

Statistics Notes

B.Tech. CSE

Gurmukh Singh

Instructor:
Mrs. Neha

Contents

1	Measures of Central Tendency	2
1.1	Mean	2
1.1.1	Properties of Mean	2
1.2	Median	3
1.3	Mode	4
1.4	The interconnection between the measures of central tendency	5
1.5	Geometric and Harmonic mean	5
1.6	Histogram	6
1.7	Ogive	6
1.8	Quartiles	6
1.9	Deciles	6
1.10	Percentiles	6
2	Measures of Spread/Dispersion	6
2.1	Coefficient of variation	7
2.2	Skewness and Kurtosis	7
2.2.1	Positive skewness	8
2.2.2	Negative skewness	8

1 Measures of Central Tendency

1. Mean
2. Median
3. Mode

1.1 Mean

It is the ratio of sum of all the observations to the total number of observations. let x_1, x_2, \dots, x_n be all the observations. then:

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

1.1.1 Properties of Mean

- The sum of deviation of observations from mean is always zero
- the sum of square of deviations of observations is minimum as compared to any other measure.
- suppose there are two sequences:

	Series 1	Series 2
Number of observations	n_1	n_2
mean of the observations	\bar{x}_1	\bar{x}_2

then

$$\bar{x} = \frac{n_1 x_1 + n_2 x_2}{n_1 + n_2}$$

Problem 1

If there are 5 and 8 number of observations of 2 series with mean 15 and 18, find the combined mean

Solution:

We can get the solution by taking the weighted mean of the two sequences.
so the required mean is :

$$\begin{aligned} & \frac{5 \times 15 + 8 \times 18}{5 + 8} \\ &= \frac{75 + 144}{13} \\ &= \frac{219}{13} \\ &= 16.846154 \end{aligned}$$

Problem 2

Class	frequency
0-10	3
10-20	5
20-30	7
30-40	4
40-50	1

Solution:

change of origin:

Class	frequency	X	d=X-A	f·d
0-10	3	5	-20	
10-20	5	15	-10	
20-30	7	25	0	
30-40	4	35	10	
40-50	1	45	20	

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

change of scale

Class	frequency	X	d=X/n	f·d
0-10	3	1	-20	
10-20	5	3	-10	
20-30	7	5	0	
30-40	4	7	10	
40-50	1	9	20	

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

1.2 Median

Steps to find Median in case of Discrete and continuous data:

1. Arrangement of data
2. if n is odd then the median is the $\frac{n+1}{2}$ th term
3. if n is even then the median is the mean of the $\frac{n}{2}$ th term and $\frac{n}{2} + 1$ th term

Problem 3

find the median for the data :

1. 9,9,10,10,12,13,15
2. 9,9,10,10,12,13,14,15

Solution:

1. 9,9,10,10,12,13,15 has 7 elements. Therefore our median will be the 4th term in the arranged order
 $\therefore \text{Median} = 10$
2. 9,9,10,10,12,13,14,15 has 8 elements. Therefore our median will be the mean of the 4th and 5th terms.
 $\therefore \text{Median} = \frac{10+12}{2} = 11$

Problem 4

Finding the median of discrete data.

X	f	cf(cumulative frequency)
1	5	5
2	8	13
3	9	22
4	12	34
5	6	40
6	7	47
7	4	51
Total	51	

find the value of x which has cumulative frequency just greater than $\frac{n}{2}$

In case of continuous data:

$$\text{Median} = l + \frac{\left(\frac{n}{2} - cf\right) h}{f}$$

where cf is the cumulative frequency and f is the frequency of the chosen class, h is the class size

1.3 Mode

The observation which occurs the most is called the mode of the data.

In more general terms, the most probable observation in a dataset is the mode of the data.

Problem 5

Find mode for the following data: 10,11,15,18,18,18,15,10,18,20

Problem 6

Find the mean, median and mode for the following data

CI	f
0-10	3
10-20	5
20-30	7
30-40	2
40-50	1
Total	51

How to find the mode for continuous data

1. Find the modal class which is having the maximum frequency.
2. based on that input the values into the following formulae:

$$mode = l + h \left(\frac{f_1 - f_2}{2f_1 - f_0 - f_2} \right)$$

1.4 The interconnection between the measures of central tendency

$$Mode = 3Median - 2Mean$$

1.5 Geometric and Harmonic mean

Defⁿ :

Geometric mean is defined as the n th root of the product of n observations

Mathematically:

$$GM = \sqrt[n]{\prod_{i=0}^n x_i}$$

Problem 7

Find the Geometric Mean for the values 2,4,8

Defⁿ :

Harmonic mean is defined as the reciprocal of arithmetic mean of the reciprocal of all the observations

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

Theorem:

The following inequality is always true:

$$AM \geq GM \geq HM$$

1.6 Histogram

Histogram can also be used to compute the value of mode.

1.7 Ogive

Ogives are used to compute the value of median. They are nothing but graphs of cumulative distribution functions.

The point of intersection of two ogives gives the median.

1.8 Quartiles

Quartiles are the values which divide the dataset into 4 equal parts. These points are called Q_1, Q_2, Q_3 .

$$Q_1 = l + \frac{\left(\frac{n}{4} - cf\right) h}{f}$$

$$Q_2 = l + \frac{\left(\frac{n}{2} - cf\right) h}{f}$$

$$Q_3 = l + \frac{\left(\frac{3n}{4} - cf\right) h}{f}$$

1.9 Deciles

Deciles are the values which divide the dataset into 10 equal parts.

1.10 Percentiles

Percentiles are the values which divide the dataset into 100 equal parts.

To find the x th percentile we can use the following formulae:

$$p_x = l + \frac{\left(\frac{xn}{100} - cf\right) h}{f}$$

2 Measures of Spread/Dispersion

Measures of spread are a numerical quantity to signify the variation in the observations.

Dispersion means the scatterment of observations.

There are a number of ways to get a gist of the Dispersion.

1. Range: it is the difference between the maximum and minimum value of the dataset.
2. Quartile Deviation: It is equal to $\frac{Q_3 - Q_1}{2}$
3. Mean Deviation: It is the arithmetic mean of absolute value of deviation of observations from average.

$$\begin{aligned} MD &= \frac{\sum |x - \bar{x}|}{n} (Discrete) \\ &= \frac{\sum f|x - \bar{x}|}{N} (Continuous) \end{aligned}$$

4. Standard deviation: It is the positive squareroot of arithmetic mean of square of deviation of observations from mean. It is denoted by σ . This is applicable when your mean is an integer.

$$\sigma = \sqrt{\frac{\sum (x - \bar{x})^2}{n}} (\text{helpful when mean is integral})$$

$$\sigma = \sqrt{\frac{1}{n} \left(\sum x^2 - \frac{(\sum x)^2}{n} \right)} (\text{helpful when mean is non - integral and } x \text{ is small})$$

This is an alternate derivation of the formulae $V = E[X^2] - E[X]^2$

$$d = x - a; \sigma = \sqrt{\frac{1}{n} \left(\sum d^2 - \frac{(\sum d)^2}{n} \right)} (\text{helpful when mean is non - integral and } x \text{ is large})$$

5. Variance: It is the square of standard deviation.

$$V = \sigma^2$$

2.1 Coefficient of variation

It is another way of measuring the spread. It is analogous to standard deviation with respect to the mean.

$$CV = \frac{\sigma}{\bar{X}} \times 100$$

Suppose there are two sequences with n_1 and n_2 observations.
 \bar{x}_1 and \bar{x}_2 are their means.

their combined variance will be:

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

where $d_i = \bar{x}_i - \bar{x}$. here \bar{x} is the combined mean of the series.

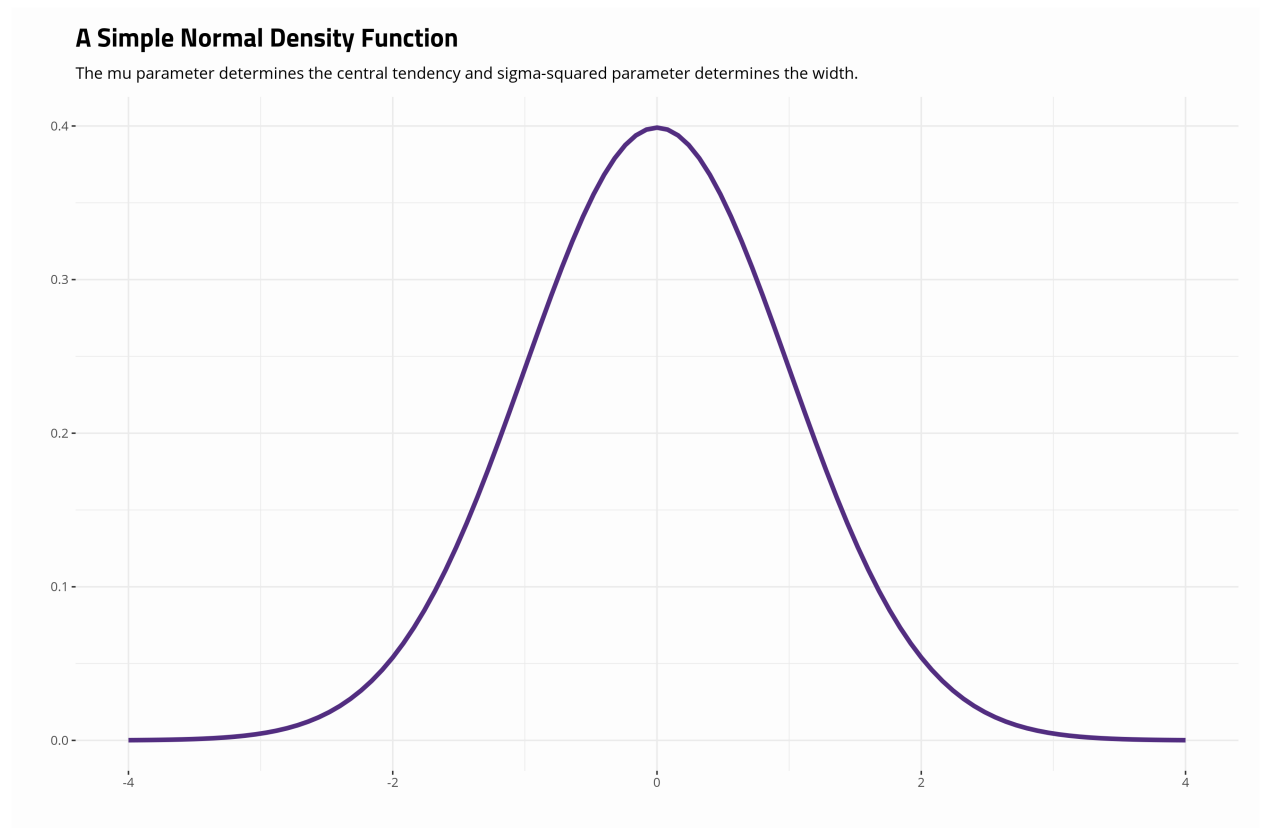
2.2 Skewness and Kurtosis

While studying a distribution we can calculate the measures of central tendency and the measures of spread. However even this information is not enough to determine the behavior of the random variable distribution. In order to further narrow down in the analysis of the behavior of the variable we study the skewness of the distribution.

A distribution is called skewed if:

- The Measures of central tendencies do not coincide.
- The curve does not follow gaussian nature.
- The quartiles are not equidistant from the median.

- Some of the positive deviations from median is not equal to the sum of negative deviations from the median.



There can be two types of skewness:

1. Positive skewness (Right skewed)
2. Negative skewness (Left skewed)

2.2.1 Positive skewness

In this most of the values lie to the right to the peak.

one way to determine this is that $mode > median > mean$

2.2.2 Negative skewness

In this most of the values lie to the left to the peak.

one way to determine this is that $mode < median < mean$