

20-07-2025

Agenda: Statistics - II

Recap - Statistics - I

- Statistics (Definition, why, when)
- Measure of central tendency
 - Mean
 - Median
 - Mode + outliers
- Measure of Dispersion
 - variance → standard deviation
- keywords of statistics
 - population, sample, variable, parameter, etc...

Types of Data:

- int, string, float, bool, list, set, dictionary, datetime, series etc. (python + pandas + ...)
- merge to create categories → numeric, boolean, string, datetime

Data types of statistics:

- Continuous (float, numeric, range)

Data that can take on any values in an interval.

- Discrete (integer, count)

Data that can take on only integers values, such as counts.

- categorical (classification problems, groups)

Nominal

- Binary (logical, boolean, indicator) →

Special case of categorical data with just two categories 0/1, false/true

- Ordinal

Categorical data that has an explicit ordering.

ex - small < medium < large

- A > B > C > D

A

B

A B

○

(A)

(B)

(C)

(D)

(E)

gender →

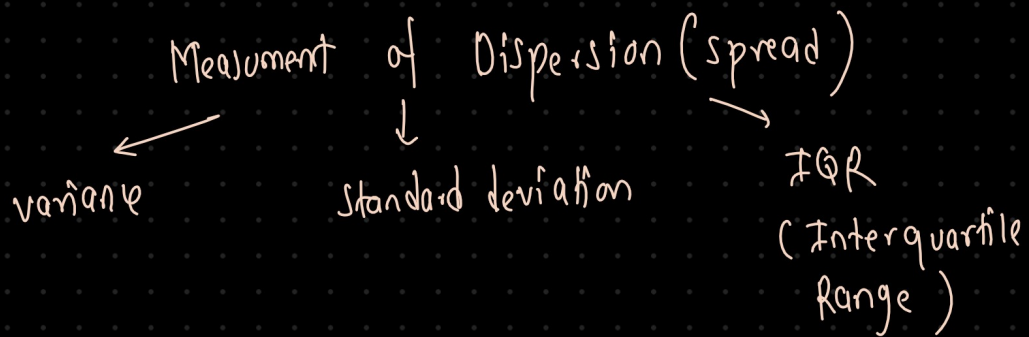
- Male
- female
- other

- Example of Continuous data - Height of a person + Weight of an object + Temperature + Distance traveled + Time taken to complete a task

- Example of Discrete data - Number of website visitors per day + Number of customers served at a restaurant + Number of errors in a software program + Word count in a document + Number of items sold in a store

- Example of Ordinal data - Education level (high school, bachelor's degree, master's degree) + Movie rating (1 star, 2 stars, 3 stars, 4 stars, 5 stars) Shirt size (small, medium, large, extra large) + Customer service rating (poor, fair, good, excellent) + Level of agreement (strongly disagree, disagree, neutral, agree, strongly agree)

- Example of Nominal data - Hair color (blonde, brown, black, red) + Product type (laptop, smartphone, tablet) + Survey responses (yes/no, agree/disagree) + Country of origin (US, China, France)



(1) variance

If population:

$$\text{population variance} = \sum \frac{(x_i - \bar{x})^2}{n} \quad \bar{x} \rightarrow \text{mean}$$

if sample:

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

population variance $\rightarrow \sigma^2$

sample variance $\rightarrow s^2$

population data : $[2, 4, 6, 8]$

$$\rightarrow \text{mean} = 5$$

$$\rightarrow (2-\underline{5})^2 + (4-\underline{5})^2 + (6-\underline{5})^2 + (8-\underline{5})^2 = 20$$

$$\rightarrow \frac{20}{n} = \frac{20}{4} = 5$$

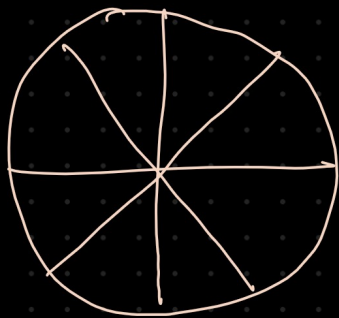
sample data = $[2, 4, 6]$

$$\rightarrow \text{mean} = 4$$

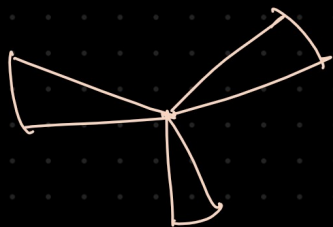
$$\rightarrow (2-4)^2 + (4-4)^2 + (6-4)^2 = 8$$

$$\rightarrow \frac{8}{2} = 4$$

$2 \rightarrow (n-1)$



$\rightarrow 8 \rightarrow \text{slice length} \rightarrow 5$



$\rightarrow 3 \rightarrow \text{slice length} \rightarrow 2.67$

ML model

A → 2024 - 2025 → ML → working
SD → 300

B → 2025 - July 2025 → ML → difference

IQR (InterQuartile Range)

$$IQR = Q3 - Q1$$

→ Jan, Feb, Mar, April, May, June, July, Aug, Sept, Oct, Nov, Dec

Q1
Q2
Q3
Q4

0 - 25% 25 - 50% 50 - 75% 75 - 100%

25th percentile

Scores = [60, 70, 75, 80, 85, 90, 95] →

Q1 → position

$$= 0.25 \times (n+1) =$$

$$\rightarrow \boxed{px(n+1) = 2 \text{ (exclusive)}} \rightarrow 70 \rightarrow$$

$$\rightarrow \rightarrow \boxed{px(n-1) + 1 \text{ (inclusive)}} \rightarrow 72.5 \rightarrow$$

