

Mean Square deviation

Suppose $d = x - a$ is a deviation taken from an arbitrary reference point a . To get rid of algebraic sign of d , we either use $|d|$ or d^2 . The measure of dispersion based on $|d|$ is mean deviation, we have already studied. Using d^2 , we can develop measure of dispersion which is better than mean deviation, clearly.

Mean square deviation $\equiv \frac{1}{n} \sum_{i=1}^n (x_i - a)^2$ for

Frequency distribution of individual observations.

$$= \frac{\sum_{i=1}^n f_i (x_i - a)^2}{\sum_{i=1}^n f_i} \quad \text{for freq^n distribution}$$

However M.S.D. is affected by choice of a . Thus it creates difficulty in measuring dispersion clearly. We try to find a measure which will overcome this difficulty.

Property :-

(i) Sum of square of deviations taken from AM is minimum.

(ii) Minimality property of MSD.

MSD is least if the deviations are taken from AM.

Since, sum of square of deviation taken from

AMI is minimum

$$\text{ie } \sum (x_i - a)^2 \geq \sum (x_i - \bar{x})^2$$

$$\frac{1}{n} \sum (x_i - a)^2 \geq \frac{1}{n} \sum (x_i - \bar{x})^2$$

MSD about $a \geq \text{MSD about } \bar{x}$

Variance, standard deviation and coefficient of variation.

Variance :-

The A.M of squares of deviations taken from mean is called as variance. It is denoted as $\text{var}(x)$ or σ^2 .

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \text{ for individual observation}$$

$$\sigma^2 = \frac{\sum f_i (x_i - \bar{x})^2}{\sum f_i} \text{ for freqn distri.}$$

The units of original items and that of variance are not same.

standard deviation (S.D)

It is positive square root of variance.

It is denoted by σ .

$$\sigma = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2} \text{ for individual observation}$$

$$= \sqrt{\frac{\sum f_i (x_i - \bar{x})^2}{\sum f_i}} \text{ for freqn distribution}$$

Simplified Formula: →

$$\sigma^2 = \frac{1}{n} \sum (x_i - \bar{x})^2$$

$$= \frac{1}{n} \sum x_i^2 - 2\bar{x} \sum x_i + \bar{x}^2$$

$$= \frac{1}{n} \sum x_i^2 - \frac{2\bar{x}}{n} \sum x_i + \bar{x}^2 \cdot \frac{1}{n} \cdot n$$

$$= \frac{1}{n} \sum x_i^2 - 2\bar{x} \cdot \bar{x} + \bar{x}^2$$

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

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

standard deviation is measure of dispersion

which satisfies most of the requisites of good measure. It is free from drawbacks present in the other measure of dispersion.

coeffⁿ of variance (C.V)-

$$C.V = \frac{S.D.}{|A.M|} \times 100 = \frac{\sigma}{|\bar{x}|} \times 100$$

C.V is always expressed in percentage

Remarks:-

R.H.S. of above includes multiplier 100 because σ is too small in many cases. Thus for

1×1

convenience it is multiplied by 100.

(ii) S.D. of weights of a group of elephants will be larger than that of group of human beings.

Suppose S.D. of weights of group of elephants is 15 kg and that of human beings is also 15 kg. In this case we cannot say both the groups have identical variations. This is because average weight of group of elephants is larger than that of average weight of a group of persons. Therefore for comparing variations between different data set, a measure based on ratio of σ & \bar{x} would be appropriate. This is achieved in CV.

* Uses -

- i) for comparing variability of two data sets.
- ii) determining uniformity of data. Smaller CV higher is uniformity.
- iii) Smaller CV data is more homogeneous.
- iv) for comparing consistency of data stability.

Properties:-

(1) Mean square deviation \geq variance

(2) Variance is invariant to the change of origin.
In other words if a constant is added (or subtracted) to each item then variance (SD) remains same.

\rightarrow Let x_1, x_2, \dots, x_n is set of observations.

$y_i = x_i - a$ where a is constant,

we have to s.t $\text{var}(y) = \text{var}(x)$.

As

$$y_i = x_i - a \Rightarrow \bar{y} = \bar{x} - a$$

we have

$$\text{var}(y) = \frac{1}{n} \sum (y_i - \bar{y})^2$$

$$= \frac{1}{n} \sum [x_i - a - (\bar{x} - a)]^2$$

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

$$= \text{var}(x).$$

(3) If $u = \frac{x-a}{h}$, a, h being constant.

$$\text{Then } \text{var}(u) = \frac{1}{h^2} \text{var}(x).$$

$$\rightarrow \text{Let } u_i = \frac{x_i - a}{h} \Rightarrow \bar{u} = \frac{\bar{x} - a}{h}$$

We have

$$\text{var}(u) = \frac{1}{n} \sum (u_i - \bar{u})^2$$

$$= \frac{1}{n} \sum \left[\frac{x_i - a}{h} - \frac{\bar{x} - a}{h} \right]^2$$

$$= \frac{1}{n} \sum \frac{(x_i - \bar{x})^2}{h^2}$$

$$\text{var}(u) = \frac{1}{h^2} \text{var}(x)$$

(4) If $y = ax + b$ Then $\text{var}(y) = a^2 \text{var}(x)$.

Use of variance & SD.

Practically in almost all advanced statistical methods such as sampling, statistical quality control, statistical inference deal with variance. Variance & SD are used in no. of situations. Some of them are as follows:

a) Precision of an instrument is inversely proportional to variance. Therefore

$$\text{precision} = \frac{K}{\text{variance}}$$

b) In portfolio analysis risk is described in terms of variance of prices of shares.

c) For the comparison of performance of two or more instruments, machines, coeff. of variation is used.

* If all observations are equal, SD is ?

* If data contains only one observation, SD is ?

- 7) Compute S.D & CV of marks scored by 10 candidates given below

54, 61, 64, 69, 58, 56, 59, 57, 55, 50

→

$$u = x - 57$$

x_i	54	61	64	69	58	56	59	57	55	50	Total
u_i	-3	4	7	12	1	-1	-8	0	-2	7	3
u_i^2	9	16	49	144	1	1	64	0	4	49	337

$$\sigma = \sqrt{\frac{1}{n} \sum u_i^2 - \left(\frac{\sum u_i}{n}\right)^2}$$

$$= 5.7974$$

To compute CV we need \bar{x}

$$\text{Here } \bar{u} = \bar{x} - 57 \Rightarrow \bar{x} = \bar{u} + 57$$

$$= \bar{u} + 57 = \frac{\sum u}{n} + 57 = 57.3$$

$$C.V = \frac{\sigma}{\bar{x}} \times 100 = \frac{5.7974}{57.3} \times 100 = 10.1176 \%$$

- ② calculate S.D & CV for freqⁿ distribution of marks of 100 candidates given below

Marks 0-20 20-40 40-60 60-80 80-100

freq? 5 12 32 40 11

→ Let $u = \frac{x-50}{20}$.

class	mid values	f_i	u_i^2	$f_i u_i$	$f_i u_i^2$
0-20	10	5	100		
20-40	30	12			
40-60	50	32			
60-80	70	40			
80-100	90	11			

0-20 10 5 100

20-40 30 12

40-60 50 32

60-80 70 40

80-100 90 11

$$U = \frac{x - 50}{20} \Rightarrow x = 50 + 20U$$

$$\bar{X} = 50 + 20\bar{U}$$

$$\bar{X} = 50 + 20 \cdot \frac{\sum f_i u_i}{\sum f_i} = 50 + 20 \cdot \frac{40}{100} = 58$$

$$\text{Var}(x) = (20)^2 \text{Var}(u)$$

$$= 400 \left[\frac{\sum f_i u_i^2}{\sum f_i} - \left(\frac{\sum u_i}{n} \right)^2 \right]$$

$$= 400 \left[\frac{116}{100} - \left(\frac{40}{100} \right)^2 \right] = 400$$

$$S.D(x) = 20$$

$$C.V(x) = \frac{S.D.}{A.m.} \times 100 = \frac{20}{58} \times 100 = 34.4828$$

(3)

The number of runs scored by cricketer A & B in 10 test matches are shown below.

A	5	20	90	76	102	90	6	108	20	16
B	40	35	60	62	58	76	42	30	30	20

- Find
 i) which cricketer is better in average
 ii) which cricketer is more consistent.

$$\rightarrow \text{mean of } A = \frac{533}{10} = 53.3$$

$$SD \text{ of } A = \sqrt{\frac{\sum x^2}{n} - \left(\frac{\sum x}{n}\right)^2} = \sqrt{40.9293} = 6.39$$

$$CV \text{ of } A = \frac{6.39}{53.3} \times 100 = 11.8\%$$

$$\text{Mean of } B = \frac{\sum y}{10} = \frac{453}{10} = 45.3$$

$$SD \text{ of } B = \sqrt{\frac{\sum y^2}{10} - \left(\frac{\sum y}{10}\right)^2} = \sqrt{16.8882} = 4.1$$

$$CV \text{ of } B = \frac{4.1}{45.3} \times 100 = 9.0\%$$

A gives better average run (\because mean A > mean B)

ii) B is more consistent ($CV \text{ of } B < CV \text{ of } A$)

Moments:

Moments give tools for comparison of data set. Using moments we can study the four aspects of comparison viz average, dispersion, symmetry & kurtosis.

There are several aspects of studying freqⁿ distribution. In earlier type we have studied two of them viz average & dispersion. In order to study few more aspects such as symmetry, shape of freqⁿ distⁿ moments are useful.

Defⁿ: -

The quantity $\sum_{n} (x_i - a)^r$ is called r^{th} moment (or moment of order r) about a .

According to choice of a we make 3 categories of moments namely raw moments, central moments, moments about a .

* Raw moments:

Raw moments of order r (r^{th} raw moment) is denoted by u_r & is

$$u_r = \sum_{i=1}^n x_i^r \quad \text{for individual observation}$$

$$= \sum_{i=1}^n f_i x_i^r \quad \text{for freq^n distribution.}$$

Note:- (1) $u_1 = \frac{1}{n} \sum x_i = \bar{x} = \text{A.M.}$

(ii) The raw moments are also called as moments about origin.

Central moments:

Central moment of order r , U_r is given by

$$U_r = \frac{1}{n} \sum (x_i - \bar{x})^r \text{ for individual observation}$$

$$= \frac{\sum f_i (x_i - \bar{x})^r}{\sum f_i} \text{ for freq^n distribution}$$

Here

$$U_1 = \frac{1}{n} \sum (x_i - \bar{x})$$

$$= \frac{1}{n} \sum x_i - \frac{\bar{x}}{n} \sum 1$$

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

$$U_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \sigma^2 = \text{variance.}$$

Mainly we need central moments. However these are difficult to compute as compared to raw moments. Therefore we require to find relation b/w them.

Relation b/w raw & central moment.

$$\therefore U_1 = 0$$

$$\therefore U_2 = \frac{1}{n} \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n} = \frac{1}{n} \sum x_i^2 - (\bar{x})^2$$

$$= \frac{1}{n} \sum x_i^2 - \left(\frac{\sum x}{n} \right)^2$$

$$u_2 = u_2 - (u'_1)^2$$

sim¹⁴

$$u_3 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^3 = \frac{1}{n} \sum x_i^3 - 3\bar{x} \sum x_i^2 + 3\bar{x}^2 \sum x_i - (\bar{x})^3$$

$$= \frac{1}{n} \left[\sum x_i^3 - \frac{3\bar{x}}{n} \sum x_i^2 + 3\bar{x}^2 \frac{1}{n} \sum x_i - \left(\frac{\sum x_i}{n} \right)^3 \right]$$

$$= u'_3 - 3u'_1 u'_2 + 3u'_1 - u'_1$$

$$= u'_3 - 3u'_1 u'_1 + 2u'_1$$

$$u_4 = \frac{1}{n} \sum (x_i - \bar{x})^4$$

$$= u'_4 - 4u'_1 u'_3 + 6u'_2 (u'_1)^2 - 3(u'_1)^3$$

In general

$$u_r = u'_r - r c_1 u'_{r-1} u'_1 + r c_2 u'_{r-2} (u'_1)^2 + \dots$$

#

Moments about a:

suppose a is any arbitrary constant.
then r^{th} moment about a is denoted by $u_r(a)$
& is

$$u_r(a) = \sum \frac{(x_i - a)^r}{n} \text{ for individual observation}$$

$$\sum f_i (x_i - a)^r$$

$$\sum f_i$$

- ① If $a = \bar{x}$ then $U_L(a) = U_R$
- ② If $a = \bar{x}$ then $U_L(a) = U_R$

Skewness & kurtosis:

In previous types we have studied two aspects of freqⁿ dist in the study of freqⁿ distribution viz average, dispersion. However in order to compare two freqⁿ distribution, average & dispersion are not adequate. Sometimes two freqⁿ distribution have same average & dispersion however they differ in symmetry so symmetry is the third aspect in the study of freqⁿ distribution. The term skewness carries opposite meaning to symmetry i.e. lack of symmetry.

Symmetry :-

A freqⁿ distribution is symmetric about a value a if the corresponding freqⁿ curve is symmetric about a . In other words the ordinate at $x=a$ divides freqⁿ curve into two equal part. For symmetric freqⁿ curve these two parts are mirror images of each other. The point a turns out to be Am, mode & median.



eg - class	0-10	10-20	20-30	30-40	40-50
freq dist	5	12	20	10	5

Here freqⁿ of first is the same as that of last class, sim^y second & second last class have equal freqⁿ. & so on.

Properties

- ① In case of bell shaped unimodal symmetric freqⁿ distribution, Am, mode, median coincide
- ② The quartiles of S.D. are equispaced
$$Q_3 - Q_2 = Q_2 - Q_1$$
- ③ The odd order

In day to day life we come across several distributions which are not symmetric

- e.g.: i) distribution of income of individual,
ii) of agricultural land holdings,
iii) of no. of misprints per page

In these situations we need to measure the extent of departure from symmetry.

Skewness:-

is lack of symmetry or departure from symmetry. If the distribution is skew, the corresponding distribution curve is elongated on either side.

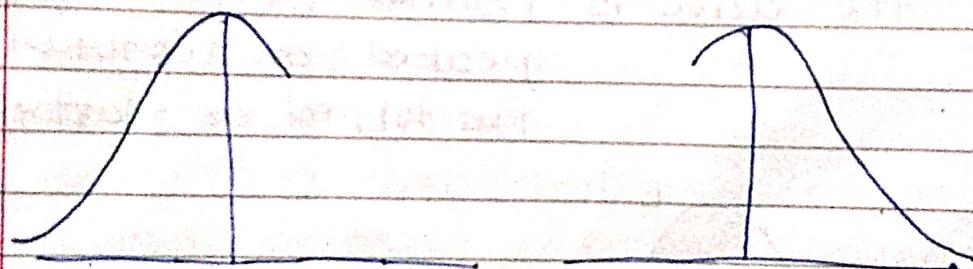
If the curve is elongated towards right side the distribution is said to possess positive skewness. On other hand if it is elongated to left side the distribution is said to possess -ve skewness.

In other words in case of +ve skewness the freqⁿ increases rapidly to reach maxm and further decreases slowly. Exactly reverse process is observed in case of the distribution with -ve skewness.

In case of +vely skew distribution we observe that Mode < median < AM.

Whereas in case of -vely skew distribution we observe that

$$AM < \text{median} < \text{mode.}$$

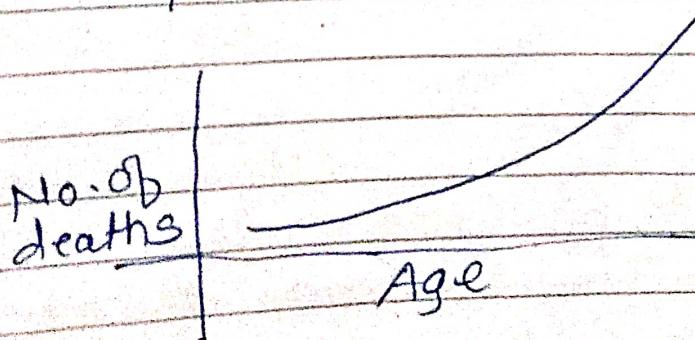


-ve skewness

+ve skew-

e.g. the freqⁿ curve of annual income is +ve skew

The freqⁿ curve of deaths among adults is -vely skew.



$$\text{coeff of skewness} = \frac{\text{mean} - \text{mode}}{\text{s.D.}}$$

Kurtosis :-

It measures the degree of peakedness of a distribution & is given by

$$\beta_2 = \frac{M_4}{M_2^2} \text{ where } M_2 = \frac{\sum (x - \bar{x})^2}{n}$$

$$M_2 = \frac{\sum (x - \bar{x})^2}{n}$$

- $\beta_2 = 3$ the curve is normal or mesokurtic
- $\beta_2 > 3$ peaked or leptokurtic
- $\beta_2 < 3$ flat topped or platykurtic.

