

① Data distribution

① What is data distribution?

→ A data distribution shows how values are spread in a data set.

→ It tells you where most value lies, how spread out they are, and whether extremes exist.

→ Population vs Sample :-

① Population :- Entire group you want information about.

- Ex:- ① All customers of amazon
- ② All students in india
- ③ All transaction in a year.

→ usually too large, expensive or impossible to measure.

② Sample :- A subset of population, used to estimate population behavior.

Ex:- 10,000 Amazon users.

500 students from different states.

Last 1 month of transaction.

→ In ML & DS we almost always work with samples.

② Measures of central tendency

Ex:- Sales

→ Mean
with

Ex:-

→

- Med
- Av

→ Pr

→ use

→ Av

dist

A S D F G H J K L

② Mean :- Mean = $\frac{\text{Sum of Values}}{\text{Num of Values}}$

Ex:- Salaries = 20, 22, 25, 30
= $(20+22+25+30)/4 = 24.25$

→ Mean becomes completely useless when encountered with outliers.

Ex:- 20, 22, 25, 30, 5,00,000

→ Completely useless.

• Mean is sensitive to outliers.

• In Skewed data mean lies towards the tail.

→ In M.L ~~pip~~ Pipelines :-

→ use mean for clean, symmetric data.

→ Avoid it for income, sales, prices without checking distribution.

③ Variance :- Mean tells center.
Variance tells spread.

→ Now after
the differen
Squaring th

→ Variance is a statistical measure that quantifies the spread or dispersion of a set of data points around their mean (avg).

→ formula

→ To find the Variance we are going to first find the mean of data.

85, 90, 95, 100, 105

$$\rightarrow \text{mean} = (85 + 90 + 95 + 100 + 105) / 5 = 95$$

→ Now after finding the mean we have to compute the difference from each point and mean and squaring them.

85, 90, 95, 100, 105

→ formula

$$\text{mean} = 95$$

$$(85 - 95)^2 + (90 - 95)^2 + (95 - 95)^2 + (100 - 95)^2 + \\ (105 - 95)^2 = 250.$$

→ now just divide the computed difference from the numbers of points in your data.

$$250 / 5 = 50.$$

G H J K L ; " Enter

Now after finding the mean we have to compute the difference from each point and mean and squaring them.

formula for population $\rightarrow \sigma^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2$

find

Population	Sample	Name
N	n	Total element
μ	\bar{x}	mean
σ^2	s^2	Variance
σ	s	STD.

formula for Sample $\rightarrow s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$

the

→ Standard deviation :-

Now to find the standard deviation you just have to take the square root of the Variance.

85, 90, 95, 100, 105

Variance = 50

$$\text{STD Deviation} = \sqrt{50} = 7.071$$

→ calculating the distance

→ formula for Population :-

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \mu)^2}$$

→ formula for Sample :-

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