

A
Project Work
On
Data Science
For

Secondary School Certificate Examination

2025-2026

By

Submitted to
Mr.Krishan Singh Sir



4C Colony Jamna Kothi,
New Loha Mandi Rd, Harmada,
Jaipur, Rajasthan 302013

Chapter -1 Use of Statistics in Data Science

What are subsets:

A **subset** is a portion of data taken from a **larger dataset**.

It contains **some (or all)** rows or columns from the original dataset, but **never anything outside it**.
Example

If you have a Table of 100 rows and 100 columns and you want to perform certain actions on the first 5 rows and the first 5 columns, you can separate it from the main table. This small table of 5 rows and 5 columns is known as a “Subset” in Data Analytics.

How do we subset the data?

1. Row-based subsetting:

Selecting specific **rows** (records) from a dataset based on some condition or criteria.

Example:

If you have a dataset of students and you only want students who scored 90 marks.

2. Column based subsetting:

Selecting specific columns (features/variables) from a dataset.

Example:

If your dataset has many columns (name, age, marks, city) but you only need name and marks.

3. Data based subsetting:

Selecting a **subset of data based on both rows and columns** — i.e., applying conditions to rows and selecting specific columns together.

Example:

Students who scored above 80 and only showing their names and marks:

Two-way frequency table:

A **two-way frequency table** shows how two categorical variables are related by counting how many times each pair of values occurs together.

Consider you are conducting a poll asking people if they like chocolates. You record the data in the below format.

Person name	Age	Like chocolates?	Age group	Like chocolates	Do not like chocolates
Person 1	6	Yes	5 - 10	2	1
Person 2	8	Yes	10 - 15	3	1
Person 3	13	Yes	15 - 20	1	2
Person 4	12	Yes			
Person 5	18	No			
Person 6	9	No			
Person 7	16	Yes			
Person 8	19	No			
Person 9	14	No			
Person 10	12	Yes			

Two-way frequency tables show how many data points fit in each category. The row category in this example is "5-10 years", "10-15 years" and "15-20 years". The column category is their choice "Like chocolates" or "Do not like chocolates". Each cell tells us the number (or frequency) of the people.

Two-way relative frequency table:

Two-way relative frequency table very similar to the two-way frequency type of table. Only difference here is we consider percentage instead of numbers.

Preference	Girls	Boys
Indoor games	70%	20%
Outdoor games	30%	80%
Total	100%	100%

Central Tendency:

Central tendency refers to a single value that attempts to describe a set of data by identifying the **central** value within that data set. The three most common measures of central tendency are: Mean, Mode, Median

Mean: Mean, also termed as the simple average, is an average value of a data set. Mean is a value in the data set around which the entire data is spread out.

Median The median like the mean is another form of central tendency. It is the middle point of a sorted data set. To calculate the median, we must order our data set in ascending or descending order, and then select middle value of the set.

Mean vs median So mean and median both represent the central tendency but Median is a more accurate form of central tendency especially in scenarios where there are some irregular values also known as outliers.

Mode: The mode is the number that appear most in a data set. A set of number may have one or more than one mode, or no mode at all.

Mean Absolute Deviation:

Mean Absolute Deviation (MAD) is the average of how far away all values in a data set are from the mean.

Let understand this with an example. Consider the below data set.

Example: Data = 12, 16, 10, 18, 11, 19

Step 1: Calculate the mean

Mean = $(12 + 16 + 10 + 18 + 11 + 19) / 6 = 14$ (rounded off)

Step 2: Calculate the distance of each data point from the mean. We need to find the absolute value. For example if the distance is -2, then we ignore the negative sign.

Value	Distance from mean (14)
12	2
16	2
10	4
18	4
11	3
19	5
Total	20

Fig 1.15 Distance from mean

Step 3: Calculate the mean of the distances.

$$\text{Mean of distances} = (2 + 2 + 4 + 4 + 3 + 5) / 6 = 3.33$$

So 3.33 is our mean absolute deviation, and the mean is 14

What is Standard Deviation?

The Standard Deviation is the measure of how spread out the numbers are. To be specific, standard deviation represents how much the data is spread out around the mean or an average.

1. Calculate the mean by adding up all the data pieces and dividing it by the number of pieces of the data.
2. Subtract mean from every value
3. Square each of the differences
4. Find the average of squared numbers calculated in point number 3 to find the variance.
5. Lastly, find the square root of variance. That is the standard deviation

For example, Take the values 1,2,3,5 and 8 to calculate Standard Deviation

Step 1: Calculate the mean $1+2+3+5+8 = 19$

$$19/5 = 3.8 \text{ (mean)}$$

Step 2: Subtract mean from every value

$$1 - 3.8 = -2.8$$

$$2 - 3.8 = -1.8$$

$$3 - 3.8 = -0.8$$

$$5 - 3.8 = 1.2$$

$$8 - 3.8 = 4.2$$

Step 3: Square each difference

$$-2.8 * -2.8 = 7.84$$

$$-1.8 * -1.8 = 3.24$$

$$-0.8 * -0.8 = 0.64$$

$$1.2 * 1.2 = 1.44$$

$$4.2 * 4.2 = 17.64$$

Step 4: Sum of squared deviations

$$7.84 + 3.24 + 0.64 + 1.44 + 17.64 = 30.8$$

Step 5: Divide by total numbers (for Population SD)

$$30.8/5 = 6.16$$

What is Z-Score?

A Z-score describes the position of a point in terms of its distance from the mean when it is measured in the standard deviation units. The z-score is always positive if the value of z-score lies above the mean and it is negative if its value is below the mean

How to calculate a Z-score?

The mathematical formula for calculating the z-score is as following:

$$Z = (x - \mu) / \sigma$$

Where, X = raw score

μ = Population mean

σ = Population Standard Deviation

Example that will illustrate the use of z-score formula. Consider that we know about a population of group of kids having weights that are normally distributed. Further to this, consider that we know that the mean of the distribution is 10 kgs and the standard deviation is 2 kgs. Now consider the below questions:

1. What is the z-score for 12 kgs?

Z-score for 12 kg

$$z = \frac{12 - 10}{2} = \frac{2}{2} = 1$$

Answer:

The z-score for 12 kg is +1

This means 12 kg is 1 standard deviation above the mean.

2. What is the z-score for 5 kgs?

Z-score for 5 kg

$$z = \frac{5 - 10}{2} = \frac{-5}{2} = -2.5$$

Answer:

The z-score for 5 kg is -2.5

This means 5 kg is 2.5 standard deviations below the mean.

3. How many kgs corresponds to a z-score of 1.25?

Weight corresponding to a z-score of 1.25

Now we rearrange the formula:

$$X = \mu + z\sigma$$

$$X = 10 + (1.25 \times 2)$$

$$X = 10 + 2.5 = 12.5$$

Answer:

A z-score of 1.25 corresponds to 12.5 kg

This means the weight is 1.25 standard deviations above the mean.

Why is a Z-score so important?

It is very helpful to standardize the values of a normal distribution by converting them into z-score because:

1. It gives us an opportunity to calculate the probability of a value occurring within a normal distribution.
2. Z-score allows us to compare two values that are from the different sample

Concept of Percentiles

A percentile can be defined as the percentage of the total ordered observations at or below it. Therefore, pth percentile of a distribution is the value such that p percentage of the ordered observation falls at or below it.

Consider the following data set: [10, 12, 15, 17, 13, 22, 16, 23, 20, 24]

Here, if we want to find the percentile for element 22, we follow the steps below:

Sort the dataset in ascending order. Once sorted, the dataset will look like [10, 12, 13, 15, 16, 17, 20, 22, 23, 24]

The number of values at or below the element 22 is 8. The total number of elements in the dataset is 10

Thus, going by the definition, 80 percent of the values are at or below the element 22. Thus, percentile for the element 22 is 80 percentiles.

Quartiles

Quartiles of dataset partitions the data into four equal parts, with one-fourth of the data values in each part. The total of 100% is divided into four equal parts: 25%, 50%, 75% & 100%. Since the median is defined as the middlemost value in the observation, the median will have 50% of the observations at or below it. Thus, the second quartile(Q2) or the 50th percentile demarcates the median.

interquartile range can be defined as the measure of middle 50% of the values when ordered from lowest to highest. The interquartile range can be calculated by subtracting first quartile(Q1) from the third quartile(Q3).

$$IQR = Q3 - Q1$$

Deciles

Just like quartiles, we have deciles. While quartiles sort the data into four quarters, deciles sort the data into ten equal parts: the 10th, 20th, 30th, 40th, 50th, 60th, 70th, 80th, 90th, 100th.

$$Di = \frac{i * (n + 1)}{10th Data}$$

Steps to calculate decile:

- a. Find out the number of data or variables in the sample or population. This is denoted by n
- b. In the next step, sort all the data or variables in the sample or population in ascending order.
- c. In the next step, based on the decile that is required, calculate the decile by using the formula:

$$Di = \frac{i * (n + 1)}{10th Data}$$

Let us look at an example to understand the concept in detail:

Suppose we have been given 23 random numbers between 20 and 80. We need to represent them as deciles.

Let's say the raw numbers are: [24, 32, 27, 32, 23, 62, 45, 77, 60, 63, 36, 54, 57, 36, 72, 55, 51, 32, 56, 33, 42, 55, 30]

first determine the number of variables in the sample (n). Here n = 23.

We then need to sort the 23 random numbers in ascending order, as shown below

SR. No	Digit
1	23
2	24
3	27
4	30
5	32
6	32
7	32
8	33
9	36
10	36
11	42
12	45
13	51
14	54
15	55
16	55
17	56
18	57
19	60
20	62
21	63
22	72
23	77

Decile	Data position	Value
1	2.4	25.2
2	4.8	31.6
3	7.2	32.2
4	9.6	36
5	12	45
6	14.4	54.4
7	16.8	55.8
8	19.2	60.4
9	21.6	68.4

Now $D_1 = 1 * (n+1)/ 10$ th data

$$= 1 * (23 + 1)/ 10$$

= 2.4th data i.e. data between digit number 2 & 3

Which is $24 + 0.4 * (27- 24) = 25.2$