5) H.M (Harmonic Mean)

Measure of central tendency or an average is a single significant figure which sums up the characteristics of a group of figures. Various measures of central tendency are arithmetic mean, median mean, mode, geometric mean, harmonic mean.

Median mean, mode, geometric mean, harmonic mean

Arithmetic Mean :-

case I - Individual Series

let x_1, x_2, \dots, x_n are n values then arithmetic mean of these numbers is given by $A.M = \frac{x_1 + x_2 + \dots + x_n}{n}$

$$A.M = \frac{\sum x}{n}$$

Q Find mean of 10, 12, 8, 14, 16, 19 ...

$$\frac{10+12+8+14+16+19}{6} = 13.167$$

Q Find mean of 2.6, 8.4, 1.5, 10.9, 4.5, 6.8, 7.5, 5.9, 6.7, 8.8

$$\frac{2.6+8.4+1.5+10.9+4.5+6.8+7.5+5.9+6.7+8.8}{10} = 6.36$$

Discrete frequency

Let x_1, x_2, \dots, x_n are n values with corresponding frequencies f_1, f_2, \dots, f_n then $A_m = \frac{f_1x_1 + f_2x_2 + \dots + f_nx_n}{f_1 + f_2 + \dots + f_n}$

Q Calculate mean from the following data

value 5 15 25 35 45 55 65 75

freq 15 20 25 24 12 31 71 52

x	f	fx
5	15	75
15	20	300
25	25	625
35	24	840
45	12	540
55	31	1705
65	71	4615
75	52	3900
	250	12600

The following data related to the distance travelled by 520 villagers to buy their weekly requirement

Case III mean in continuous frequency distribution

In the case of continuous frequency distribution we can take mid value of the class & then $A_m = \frac{\sum f_n}{N}$

Q Find the Arithmetic mean from the following data

Age : 0-10	10-20	20-30	30-40	40-50	50-60	60-70	70-80
No. of : 15	30	53	75	100	110	115	125

Class	Frequency (f)	Midvalue (x)	$f \times f$
0-10	15	5	75
10-20	30	15	450
20-30	53	25	1325
30-40	75	35	2625
40-50	100	45	4500
50-60	110	55	6050
60-70	115	65	7475
70-80	125	75	9375
	623		31,875

$$\text{Arithmetic Mean} = \frac{31875}{623} = 51.16$$

∴ arithmetic mean (approx) is 51.16
 \therefore mid value \times freq. of given bin gives us

Median

Median is the value of that item which occupies the central position when the items are arranged in the ascending or descending order of their magnitude. Therefore median is the value of that item which has equal number of items above & below that is, number of items greater than median and number of items less than median are equal.

(a) Median in Individual Series

Arrange the values in the data in the ascending or descending order of their magnitude and find out the value of middle item. It is the median. It is also called size of $(\frac{n+1}{2})^{\text{th}}$ item.

Q Find the median for the following values?

4, 45, 60, 20, 83, 19, 26, 11, 27, 12, 52

- writing the values in the ascending order

4, 11, 12, 19, 20, 26, 27, 45, 52, 60, 83

Median = Size of $(\frac{n+1}{2})^{\text{th}}$ item = 6th item = 26

when 'n' is even, there are two middle items.
Take the average of them.

Q calculate median

35, 23, 45, 50, 80, 61, 92, 40, 52, 61

- writing the values in the ascending order

23, 35, 40, 45, 50, 52, 61, 61, 80, 92

median = size of $(\frac{n+1}{2})^{\text{th}}$ item = size of $(\frac{(n+1)}{2})^{\text{th}}$

item = size of 5. 5th item = $\left(\frac{50+52}{2}\right) = 51$

(b) Median in Discrete Frequency distribution

Formula : Median = size of $(\frac{N+1}{2})^{\text{th}}$ Item

where $N = \sum f$

Q calculate Median

Size : 5 8 10 15 20 25

Frequency : 3 12 8 7 5 4

- Size Frequency (f) Cumulative Frequency

5

3

3

8

12

$3 + 12 = 15$

10

8

$15 + 8 = 23$

$\frac{N+1}{2}$

15

7

$23 + 7 = 30$

20

5

$30 + 5 = 35$

25

4

$35 + 4 = 39$

$N = 39$

Median = size of $\left(\frac{N+1}{2}\right)^{\text{th}}$ item = size of $\left(\frac{10}{2}\right)^{\text{th}}$
item = size of 20th item

Median corresponds to 20th item of the series.

The first cumulative frequency which includes 20 is 33.

The size of item for which cumulative frequency is 23, is 10 \therefore median = 10

(a) Median in continuous frequency distribution
In the case of continuous frequency distribution, median class corresponds to the cumulative frequency which includes $N/2$. After getting median class, find median by using the following interpolation formula

$$\text{Median} = l_i + \frac{\left(\frac{N}{2} - cf\right)}{f} \times c$$

where l_i is the lower limit of median class.

'cf' is the cumulative frequency of the class just preceding the median class. 'f' is the frequency of the median class. 'c' is the interval of median class.

$$3 \times \frac{(N-1)u}{f}$$

Q Compute median

- class : 0-10 10-20 20-30 30-40 40-50 50-60 60-70	$\frac{N}{2} = \frac{1+1}{2} = 1$
Frequency : 8 12 20 23 18 7 2	

class	frequency	cumulative frequency	$\frac{N}{2}$
0-10	8	8	
10-20	12	20	
20-30	20	40	
30-40	23	63	$\frac{N}{2}$
40-50	18	81	
50-60	7	88	
60-70	2	90	
	90		

Median = Size of $(\frac{N}{2})^{\text{th}}$ item = size of $(\frac{90}{2})^{\text{th}}$ item = size of 45th item. 45 is included in the cumulative frequency 63. The class having cumulative frequency 63 is 30-40.

∴ 30-40 is the median class.

Applying interpolation formula,

$$\text{Median} = l_1 + \frac{(\frac{N}{2} - cf)}{f} \times c$$

Here $l_1 = 30$, $N = 2 = 45$, $c_f = 40$, $f = 23$, $c = 10$

$$\text{median} = 30 + \frac{(45-40)}{23} \times 10 = 30 + \frac{50}{23} = 30 + 2.7 \\ = 32.7$$

Q Find the median of the values

17, 32, 35, 33, 15, 21, 41, 32, 11, 10, 20

- 10, 11, 15, 17, 20, 21, 32, 32, 33, 35, 41

medium = size of $(\frac{n+1}{2})^{\text{th}}$ item = 6^{th} item = 21

Q The following table shows age of 8 student find median

SINo:	1	2	3	4	5	6	7	8
Age :	18	16	14	11	13	10	9	20

class	Frequency	Cumulative Frequency
1	18	18
2	16	34
3	14	48
4	11	59
5	13	72
6	10	82
7	9	91
8	20	101

$$N = 101 \quad \frac{N+1}{2} = \frac{8+1}{2} = \frac{9}{2} = 4.5 \quad = \frac{13+14}{2} \\ = 13.5$$

9, 10, 11, 13, 14, 16, 18, 20

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

$$= \frac{101+1}{2} = 51^{\text{th}} \text{ item}$$

The 1^{st} cumulative frequency which include 59
 is the size of item for which cumulative frequency is 51

$\therefore \text{Median} = 10$

Q Find the average no. of person per house using median

No. of person in house	1	2	3	4	5	6	7	8	9
No. of house	26	113	120	95	60	42	21	14	54

class	Frequency	Cumulative Frequency
1	26	26
2	113	139
3	120	259
4	95	354
5	60	414
6	42	456
7	21	477
8	14	491
9	54	545
	$N = 545$	

Mode Median = size of $(\frac{N+1}{2})^{\text{th}}$ item = size of $(\frac{546}{2})^{\text{th}}$

item = size of 273th item.

The 1st C.F which includes 273 is 354, 4

$$\therefore \text{Median} = 4$$

Q Find the median from the following data

Marks : 15-25	25-35	35-45	45-55	55-65	65-75	75-85
no. of std : 3	5	12	15	9	9	7

Class	F	Cumulative Frequency
15-25	3	3
25-35	5	8
35-45	12	20
45-55	15	35 $\leftarrow \frac{N}{2}$
55-65	9	44
65-75	9	53
75-85	7	60
	<u>60</u>	

Median = size of $(\frac{N}{2})^{\text{th}}$ item = size of $(\frac{60}{2})^{\text{th}}$ item
 = size of 30th item 30 is included in the C.F
 35. The class having C.F 35 is 45-55

$\therefore 45-55$ is the median class applying interpolation formula.

$$\text{Median} = l_1 + \frac{\left(\frac{N}{2} - CF \right) \times c}{f}$$

$$l_1 = 45, \frac{N}{2} = 30, CF = 20, f = 15, c = 10$$

$$45 + \left(\frac{30 - 20}{15} \right) \times 10$$

$$= \frac{155}{3}$$

$$= \underline{\underline{51.6667}}$$

Q Calculate median from the following data

class : 10-20 20-40 40-70 70-120 120-140

F : 4 10 26 8 2

class	F	CF
10-20	4	4
20-40	10	14
40-70	26	40 $< \frac{N}{2}$
70-120	8	48
120-140	2	50
		50

Median = size of $(\frac{N}{2})^{\text{th}}$ item = size of $(\frac{50}{2})^{\text{th}}$ item
 = 25th item. 25 is included in the C.F 40. The class having C.F 40 is 40-70
 \therefore 40-70 is the median class interpolation formula

$$\text{media} = l_1 + \frac{\left(\frac{N}{2} - C.F\right)}{C} \times C$$

$$l_1 = 40$$

$$= 40 + \frac{(25 - 14)}{30} \times 30$$

$$= \underline{\underline{50.69}}$$

$$F = 26$$

$$C = 30$$

Mode

The value of the variable which occurs most frequently in a distribution is called mode.

Mode in Individual Series :-

In the case of individual series the value which occurs more no. of times is mode

Q Find mode 23 35 28 42 62 53 35 28 42 35 23 42 35

Ans 35 is the mode

NOTE :

when no item appears more no.of times than other we say mode is ill-defined in that case mode is obtain by the formula mode = 3 median - 2 mean

Q Find mode from the values

Ans 40, 25, 60, 35, 81, 75, 90, 10

Here all items appear equal no.of times so mode is ill-defined mode is obtain by the formula 3median - 2 mean

$$\text{mean } \frac{4+6}{2} = 5$$

$$\text{Mean } \frac{4+6}{8} = 5$$

10, 25, 35, 40, 60, 75, 81, 90

$$\text{median} = \frac{n+1}{2} = \frac{8+1}{2} = \frac{9}{2} = 4.5$$

$$\frac{(40+60)}{2} = 50$$

$$\text{mode} = 3\text{median} - 2\text{mean}$$

$$= 3 \times 50 - 2 \times 5 = 46$$

Mode is discrete frequency distribution

In the case of discrete frequency distribution the value having highest frequency is taken as mode.

Find mode

size : 5 8 10 12 19 35 40 46

no of item : 3 12 25 40 31 20 18 7

Ans Mode is 12

The value has the highest Frequency so mode is equal to N

Mode in continuous frequency distribution. In the case of continuous frequency distribution mode lies in the class having the highest frequency this class is known as mode class from the model class mode is calculated by the formula

$$l + \frac{(f_1 - f_0)c}{2f_1 - f_0 - f_2}$$

where l = lower limit of model class

f_0 = frequency of class just preceding model class

f_2 = frequency of the class just succeeding model class

f_1 = frequency of model class

c = Interval of the model class

Q calculate mode from the following data

Size	10-15	15-20	20-25	25-30	30-35	35-40	40-45	45-50
Freq:	4	8	18	30	20	10	3	2

25-30 is the model class then

$$l = 25$$

$$P_0 = 18$$

$$f_1 = 30$$

$$f_2 = 20$$

C = 5

$$\text{mode} = l + \frac{(f_1 - f_0)c}{2f_1 - f_0 - f_2}$$

$$= l + \frac{(f_{r_1} - f_0)c}{2f_{r_1} - f_0 - f_2}$$

$$= 25 + \frac{(30 - 18)5}{}$$

$$2 \times 30 - 18 - 20$$

(sub) Libom prisosejva
= 47.73

↳ lesson progressiv ~~→ mol~~ und -lo- passivum = st

Experiments of modern

Essential Properties (or characteristics) of Good Average

- 1) An average is regarded as a good average if it has the following properties.
- 1) It should be clearly defined.
 - 2) It should be based on all the observations of the data.
 - 3) It should be easy to calculate and simple to follow.
 - 4) It should not be influenced by sampling fluctuations.
 - 5) It should be amenable to further algebraic treatment.
 - 6) It should not be affected by extreme items.

MERITS and DEMERITS of measures of Central Tendency

merits of Arithmetic mean

- 1) Arithmetic mean is simple to understand.
- 2) Arithmetic mean can be easily calculated.
- 3) Arithmetic mean can be determined in most of the cases.
- 4) It is based on all the observations of the series.
- 5) It is capable of more algebraic treatment.
- 6) Arithmetic mean is stable. It does not differ from sample to sample.

Demerits of Arithmetic Mean :-

- 1) Arithmetic mean is affected by extreme values.
- 2) Usually mean does not coincide with any of the observed values.
- 3) It is not suitable for averaging ratios and percentages.
- +) It cannot be calculated for qualitative data which cannot be measured numerically.
- 5) It may offer misleading and absurd results in some cases. For example average number of children per family in a village may work out as 3.2 which is impossible.
- 6) In the case of Frequency distributions with open end classes, the mid values of all classes cannot be obtained. Therefore in such distributions, mean cannot be calculated accurately.

Merits of Median :-

- 1) It is very simple measure.
- 2) Sometimes it can be located even by a mere glance.
- 3) It is not affected by extreme items.
- 4) It is suitable for even such data which are not capable of numerical expression.
- 5) It can be determined graphically also.
- 6) Median is the suitable average in the case of

Open end class Series,

Demerits of Median :-

- 1) Median is not based on all the observations.
- 2) It is not capable of algebraic treatments.
- 3) It requires arranging (ie writing in the ascending or descending order) also.
- 4) In the case of continuous Series interpolation formula is to be used. The value thus obtained may give only an approximate value.

Merits of Mode :-

- 1) Mode is a very simple among measures of central tendency. Even a glance at the series is enough to locate the modal value. It's a popular average.
- 2) Mode is less affected by extreme values in the series.
- 3) Mode can be located graphically also.
- 4) Usually mode coincides with one of the values of the series.
- 5) Mode is the value with which occurs most frequently. So mode is the best representative.

~~Elements of Mode :-~~

Q Why is Arithmetic mean considered to be the best average?

- Arithmetic mean can be easily calculated. It is simple to understand. It is well defined. It can be determined in most of the cases. It is Capable of more algebraic treatment. It is based on all observations of the data. It can be located even without arrangement of the data. It is stable as it does not differ from Sample to sample when the sample selected is sufficiently large. Thus mean satisfies many of the properties of a good average. That is why Arithmetic mean is considered to be the best among commonly used averages. It has become an average of every day life. It is used in the study of many social and economic problems.

Compare Mean, Median and Mode :-

Mean is a mathematical average while Median and mode are positional averages. Mean is based on all observations while median is the middle item and mode is commonly found item. Mean is capable of further mathematical treatment while median and mode are not capable of more mathematical treatment. Median and mode are values found in the series while mean is not necessarily be such value. Median and mode can be graphically located.

But mean cannot generally be located graphically.
Median and Mean can be obtained by mere inspection
while mean can be obtained by calculation. Mean is
affected by extreme items while median and mode are
not much affected by extreme items. Mode is ill-defined
in some cases but median and mean are well defined
in all cases.

Demerits of Mode :-

- 1) mode is uncertain and vague measure of central tendency. It is illdefined in some cases.
- 2) mode is not capable of further algebraic treatment
- 3) sometimes grouping becomes necessary to identify the mode value.
- 4) mode is not based on all the values of the series
- 5) when the extreme value appears most frequently the extreme value becomes mode. For eg:- If zero appears more number of times then mode is zero. This is not a representative value.

$$\text{Mode} = \frac{\sum f_i M_i}{\sum f_i}$$

Geometric Mean :-

If there are n values in a series then their G.M is defined as the n th root of the product of those n values G.M is a mathematical average.

Geometric Mean in Individual Series :-

If x_1, x_2, \dots, x_n are n values then

$$G.M = \sqrt[n]{x_1, x_2, \dots, x_n}$$

$$\log n^{\frac{1}{n}} = \frac{1}{n} \log n$$

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

$$\sqrt[n]{n} = n^{\frac{1}{n}}$$

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

$$\sqrt[3]{n} = n^{\frac{1}{3}}$$

$$\sqrt[4]{n} = n^{\frac{1}{4}}$$

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

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

$$\text{or } \log G.M = \frac{\log x_1 + \log x_2 + \dots + \log x_n}{n}$$

$$\log G.M = \frac{\sum \log x_i}{n}$$

$$G.M = \text{Antilog} \left[\frac{\sum \log n}{n} \right]$$

Q calculate G.M of the following figure

57.5, 87.75, 53.5, 73.5, 81.75

x	$\log n$
57.5	1.7597
87.75	1.9432
53.5	1.7283
73.5	1.8663
81.75	1.9125
	9.21

$$G.M = \text{Antilog} \left[\frac{\sum \log n}{n} \right]$$

$$= \text{Antilog} \left[\frac{9.21}{5} \right]$$

$$= \text{Antilog} [1.842]$$

$$= 69.51$$

G.M. in Frequency distribution :-

For Discrete Frequency distribution

$$G.M = \text{Antilog} \left[\frac{\sum f \log n}{N} \right]$$

Q Find G.M from the following data

Size : 5 8 10 12

Freq : 2 3 4 1

x	f	$\log x$	$f \log x$
5	2	.6980	1.396
8	3	.9031	2.7093
10	4	1	4
12	1	1.0792	1.0792
$N = 10$			9.1845

$$G.M = \text{Antilog} \left(\frac{\sum f \log x}{N} \right)$$

$$= \text{Antilog} [0.91845]$$

$$G.M = 8.287$$

$$\boxed{\frac{\sum f \log x}{N}} \text{ for } N = 10$$

G.M in Continuous frequency distribution

In continuous frequency distribution mid value of the class is taken as x

Find G.M

$$Ed - 2.5 =$$

Mark	0-10	10-20	20-30	30-40	40-50
no. of stud	5	7	15	25	8

class	freq	midvalue x	log n	Hogn
0-10	5	5	• 6990	3.495
10-20	7	15	1.1761	8.2327
20-30	15	25	1.3980	20.97
30-40	25	35	1.5441	38.6025
40-50	8	45	1.6532	13.2256
	60			89.5258

$$G.M = \text{Antilog} \left[\frac{\sum f \log u}{n} \right]$$

$$= \text{Antilog} \left[\frac{84.5258}{60} \right]$$

$$= \text{Antilog} [1.4088]$$

$$= \underline{\underline{25.63}}$$

Q calculate G.M for the following data

Data	100-104	105-109	110-114	115-119	120-124	125-129
no. of items	29	30	45	65	72	84

F.S.E.S - 8	12 F.I. I	21	6	05 - 01
F.P. - 05	130 - 134	135 - 139	5	21
2500 - 25	124	58	8	02 - 08

22.55 - 81 5.66 - 01 2.2 8 02 - 02

82.32 - 28 6.0 6.0 6.0 6.0

class	frequency	midvalue x	log n	f log n
100 - 104	24	102	2.0086	48.2064
105 - 109	30	107	2.0294	60.882
110 - 114	45	112	2.0492	92.214
115 - 119	65	117	2.0682	134.433
120 - 124	72	122	2.0863	150.2136
125 - 129	84	127	2.1038	176.7192
130 - 134	124	132	2.1206	262.9544
135 - 139	58	137	2.1367	123.9286
N = 502				1049.55

$$G.M = \text{Antilog} \left[\frac{\sum f \log x}{N} \right]$$

$$= \text{Antilog} \left[\frac{1049.55}{502} \right]$$

$$= \text{Antilog } [2.091]$$

$$= \underline{\underline{123.3}}$$

Harmonic Mean (H.M)

Harmonic mean of a set of n values is defined as the reciprocal of the mean of reciprocals of the values.

H.M in Individual Series

Let x_1, x_2, \dots, x_n are n values then

$$H.M = \frac{n}{\frac{1}{x_1} + \frac{1}{x_2} + \dots + \frac{1}{x_n}}$$

i.e.,

$$H.M = \frac{n}{\sum \frac{1}{x_i}}$$

Q Find H.M

2, 3, 4, 5

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

$$\sum \frac{1}{n} = \frac{77}{60} = 1.283$$

$$H.M = \frac{4}{\sum \frac{1}{x_i}} = \frac{4}{1.283} = 3.125$$

Q Find H.M 10, 20, 40, 25

x	$1/x$
10	0.1
20	0.05
40	0.025
25	0.04
$n=4$	0.215

$$H.M = \frac{n}{\sum 1/x} = \frac{4}{0.215} = 18.604$$

H.M in Frequency distribution

In the case of Discrete and continuous Frequency distribution

$$H.M = \frac{N}{\sum f 1/x}$$

Q Find H.M

Size	60	10	14	18
Freq	20	40	30	10

08

n	f	$\frac{f}{n}$	$f \times \frac{f}{n}$
6	20	0.167	3.34
10	40	0.1	4
14	30	0.071	2.13
18	10	0.056	.56
100			10.03

$$H.M = \frac{N}{\sum f \frac{f}{n}} = \frac{100}{10.03} = 9.97$$

Q Find H.M from the following data

Class	10-20	20-30	30-40	40-50	50-60
f	4	6	10	7	3

class	f	midvalue x	$\frac{f}{n}$	$f \times \frac{f}{n}$
10-20	4	15	0.133	0.268
20-30	6	25	0.1	0.24
30-40	10	35	0.067	0.268
40-50	7	45	0.047	0.154
50-60	3	55	0.025	0.059
30				0.996

$$H.M = \frac{N}{\sum f/n} = \frac{30}{0.996} = 29.88$$

Q Find H.M of

$x : 2 \ 4 \ 8 \ 10 \ 12$

$f : 1 \ 3 \ 5 \ 4 \ 2$

x	f	f/n	f^2/n
2	1	.5	.5
4	3	.25	.75
8	5	.125	.625
10	4	.1	.4
12	2	.083	.166
	15		2.441

$$H.M = \frac{N}{\sum f/n} = \frac{15}{2.441} = 6.145$$

Q Find H.M

Size	0-10	10-20	20-30	30-40	40-50
Freq	2	3	7	5	2

class	freq	midvalue x	$\frac{1}{m}$	$f \cdot \frac{1}{m}$
0-10	2	5	0.2	0.4
10-20	3	15	0.067	0.201
20-30	7	25	0.04	0.28
30-40	5	35	0.028	0.14
40-50	2	45	0.022	0.044
	19		1.065	

$$H.M = \frac{N}{\sum f \cdot \frac{1}{m}} = \frac{19}{1.065} = 17.84$$

Measure of Dispersion

Dispersion refers to the variability in the size of items

Measure of dispersion are classified into:

- 1) Absolute measures
- 2) Relative measures

Absolute Measures of dispersion are

- 1) Range
- 2) Quartile Deviation (QD)
- 3) Mean Deviation (MD)
- 4) Standard Deviation (SD)

The various relative measures of dispersion are

- 1) Coefficient of Range
- 2) Coefficient of quartile deviation
- 3) Coefficient of standard deviation
- 4) Coefficient of variation

Range

- It is the simplest possible measure of dispersion.
- It is the difference between highest & lowest value in a series.

$$\text{Range} = H - L$$

$$\text{Coefficient of Range} = \frac{H-L}{H+L}$$

Q Find range and coefficient of Range for the following values

25, 32, 85, 32, 42, 10, 20, 18, 28

- Here $H = 85$ $L = 10$

$$\text{Range} = H - L = 75$$

$$\text{Co of Range} = \frac{H-L}{H+L} = \frac{75}{95} = 0.789$$

Q Find Range

wage	10	15	18	20	25
no. of employee	3	5	12	8	6

- $H = 25$ $L = 10$

$$\text{Range} = H - L = 15$$

$$\text{Co of Range} = \frac{H-L}{H+L} = \frac{15}{35} = 0.428$$

Q Five students obtained the following marks in statistics 20, 35, 25, 30, 15. Find out Range and coefficient of Range of the marks

- Here $H = 35$ & $L = 15$

$$\text{Range} = 35 - 15 = 20$$

$$\text{C. of Range} = \frac{H-L}{H+L} = \frac{20}{50} = \underline{\underline{0.4}}$$

Q Find out Range and coefficient of Range of the following series

Marks : 20-29 30-39 40-49 50-59 60-69

No. of Stud : 8 12 20 7 3

[Hint : Read the problem as exclusive series]

- we have to make series exclusive
∴ It is inclusive

1st find interval b/w upper limit of a class and lower limit of next class and find its average and subtract from every lower limit and add to every upper limit

Marks	19.5-29.5	29.5-39.5	39.5-49.5	49.5-59.5	59.5-69
No. of std	8	12	20	7	3

$$H = 69.5 \quad L = 19.5$$

$$\text{Range} = H - L = 50$$

$$\text{C. of Range} = \frac{H-L}{H+L} = \frac{50}{89} = 0.562$$

Quartile Deviation (Semi Interquartile Range)

It is defined as half the distance between 3rd and 1st quartiles.

$$\text{i.e. } Q.D = \frac{Q_3 - Q_1}{2}$$

$$\text{Coefficient of Q.D} = \frac{Q_3 - Q_1}{Q_3 + Q_1}$$

Q.D in Individual Series

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

$$Q_3 = \text{size of } \left(3 \times \frac{n+1}{4} \right)^{\text{th}} \text{ item}$$

Q Compute Q.D and coefficient of Q.D

23, 25, 8, 10, 9, 29, 45, 85, 10, 16, 24

Arrange Values in Ascending Order

8, 9, 10, 10, 10, 16, 23, 24, 25, 29, 45, 85

$$Q_1 = \left(\frac{n+1}{4} \right)^{\text{th}} \text{ item} = \left(\frac{11+1}{4} \right) = 3^{\text{rd}} \text{ item} = 10$$

$$Q_3 = \frac{3 \times (n+1)}{4} \text{ item} = \frac{3 \times 11}{4} = 9^{\text{th}} \text{ item} = 29.5$$

$$\text{Coefficient of Q.D} = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{29.5 - 10}{29.5 + 10} = \frac{19.5}{39.5} = 0.48$$

Q Find Q.D and coefficient Q.D

170, 82, 110, 100, 150, 120, 200, 116, 250

Ascending order

82, 100, 110, 116, 120, 150, 170, 200, 250

$$Q_1 = \left(\frac{n+1}{4} \right)^{\text{th}} \text{ item} = \left(\frac{9+1}{4} \right)^{\text{th}} = \left(\frac{10}{4} \right)^{\text{th}} = (2.5)^{\text{th}} \text{ item}$$

$$2^{\text{th}} \text{ item} + 0.5 \times (3^{\text{rd}} \text{ item} - 2^{\text{th}} \text{ item}) = 100 + 0.5(110 - 100) = 105$$

$$Q_3 = \frac{3 \times (n+1)}{4}^{\text{th}} \text{ item} = \left(\frac{3 \times 10}{4} \right)^{\text{th}} \text{ item} = \left(\frac{30}{4} \right)^{\text{th}} = (7.5)^{\text{th}} \text{ item}$$

$$7^{\text{th}} \text{ item} + 0.5 \times (8^{\text{th}} \text{ item} - 7^{\text{th}} \text{ item}) = 170 + 0.5(200 - 170) = 185$$

$$Q_3 = \underline{\underline{185}}$$

$$Q.D = \frac{Q_3 - Q_1}{2} = \frac{185 - 105}{2} = \frac{80}{2} = 40$$

$$\text{Coefficient of } Q.D = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{185 - 105}{185 + 105} = \frac{80}{290} = 0.275$$

Q.D in Discrete Frequency distribution

In the case of Discrete Frequency distribution

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

Q_3 = size of $\left(\frac{N+1}{4} \times 3\right)^{\text{th}}$ item, where $N=27$

This formula gives in which the quartiles lie

Q Find Q.D and coefficient of Q.D

Size : 5 8 10 12 19 20 32

Frequency : 3 10 15 20 8 7 6

Size	Frequency	Cumulative Frequency
5	3	3
8	10	13
10	15	28
12	20	48

19

8

56

20

7

63

32

6

69

 $N = 69$

$$\frac{N+1}{4} = 17.5$$

$$Q_1 = \left(\frac{N+1}{4} \right)^{\text{th}} = \left(\frac{70}{4} \right)^{\text{th}} = 17.5 \Rightarrow 10$$

$$Q_3 = \left(\frac{N+1}{4} \times 3 \right)^{\text{th}} = 3 \times 17.5 = 52.5 \Rightarrow 19$$

$$Q.D = \frac{Q_3 - Q_1}{2} = \frac{19 - 10}{2} = \frac{9}{2} = 4.5$$

$$\text{Coef of Q.D} = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{9}{29} = 0.3103$$

The specified value of the quartiles are obtained by the interpolation formula

$$Q_1 = l_1 + \frac{\left(\frac{N}{4} - cf \right)}{f} \times c, \quad Q_3 = l_1 + \frac{\left(\frac{3N}{4} - cf \right)}{f} \times c$$

where l_1 is the lower limit of the quartile class

f is the frequency of the class

CF is the cumulative frequency of the preceding class

c is the class interval of the quartile class

Q obtain the quartile measure measure of dispersion and its coefficient for the data given below

Age	0-10	10-20	20-30	30-40	40-50	50-60	60-70
No. of person	15	30	53	75	100	110	115

$$\frac{2}{2} = \frac{10 - 5}{5} = 0.2$$

$$Q_1 = \frac{10 - 5}{5} = 0.5 \text{ or } 30$$

class	frequency	cf
0-10	15	15
10-20	30	45
20-30	53	98
30-40	75	173
40-50	100	273
50-60	110	383
60-70	115	498
70-80	125	623
	623	

$$Q_1 = \text{size of } (N/4)^{\text{th}} \text{ item} = \left(\frac{623}{4}\right)^{\text{th}} \text{ item} = 155.75$$

$$Q_3 = \text{size of } (3N/4)^{\text{th}} \text{ item} = \left(\frac{3 \times 623}{4}\right)^{\text{th}} \text{ item} = 467.25$$

interpolation formula

$$Q_1 = l_1 + \frac{\left(\frac{N}{4} - cf\right)}{f} \times c$$

$$= 30 + \frac{(155.75 - 98)}{75} \times 10$$

$$= 30 + \frac{57.75}{75}$$

$$= 37.7$$

$$Q_3 = l_1 + \frac{\left(\frac{3N}{4} - cf\right)}{f} \times \frac{c}{N}$$

$$= 60 + \frac{(467.25 - 383)}{115} \times 10$$

$$= 60 + \frac{84.25}{115}$$

$$= 67.33$$

$$Q.D = \frac{Q_3 - Q_1}{2} = \frac{67.33 - 37.7}{2} = 14.815$$

$$\text{Coef of Q.D} = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{67.33 - 37.7}{67.33 + 37.7} = 0.28$$

Mean Deviation (M.D.)

It is defined as the arithmetic mean of deviation of all the values in a series from their average counting all such deviations as positive.

The average may be mean, median or mode

$$\text{The coefficient of M.D} = \frac{\text{M.D}}{\text{Average from which M.D is calculated}}$$

M.D in Individual Series

If $|d|$ represent deviation from an average then $\text{M.D} = \frac{\sum |d|}{n}$

Q Find M.D from mean and its coefficient from the following data

25, 63, 85, 75, 62, 70, 83, 28, 30, 12

x	$ d $ $ x - 53.3 $
25	28.3
63	9.7
85	31.7
75	21.7

62	8.7		
70	16.7		
83	29.7		
28	25.3		
30	23.3		
12	41.3		
Σx	236.4		

$$\text{Mean} = \frac{\sum x}{n} = \frac{533}{10} = 53.3$$

$$M.D = \frac{\sum |d|}{n} = \frac{236.4}{10} = \underline{\underline{23.64}}$$

$$\text{Coefficient of M.D} = \frac{M.D}{\text{mean}} = \frac{23.64}{53.3} = 0.443$$

M.D is Discrete Frequency distribution

In the case of discrete frequency distribution

$$M.D = \frac{\sum f|d|}{N}$$

Q Compute the mean deviation about mean and its coefficient for the following data

No. of children : 0 1 2 3 4 4 5 6

No. of families : 171 82 50 25 13 7 2

PP, AP, F8, E8, SP, PP, ZE, EE, G, ZE, 2

X	F	FX	d (x - 1.02)	fd
0	171	0	1.02	174.42
1	82	82	0.02	1.64
2	50	100	1.98	49
3	25	75	1.98	49.5
4	13	52	2.98	58.74
5	7	35	3.98	27.86
6	2	12	4.98	9.96
	350	356	0.02	351.12

$$\text{mean} = \frac{\sum Fx}{N} = \frac{356}{350} = 1.02$$

$$M.D = \frac{\sum |fd|}{N} = \frac{351.12}{350} = 1.0032$$

Q calculate mean deviation from median and its coefficient for the following values.

Values: 5, 28, 33, 44, 83, 87, 96, 99, 25, 35, 82

Asc order 5, 25, 28, 33, 35, 44, 82, 83, 87, 96, 99

Asc order

X	$ d (x - 44)$
5	39
25	19
28	16
33	9
35	0
44	0
52	38
53	39
57	43
59	52
61	55
	321

$$\text{Median} = \left(\frac{n+1}{2} \right)^{\text{th}} \text{ item} = \left(\frac{11+1}{2} \right)^{\text{th}} \text{ item} = 6^{\text{th}} \text{ item} = 44$$

$$M.D = \frac{\sum |d|}{n}$$

$$= \frac{321}{11} = 29.18$$

$$\text{Coefficient of M.D} = \frac{M.D}{\text{median}} = \frac{29.18}{44} = 0.66$$

M.D is continuous frequency distribution

In continuous Frequency distribution =

$$M.D = \frac{\sum f |d|}{N}$$

here midvalue of the class is chosen as m

Q calculate mean deviation about medium for the following data and coefficient of mean deviation

Class	0-10	10-20	20-30	30-40	40-50	50-60	60-70	70-80
Freq	18	16	15	12	10	5	2	2

class	F	cf	midvalue	lal	$\sum f d $
0-10	18	18	5	19	342
10-20	16	34	15	9	144
20-30	15	49	25	1	15
30-40	12	61	35	11	132
40-50	10	71	45	21	210
50-60	5	76	55	31	155
60-70	2	78	65	41	82
70-80	2	80	75	51	102
	80				

$$\text{M.D} = \frac{\sum f |d|}{N} = \frac{1182}{80} = 14.75$$

$$\text{Median} = \left(l_1 + \frac{(N_2 - cf) \times c}{f} \right)$$

$$l_1 = 20$$

$$N_2 = 40$$

$$cf = 34$$

$$c = 15$$

$$C = 10$$

$$= 20 + \frac{40 - 34}{15} \times 10$$

$$= 20 + \frac{6}{15} \times 10$$

$$= 20 + \frac{60}{15}$$

$$= \underline{\underline{24}}$$

$$M.D = \frac{\sum f_i d_i}{N}$$

$$= \frac{1182}{80} = \frac{1182}{N} = M.D/M$$

$$= 14.775$$

$$P.E = \frac{O.P.E}{M.D} =$$

$$\text{Coefficient of } M.D = \frac{M.D}{\text{median}} = \frac{14.775}{24} = \underline{\underline{0.615}}$$

Q calculate M.D about mean of the following data

$$21, 29, 35, 40, 42, 73, 50, 30, 18, 80 = 10$$

X	(d)	$\sum d^2 = 182$
21	18	$21 - 9$
29	10	$01 = 1$
35	4	
40	$\frac{4+2-0+1}{10} = 0.5$	
42	3	
73	$\frac{3}{10} = 0.3$	
50	11	$00 + 0.5 = 0.5$
30	9	21
18	21	$18 = 0.5$
80	41	
390	182	$\frac{182}{10} = 18.2$

$$\text{Mean} = \frac{\sum X}{n}$$

$$= \frac{390}{10} = 39$$

$$\text{M.D} = \frac{\sum |d|}{n}$$

Q Find mean deviation about the median

SI. NO	1	2	3	4	5	6	7	8	9	10
marks	50	27	38	40	90	42	70	60	15	40

[Hint This is an individual series]

- 15, 27, 38, 40, 40, 42, 50, 60, 70, 90

X	(d ₁) (x - M)
15	26
27	14
38	3
40	1
40	1
42	1
50	9
60	19
70	29
90	49
	152

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

$$\left(\frac{11}{2}\right)^{\text{th}} \text{ item} = (5.5)^{\text{th}} \text{ item}$$

$$= \frac{40+42}{2} = \frac{82}{2}$$

$$\text{median} = \underline{\underline{41}}$$

$$M.D = \frac{\sum |d_1|}{n}$$

$$= \frac{152}{10}$$

$$= \underline{\underline{15.2}}$$

Q calculate M.D about mean from the following data

X :	10	11	12	13	14
F :	3	12	18	12	3

[using formula no. 2 in text]

X	F	FX	d	F d
10	3	30	2	6
11	12	132	1	12
12	18	216	0	0
13	12	156	1	12
14	3	42	2	6
	48	576		36

$$\text{Mean} = \frac{\sum Fx}{N}$$

$$= \frac{576}{48} = 12$$

$$M.D = \frac{\sum F|d|}{N} = \frac{36}{48} = 0.75$$

Q Compute mean deviation about mean from the following frequency distribution

class : 0-10 10-20 20-30 30-40 40-50

freq : 3 8 12 9 8

class	freq	midvalue x	Fx	$ d $ $x - 27.75$	$f d $
0-10	3	5	15	22.75	68.25
10-20	8	15	120	12.75	102
20-30	12	25	300	2.75	33
30-40	9	35	315	7.25	65.25
40-50	8	45	360	17.25	138
Σf	40		1110		406.5

$$\text{Mean} = \frac{\sum Fx}{N}$$

$$= \frac{1110}{40}$$

$$= 27.75$$

$$M.D = \frac{\sum f|d|}{N}$$

$$= \frac{406.5}{40} = 10.16$$

Q calculate M.D from mean and from median
 for the following distribution of the scores of
 50 college student is

Scores : 140 150 160 170 180 : 190

Freq : 4 6 10 18 9 3

[Hint : Read close as 140, 150, 160, 160 ...]

M.D about mean

class	F	midvalue x	$\sum Fx$	$\sum fd$ $x - 171.2$	$\sum fd^2$
140-150	4	145	580	26.2	104.8
150-160	6	155	930	16.2	97.2
160-170	10	165	1650	6.2	62
170-180	18	175	3150	3.8	68.4
180-190	9	185	1665	13.8	124.2
190-200	3	195	585	23.8	71.4
	50		8560		528

$$\text{Mean} = \frac{\sum Fx}{N} = \frac{8560}{50} = 171.2$$

$$M.D = \frac{\sum f|d|}{N}$$

$$= \frac{528}{50} = \underline{\underline{10.56}} + 0.51 =$$

M.D about Mean

class	f	cf	midvalue n	d $x - 171.78$	$f d $
140-150	4	4	145	27.78	111.12
150-160	6	10	155	17.78	106.68
160-170	10	20	165	7.78	77.8
N/3 170-180	18	38	175	2.22	39.96
180-190	9	47	185	12.22	109.98
190-200	3	50	195	22.22	66.66
		50			512.2

$$\text{Median} = l_1 + \frac{N/2 - cf}{f} \times c$$

$$l_1 = 170$$

$$N/2 = 25$$

$$cf = 20$$

$$f = 18$$

$$c = 10$$

$$\text{Median} = 170 + \frac{25 - 20}{18} \times 10 = 172.78$$

$$M.D = \frac{\sum F | d |}{N}$$

$$= \frac{512 - 2}{50}$$

$$= 10.24$$

$$O.F.I = \frac{f_1 - M}{M} \times 100$$

$$O.F.I = 0$$

$$Z.S = \frac{M - \bar{x}}{S}$$

$$Z.S = 4.5$$

$$Z.T = 3$$

$$O.T = 0$$

Standard Deviation

It is the square root of mean of squares of the deviation of all values of a series from their arithmetic mean

S.D in individual Series

$$S.D = \sqrt{\frac{\sum x^2}{n} - \left(\frac{\sum x}{n}\right)^2}$$

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

$$\text{Coefficient of variation } CV = \frac{SD}{\text{mean}} \times 100$$

Q For the following values find S.D, variance coefficient of variance

5, 8, 7, 11, 9, 10, 8, 2, 4, 6

$$\sqrt{\left(\frac{\sum x^2}{n}\right) - \left(\frac{\sum x}{n}\right)^2}$$

X	x^2
5	25
8	64
7	49
11	121
9	81
10	100
8	64
2	4
4	16
6	36
70	560

$$S.D = \sqrt{\frac{\sum x^2}{n} - \left(\frac{\sum x}{n}\right)^2}$$

$$\text{Mark is } 560 \quad S.D = \sqrt{\frac{560}{10} - \left(\frac{70}{10}\right)^2}$$

$$= \sqrt{56 - 49}$$

$$= \sqrt{7} \quad \text{Ansatz der Varianz in Q. 2}$$

$$= 2.64$$

$$V \rightarrow o \leftarrow \left(\frac{\sum x^2}{n} - \left(\frac{\sum x}{n} \right)^2 \right) = (Q. 2)$$

$$= 7$$

$$CV = \frac{S.D}{\text{mean}} \times 100$$

$$= \frac{2.64}{7} \times 100$$

$$= 37.71$$

S.D in continuous and Discrete frequency distribution

$$S.D = \sqrt{\frac{\sum Fx^2}{N} - \left(\frac{\sum Fx}{n}\right)^2}$$

Q For the following data calculate S.D variance and C.V

Mark 2 4 6 8 10

no. of std 8 10 6 9 7

X	F	FX	X^2	FX^2
2	8	16	4	32
4	10	40	16	160
6	16	96	36	576
8	9	72	64	576
10	7	70	100	700
	50	294		2044

$$S.D = \sqrt{\frac{\sum FX^2}{N} - \left(\frac{\sum FX}{N}\right)^2}$$

$$= \sqrt{\frac{2044}{50} - \left(\frac{294}{50}\right)^2}$$

$$\bar{x} = \sqrt{40.88 - 35.28}$$

$$= \sqrt{5.6}$$

$$= 2.36$$

$$V = \frac{\sum Fx^2}{N} - \left(\frac{\sum Fx}{N} \right)^2$$

$$= 5.6$$

$$CV = \frac{S.D}{\text{mean}} \times 100$$

$$= \frac{2.36}{35.88} \times 100$$

$$= \frac{2.36}{5.88}$$

$$= 40.14$$

Q calculate SD, Variance & CV for the following data

class : 0-2 2-4 4-6 6-8 8-10 10-12

Freq : 2 4 6 4 2 6

class	freq	midvalue	Fx	x^2	Fx^2
0-2	2	1	2	1	2
2-4	4	3	12	9	36
4-6	6	5	30	25	150
6-8	4	7	28	49	196
8-10	2	9	18	81	162
10-12	6	11	66	121	726
	24		156		1272

$$S.D = \sqrt{\frac{\sum Fx^2}{N} - \left(\frac{\sum Fx}{N}\right)^2}$$

$$S.D = \sqrt{\frac{1272}{24} - \left(\frac{156}{24}\right)^2}$$

$$= \sqrt{53 - 42.25}$$

$$= \sqrt{10.75}$$

$$= \underline{3.28}$$

$$V = \frac{\sum Fx^2}{N} - \left(\frac{\sum Fx}{N}\right)^2 = \underline{10.75}$$

$$C.V = \frac{S.D}{\text{mean}} \times 100$$

$$\therefore C.V = \frac{3.28}{3.5} \times 100 = \underline{50.46}$$

Merits and Demerits of Measure of Dispersion

Merits of Range

- 1) Range is the simplest measure of dispersion.
- 2) It can be easily calculated.
- 3) It can be understood even by a layman.

Demerits of Range

- 1) Range is not based on all items of the series
- 2) It is highly affected by Sampling Fluctuations
- 3) It cannot be computed on the case of open end distribution

Uses of Range

Range is used in certain fields to measure variability particularly those data where variation is not much. For example, doctors are generally interested in the range of the fluctuating temperatures of patients. It is also used in Quality Control. Range is used to study variation in prices of shares and interest rates. In weather forecasts, the minimum and the maximum temperature of every day are studied from this, they can forecast the range within which the temperature may likely vary.

Merits of Quantile Deviation

- 1) It is easy to understand and to calculate.
- 2) It is not affected by extreme values.

Demerits of Quantile Deviation

- 1) It ignores extreme items.
- 2) It is not capable of more algebraic treatments.

Merits of Mean Deviation

- 1) Mean Deviation is a very simple and an easy measure of dispersion. It is easily understood.
- 2) It is based on all the items of the series. So it is more representative.
- 3) Mean deviation is less affected by extreme values.

Demerits of Mean Deviation

- 1) Mean Deviation suffers from inaccuracy because '+' or '-' signs are ignored.
- 2) Mean deviation is not capable of any further algebraic treatment.

3) Mean deviation is not reliable measure when calculated from mode as the mode is uncertain in some cases.

Uses of Mean deviation

Mean Deviation is significantly used for measuring variability of the series relating to Economic and social phenomena. Variability in the distribution of wealth and income is generally measured in terms of mean deviation.

Merits of Standard Deviation

- 1) Standard deviation is based on all the values of a series it does not ignore any value.
- 2) Standard deviation is a clear and definite measure of dispersion so that it can be measured from all series.
- 3) It is not very much affected by Sampling fluctuators.

4) It is capable of further algebraic treatment.

Demerits of Standard Deviation

- 1) Standard is difficult to calculate
- 2) It gives more weight to extreme values

Features of Standard Deviation

- 1) Deviations are measured from Arithmetic mean
- 2) Signs of the deviations are not ignored.

Q Why is Standard Deviation considered to be the best measure of dispersion?

- Standard Deviations possess most of the important characteristics which an ideal measure of dispersion should have

- I Standard Deviation is rigidly defined
- II It is based on all the observations of the data
- III It is amenable to more algebraic treatment
- IV It possesses many mathematical properties
- V It is not much affected by sampling fluctuations
- VI It does not ignore the signs of the deviations
- VII It is possible to find out S.D of two or more groups

VIII coefficient of variation is based on standard deviation and it can be used to compare variability of 2 series.

IX Standard Deviation is used to find out statistical measures like coefficient of skewness, coefficient of correlation Regression equations etc...
So Standard Deviation is best measure of dispersion

Difference between Standard Deviation and mean Deviation

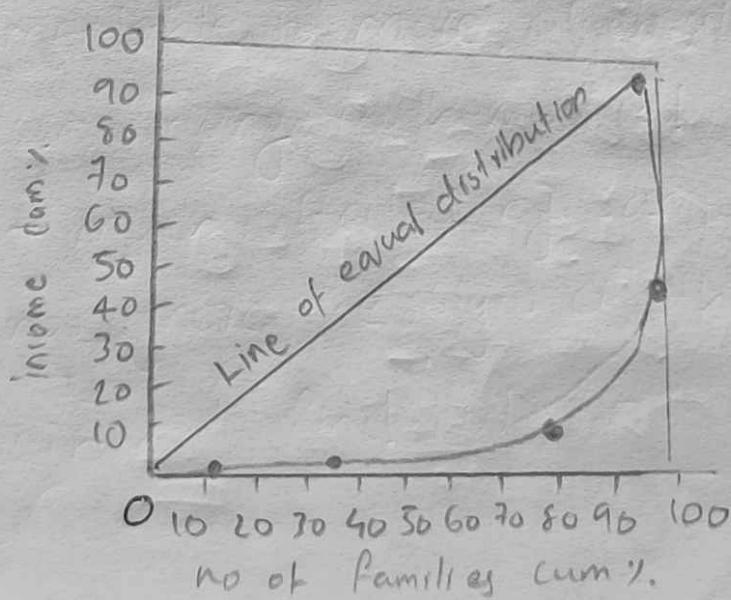
- 1) In standard Deviation, deviations are taken only from the mean of the values of the series. But in mean Deviation, deviations may be taken from mean or median or mode.
- 2) In standard deviation, signs of deviations are not ignored. But in mean deviation, signs of deviations (+ or -) are ignored.

Lorenz Curve

Lorenz curve is a graphic method of studying dispersion in a series. It is used in business to study the disparities of the distribution of wages, turnover, production, population etc... In Economics it is useful to measure inequalities in the distribution of income between different countries or between different periods of time.

Lorenz curve is a graph drawn to frequency distribution, taking the cumulated percentage values of frequencies along the x-axis and cumulated percentage values of the variable along the Y-axis.

Given below is an example of Lorenz curve



Q How to measure variability (or inequality) from a Lorenz presentation

- If there is no inequality in the distribution, the Lorenz curve will coincide with the line of equal distribution. The more the Lorenz curve is away from the line of equal distribution, the greater is the inequality (or variability).
Lorenz curve is useful to

1) study the variability in a distribution.

2) compare the variability relating to a phenomenon for two regions.

3) study the changes in variability over a period.

unit IV

curve fitting

Fitting of a straight line of the form $y = a + bx$

Suppose we want to find a straight line as a best approximation to the actual curve $y = f(x)$ passing through n given points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ this line will be referred as the line of best fit and we take its equation as $y = a + bx - ①$ where a & b are the parameters to be determined. Values of a and b are obtained by solving two equations known as Normal equations

$$\begin{cases} \sum y = n a + b \sum x \\ \sum xy = a \sum x + b \sum x^2 \end{cases}$$

After determine the values of a & b put in these values in ① we get the equation of the line of best fit for the given data.

Q Fit a straight line to the following data

x 1 2 3 4 5

y 14 13 4 5 2

Estimate the value of y at $x = 3.8$

$$n \quad y \quad ny - n^2$$

1

14

14

1

2

13

26

4

$\Sigma n = 15$ $\Sigma y = 38$ $\Sigma ny = 114$ $\Sigma n^2 = 55$

$\Sigma y^2 = 82$ $\Sigma xy = 82$ $\Sigma nxy = 55$

$\therefore (-14, 14) (-13, 26) (-12, 20) (-11, 16)$

$(-10, 5) (-9, 5) (-8, 2) (-7, 10) (-6, 25)$

$\therefore \sum n = 15 \quad \sum y = 38 \quad \sum ny = 114 \quad \sum n^2 = 55$

The normal equation are

$$\sum y = na + b \sum n \quad \text{or} \quad \sum y = a \sum n + b \sum n^2$$

$$\sum ny = a \sum n + b \sum n^2$$

$$\text{i.e., } \sum y = na + b \sum n \quad \text{or} \quad \sum y = a \sum n + b \sum n^2$$

$$38 = 15a + 13b$$

$$82 = 15a + 55b$$

$$15a + 13b = 38 \quad \text{--- (1)}$$

$$15a + 55b = 82 \quad \text{--- (2)}$$

$$3 \times (1) \Rightarrow 45a + 39b = 114 \quad \text{--- (3)}$$

$$(2) - (3) \Rightarrow 16b = -32$$

$$b = \frac{-32}{16} = -2$$

$$8.2 = a + b \quad \text{or} \quad a = 8.2 - b$$

put $b = -3.2$ in eqn ①

$$5a + 15x - 3.2 = 38$$

$$5a - 48 = 38$$

$$5a = 38 + 48$$

$$5a = 86$$

$$a = \frac{86}{5}$$

$$= 17.2$$

$$y = a + bx$$

$$y = 17.2 - 3.2x$$

$$\text{at } x = 3.5$$

$$y = 17.2 - 3.2 \times 3.5$$

$$y = 17.2 - 10.2$$

$$y = 6$$

Q Fit a straight line to the following data

$$x \quad 1 \quad 2 \quad 3 \quad 4 \quad 6 \quad 8$$

$$y \quad 2.4 \quad 3 \quad 3.6 \quad 4 \quad 5 \quad 6$$

$$\sum f_i = 14 \quad \sum x_i = 28$$

$$\frac{\sum f_i}{6} = 2$$

$$\frac{\sum x_i}{6} = 4.67$$

x	y	xy	x^2	as per given standard form
1	2.4	2.4	1	$8E = 8 \cdot E - xE + nE$
2	3	6	4	$8E = xE - nE$
3	3.6	10.8	9	$8E + 8E = nE$
4	4	16	16	$32E = nE$
6	5	30	36	$180E = nE$
8	6	48	64	$24E = nE$

$$\sum x = 24 \quad \sum y = 24 \quad \sum xy = 113.2 \quad \sum x^2 = 130$$

The normal equation are

$$\sum y = na + b \sum x$$

$$\sum xy = a \sum x + b \sum x^2$$

i.e,

$$24 = 6a + b \cdot 24$$

$$113.2 = a \cdot 24 + b \cdot 130$$

~~$$6a + 24b = 24 \quad \text{--- ①}$$~~

~~$$24a + 130b = 113.2 \quad \text{--- ②}$$~~

$$1 \times ① \Rightarrow 24a + 96b = 96 \quad \text{--- ③}$$

$$② - ③ \Rightarrow 34b = 17.2$$

$$b = \frac{17.2}{34}$$

$$= 0.50$$

put $b = 0.50$ in eqn ①

$$6a + 24 \times 0.50 = 24$$

$$6a + 12 = 24$$

$$6a = 24 - 12 \rightarrow \text{solution of } ①$$

$$6a = 12$$

$$a = 12/6$$

$$a = 2$$

$$y = a + bx$$

$$y = 2 + 0.50x$$

Correlation And Regression

Correlation

It is a statistical measure for finding out degree of association between two or more variables.

→ Suppose X & Y are two variables if the movement of X & Y be in the same direction ie, either both X & Y increase or both decrease

Then we say that X & Y are in positive correlation

→ If the movement are in the opposite direction ie, if X increase then Y decrease or viceversa then we say that there is a negative correlation between X & Y .

→ If y is unaffected by any change in x then y are said to be uncorrelated or zero correlation.

Determination of correlation

Correlation between two variables can be determined by the following methods

1 → Scatter diagram

2 → Karl Pearson's method

3 → Rank method

1) Scatter Diagram

The existence of correlation can be shown graphically by means of a Scatter diagram

The data relating to two variables can be graphically represented by points taking one of the variables x along horizontal axis and second variable y along vertical axis. All the pair of value x and y are shown by points on the graph paper. This diagrammatic representation of bivariate data is known as Scatter diagram.

Y & X measured

Wiederholung der Verteilung von Na^+ und K^+ im Zellinneren

bottom ③ now $r=1$ & j bottom ④ $r=-1$ last. S

In Figure ① & ③ the movement of the 2 variable are in the same direction in these the correlation +ve or direct.

In Figure ② & ④ the movement of the variable are in the opposite direction so. the correlation is -ve or direct.

In Figure ⑤ the points lie around a curve in this case the correlation is very small and we take $r=0$.

In figure 3 all points are lie in a straight line. In this case the correlation is perfect. So we can take $r=+1$.

In figure 4 all the points lie on a straight line in this case correlation is perfect and also negative. So we can take $r=-1$.

2. Karl Pearson's method [Co-variance method]

Co-variance

If $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be n pair of observation of two variable x & y. then the co-variance of x and y is denoted as $\text{Cor}(xy)$ and defined as

$$\text{definition} \ L \ \text{Cor}(xy) = \frac{1}{n} \sum (x - \bar{x})(y - \bar{y})$$

Karl Pearson's Co-efficient of Correlation

The correlation coefficient between two variables x and y are denoted by r and it is given by

$$r = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

where σ_x and σ_y are standard deviation of x & y respectively.

The formula can be written as

$$r = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y}$$

$$r = \frac{\frac{1}{n} \sum (x - \bar{x})(y - \bar{y})}{\sqrt{\frac{\sum (x - \bar{x})^2}{n}} \sqrt{\frac{\sum (y - \bar{y})^2}{n}}}$$

take $x = (x - \bar{x})$ and $y = (y - \bar{y})$

$$\text{then } r = \frac{\frac{1}{n} \sum xy}{\sqrt{\frac{\sum x^2}{n}} \sqrt{\frac{\sum y^2}{n}}}$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

Q Find the Correlation co-efficient from the following data

x	1	2	3	4	5	6	7
y	6	8	11	9	12	10	14

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

$$x \quad y \quad x = (x - \bar{x}) \quad y = (y - \bar{y}) \quad x^2 \quad y^2 \quad xy$$

$$1 \quad 6 \quad -3 \quad -4 \quad 9 \quad 16 \quad 12$$

$$2 \quad 8 \quad -2 \quad -2 \quad 4 \quad 4 \quad 4$$

$$3 \quad 11 \quad -1 \quad 1 \quad 1 \quad 1 \quad -1$$

$$4 \quad 9 \quad 0 \quad -1 \quad 0 \quad 1 \quad 0$$

$$5 \quad 12 \quad 1 \quad 2 \quad 1 \quad 4 \quad 2$$

$$6 \quad 10 \quad (\bar{x}-\bar{y})^2 \quad (\bar{x}-\bar{y})^2 \quad 0 \quad 4 \quad 0 \quad 0$$

$$7 \quad 14 \quad \frac{3}{(\bar{x}-\bar{y})^2} \quad \frac{4}{(\bar{x}-\bar{y})^2} \quad 9 \quad 16 \quad 12$$

$$\sum x^2 = 28 \quad \sum y^2 = 42 \quad \sum xy = 29$$

$$\bar{x} = \frac{\sum x}{n} = 4$$

$$\bar{y} = \frac{\sum y}{n} = 10$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

$$= \frac{29}{\sqrt{28} \cdot \sqrt{42}}$$

$$= 0.8456$$

Q calculate the Karl Pearson's correlation between advertisement cost & sales as per the data given below.

cost	39	65	62	90	82	75	25	98	36	78
sales	47	53	58	86	62	68	60	91	51	84

x	y	$x = (x - \bar{x})$	$y = (y - \bar{y})$	x^2	y^2	xy
39	47	-26	-19	676	361	494
65	53	0	-13	0	169	0
62	58	-3	-8	9	64	24
90	86	25	20	625	400	500
82	62	17	-4	289	16	-68
75	68	10	2	100	4	20
25	60	-40	-6	1600	36	240
98	91	33	25	1089	81	2704
36	51	-29	-15	841	625	825
78	84	13	18	169	324	435

$$\bar{x} = \frac{\sum x}{n} \quad \bar{y} = \frac{\sum y}{n} \quad \sum x^2 = 5398 \quad \sum y^2 = 2224 \quad \sum xy = 2704$$

$$= 65 \quad = 66$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \cdot \sqrt{\sum y^2}}$$

$$= \frac{2704}{\sqrt{5398} \sqrt{2224}} = \underline{\underline{0.7804}}$$

Rank Correlation Coefficient

Simple correlation coefficient is based on the magnitude of the variables. But in some situations it is not possible to find magnitude of all the variables.

For example we cannot measure beauty or intelligence qualitatively. But in this case it is possible to rank the intelligent in some order. Rank correlation is based on the order of the rank, Spearman formula for rank coefficient correlation

$$R = 1 - \frac{6 \sum d^2}{n(n^2-1)}$$

where $d \rightarrow$ the difference between the rank.

Q Calculate rank correlation coefficient

mark in
math 78 36 98 25 75 82 90 62 65 69

marks
in
statistic 84 51 91 60 68 62 86 58 53 47

X	Y	$R(X)$	$R(Y)$	d	d^2
78	84	4	3	1	1
36	51	9	9	0	0
98	91	1	1	0	0
25	60	10	(A) 5	6(x) 34	16
75	68	5	4	1	1
82	62	3	5	-2	25
90	86	2	2	0	4
62	58	8	7	1	0
65	53	7	8	-1	1
69	47	6	10	-4	16
28	32	8	5	3	14
1	1	5	8	8	64

$$\sum d^2 = 40$$

$$R = 1 - \frac{6 \sum d^2}{n(n^2 - 1)}$$

$$R = 1 - \frac{6 \times 40}{10(100 - 1)}$$

$$= \frac{1 - 240}{990} = \frac{1 - 24}{110} = \frac{1}{5}$$

$$R = 0.7575$$

Q Find the rank co-efficient correlation

X 28 45 40 35 33 41 32 36 64

Y 23 34 33 36 36 26 28 31 37

X	Y	$R(X)$	$R(Y)$	d	d^2
28	23	9	9	0	0
45	34	1	5	-2	4
40	33	3	4	-1	1
35	36	5	2	3	9
33	30	7	6	1	1
41	26	2	8	-6	36
32	28	8	7	1	1
36	31	4	5	-1	1
34	37	6	7	5	25

$$n \times d - 1 \leq d^2 = 78$$

$$R = \frac{1 - 6 \sum d^2}{n(n^2 - 1)}$$

$$= \frac{1 - 6 \times 78}{9(81 - 1)}$$

$$r_b = 1 - \frac{468}{720}$$

$$= 0.35$$

Spearman's Formula for Repeated Ranks

If in a series 2 or more individual or items having the same score then we find the average of the rank of these individuals and allowed this rank to each of them. In such a case Spearman's modified formula is

$$R = 1 - \frac{6(\sum d^2 + \sum (t^3 - t))}{n(n^2 - 1)}$$

where t is the no. of individuals involved in a tie in the 1st or 2nd series.

Q Find the rank correlation coefficient of the following data

Series A 115 109 112 87 98 120 98 100 98 118

Series B 75 73 85 70 76 82 65 73 68 80

OPP

$$\frac{O-S}{OPP} - 1$$

x	y	R(x)	R(y)	d	d^2
115	75	3	5	-2	4
109	73	5	6.5	-1.5	2.25
112	85	4	1	3	9
87	70	8	2	2	4
98	76	4	4	0	16
120	82	2	2	-1	1
98	65	8	-2	4	16
100	73	6	6.5	-0.5	0.25
98	68	8	9	-1	1
118	80	2	3	-1	1

$$\sum d^2 = 42.5$$

$$R = 1 - \frac{\sum d^2 + \sum (f^3 - f)}{n(n^2 - 1)}$$

$$= 1 - \frac{42.5 + 2.5}{990}$$

$$= 1 - \frac{45}{990} = 0.905$$

$$= 1 - \frac{42.5 + 2.5}{990} = 0.905$$

$$= 1 - \frac{2.5}{990}$$

$$= \frac{1 - 0.272}{(1 - 0.272)^2} = \underline{\underline{0.728}}$$

Q Find the co-efficient of correlation of the marks of some students in subjects

Paper 1 80 45 55 56 58 60 65 68 70 75 85
 Paper 2 82 56 50 48 60 62 64 65 70 64 90

x	y	R(x)	R(y)	d	d^2
80	82	2	2	0	0
45	56	11	9	-2	4
55	50	10	10	0	0
56	48	9	11	-2	4
58	60	8	8	0	0
60	62	7	7	0	0
65	64	6	5.5	0.5	0.25
68	65	5	4	1	1
70	70	4	3	1	1
75	64	3	5.5	-2.5	6.25
85	90	1	1	0	0

$$\sum d^2 = 16.5$$

$$R = \frac{1 - 6 \left(\sum d^2 + \sum (t^3 - t) \right)}{n(n^2 - 1)}$$

$$= \frac{1 - 6(16.5 + 6/12)}{11 \times 10^2}$$

$$\frac{= 1 - G(16.5 + 10.5)}{1320}$$

$$= 1 - \frac{6 \times 17}{1320}$$

$$+ = 0.9227$$

Regression Analysis

Suppose we are given n pairs of values $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ of two variables x and y . If we fit a straight line to this data by taking x as independent variable and y as dependent variable, then the straight line obtained is called regression line of y on x . Its slope is called the regression coefficient of y on x . Similarly if we fit a straight line to the data by taking y as independent variable and x as dependent variable, then the straight line obtained is called regression line of x on y . The reciprocal of its slope is called regression coefficient of x on y .

Equation for regression line

Regression equation of y on x is

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

where b_{yx} is the regression coefficient of y on x .

Regression equation of y on x is

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

where b_{xy} is the regression co-efficient x on y .

Here $b_{xy} = \frac{\text{cov}(x, y)}{\sigma_x^2}$

$$= \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$$

$$b_{xy} = \frac{\sum xy}{\sum x^2}$$

$$b_{xy} = \frac{\text{cov}(x, y)}{\sigma_y^2}$$

$$b_{xy} = \frac{\sum xy}{\sum y^2}$$

$$\bar{x} = 0.55, \bar{y} = 10.15$$

Q You are given the data relating to purchase and sale obtain the regression equation by the method of least squares and estimated the value of the sales when the purchase = 100

Purchase 62 72 98 76 81 56 76 92 88 49

Sales 112 124 131 117 132 96 120 136 97 85

$$(P - \bar{P})(S - \bar{S}) = -$$

x	y	$x = x - \bar{x}$	$y = y - \bar{y}$	x^2	y^2	xy
62	112	-13	-3	169	9	39
72	124	-3	9	9	81	-27
98	131	23	(8.16)	529	256	368
76	117	1	2	1	4	2
81	132	6	17	36	289	102
56	96	-19	-19	361	361	361
76	120	1	5	1	25	5
92	136	17	21	289	441	357
88	97	13	-18	169	324	-234
49	85	-26	-30	676	900	780

$$\bar{x} = 75 \quad \bar{y} = 115 \quad \sum xy = 1753$$

$$\sum y^2 = 2690 \quad \sum n^2 = 2240$$

$$2. P31 = 6$$

$$b_{yx} = \frac{\sum xy}{\sum n^2} = \frac{1753}{2240} = 0.78$$

$$b_{xy} = \frac{\sum xy}{\sum y^2} = \frac{1753}{2690} = 0.651$$

Regression equation of y on n ,

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

$$y - 115 = 0.78(n - 75)$$

$$y - 115 = 0.78n - 58.5$$

$$y = 0.78n - 58.5 + 115$$

$$y = 0.78n + 56.5$$

Regression equation of n on y

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

$$n - 75 = 0.65(y - 115)$$

$$n - 75 = 0.65y - 74.75$$

$$n = 0.65y - 74.75 + 75$$

$$n = 0.65y + 0.25$$

when purchase $n=100$, then sales

$$y = 0.78 \times 100 + 56.5$$

$$y = 134.5$$

Q For the given data obtain the two regression lines

x	8	6	4	7	5
y	9	8	5	6	2

$$x = x - \bar{x} \quad y = y - \bar{y} \quad x^2 \quad y^2 \quad xy$$

$$(n-1) \sum d^2 = \bar{d}^2 - \bar{d}^2$$

$$(2F-x)(2F-y) = 211 - \bar{d}^2$$

$$2 \cdot 32 - x \cdot 2F = 211 - \bar{d}^2$$

$$211 + 2 \cdot 32 - 2F \cdot x = \bar{d}^2$$

$$2.02 + 2F \cdot x = \bar{d}^2$$

for no. of no. of observations

$$(F-d)^2 = \bar{d}^2 - \bar{d}^2$$

$$211 - \bar{d}^2 \cdot 20.0 = 2F - x$$

$$2F \cdot 2F - \bar{d}^2 \cdot 20.0 = 2F - x$$

$$2F + 2F \cdot 2F - \bar{d}^2 \cdot 20.0 = x$$

$$25.0 + \bar{d}^2 \cdot 20.0 = x$$