

20/01/25

Day-2

Page No.

Date

- 1) Measure of Central Tendency
- 2) Measure of Dispersion
- 3) Gaussian Distribution
- 4) Z score
- 5) Standard Normal Distribution

⇒

1. Arithmetic mean for population & Sample :

i) mean (Average)
population (N)

$x = \{1, 1, 2, 2, 3, 3, 4, 5, 5, 6\}$

$$\text{Formula} = \sum x_i = \frac{\sum x}{N}$$

$$= \frac{1+1+2+2+3+3+4+5+5+6}{10}$$

$$= 3.2$$

Sample (n)

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

$$= 3.2$$

* Central tendency

i) mean

ii) median

iii) mode

→ What is Central Tendency?

⇒ Refers to the measures used to determine the centre of the distribution of data.

→ Center of the data

$x = 1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100$

mean = 13.2

Median

i) Sort the numbers

odd number = 11 (Central Element)

median = 3

mode = 3

Median - [1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 100, 112]

$$\frac{3^{th} + 4^{th}}{2} = \frac{3 + 4}{2} = 3.5$$

Median mode = 3.5

→ Median works with well outliers.

* Mode :-

[1, 2, 2, 3, 4, 5, 6, 6, 6, 7, 8, 100, 200]

Mode = 6 {most frequent Element}

↳ Measure of Central Tendency

Ex:- Dataset :-

Type of flower	Petal length	Petal width	Dataset
Rose			100
Lilly			107 - missing data
Sun-Flower			

missing value → most frequent occurring element



Categorical variable

Age of student

Age

25

26

...

...

32

30

mean ? ✓✓✓

median ?

mode ?

* Measure of dispersion Dispersion means Spread

1. Variance
2. Standard deviation

⇒ 1. Variance

Ex: {1, 1, 1, 2, 2, 2, 2, 2} = $\frac{10}{5} = 2$
Avg {2, 2, 2, 2, 2} = $\frac{10}{5} = 2$

both are same but distribution is different.
To identify how to distribute are diff on that time may use variance & sd

Population variance

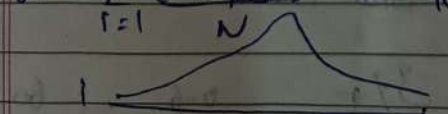
$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N}$$

Sample variance

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

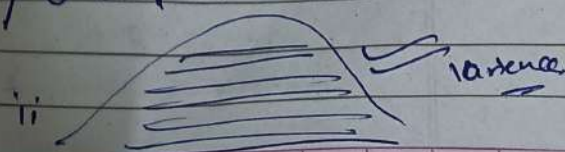
X	-M	X-M	(X-M) ²
1	2.83	-1.83	3.3489
2	2.83	-0.83	0.6889
2	2.83	-0.83	0.6889
3	2.83	0.17	0.0289
4	2.83	1.17	1.3689
5	2.83	2.17	4.7089
<u>N = 14/6 = 2.83</u>			<u>10.84</u>

$$\sigma^2 = \frac{\sum_{i=1}^N (x_i - \mu)^2}{N} = 10.84 / 6 = 1.81$$



which is maximum variance

ii → Variance more



variance

$$\text{Variance} = \sigma^2 = \frac{\sum (x_i - \mu)^2}{n}$$

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

Sample Variance Example

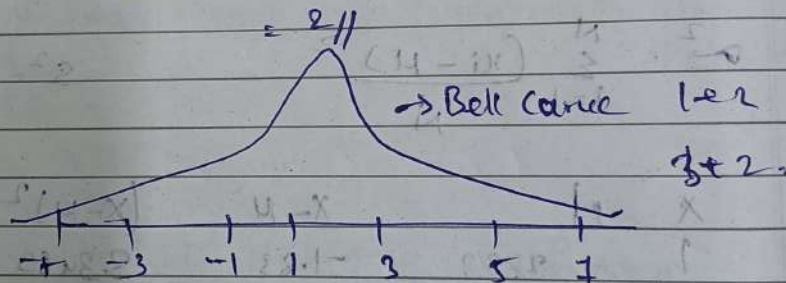
4, 8, 6

$$\bar{x} = \frac{4+8+6}{3} = 6$$

$$\frac{(4-6)^2 + (8-6)^2 + (6-6)^2}{3-1} = \frac{4+4+0}{2} = 4$$

$$S^2 = \frac{8}{3-1} = 4$$

$$\text{Standard deviation} = \sqrt{4} = 2$$



* Percentiles and Quartile } First step to find outliers

Data - 1, 2, 3, 4, 5

% of the numbers that are odd

$\frac{\text{No of numbers that are odd}}{\text{Total Numbers}}$

$$= \frac{3}{5} = 0.6 = 60\%$$

Percentile :- Definition :-

→ A Percentile is a value below which a certain percentage of observation lie.

Eg:-

Data :- 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

Q) what is the Percentile ranking of 10 p

Percentile rank of 10

$$P = \frac{\text{No of values below } x}{n} \times 100$$

$$n = 20$$

$$= \frac{16}{20} \times 100 = 80 \text{ Percentile}$$

Q) What is percentile rank of 11

$$P = \frac{\text{No of values below } x}{n} \times 100$$

$$= \frac{17}{20} \times 100 = 85\%$$

Q) What value exists at Percentile ranking of 25%?

Formula: Value = $\frac{\text{Percentile} \times (n+1)}{100}$

$$= \frac{25}{100} \times (20+1)$$

$$= \frac{25}{100} \times 21 = 5.25 //$$

This 5.25 is the index position

Take 5th and 6th Index
we do avg

$$5+5/2 = 10/2 = 5$$

5 → 25% (percentile)

(9) what is 75%?

$$\frac{75}{100} \times (n+1) = \frac{75}{100} \times (20+1)$$

$$\frac{75}{100} (21) = 15.75 \quad \left(\text{Index value} \right)$$

we take 15 and 16 term

$$9+9/2 = 18/2 = 9$$

* Five Number Summary :-

i) Minimum

ii) First quartile (Q1)

iii) median

iv) Third quartile (Q3)

v) maximum

⇒ ⇒ Removing the outliers :-

Eg:- {1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 7, 8, 8, 9, 27}

27 is the outlier

{ Lower Fence ← → Higher Fence }

IQR: Inter Quartile

Page No.	
Date	

$$\text{Lower Fence} = Q_1 - 1.5(IQR)$$

$$\text{Upper Fence} = Q_3 + 1.5(IQR)$$

$$IQR = Q_3 - Q_1$$

$$\text{Ex: } Q_3 = 75.1 \rightarrow 75/100 \times (19+1) = 15(\text{Index}) = 71$$

$$Q_1 = 25.1 \rightarrow 25/100 \times (19+1) = 5(\text{Index}) = 31$$

$$IQR = Q_3 - Q_1$$

$$= 71 - 31$$

$$= 40$$

Ex

$$\text{Lower Fence} = Q_1 - 1.5(IQR)$$

$$= 3 - 1.5(4)$$

$$= -3$$

$$\text{Higher Fence} = Q_3 + 1.5(IQR)$$

$$= 7 + 1.5(4)$$

$$= 7 + 6$$

$$= 13$$

$$\begin{array}{ccc} \text{Lower Fence} & \longleftrightarrow & \text{Higher Fence} \\ [-3] & \longleftrightarrow & [13] \end{array}$$

Outlier 13 > value

27 outlier

Remaining data

1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 9

Minimum = 1

Q1 = 3

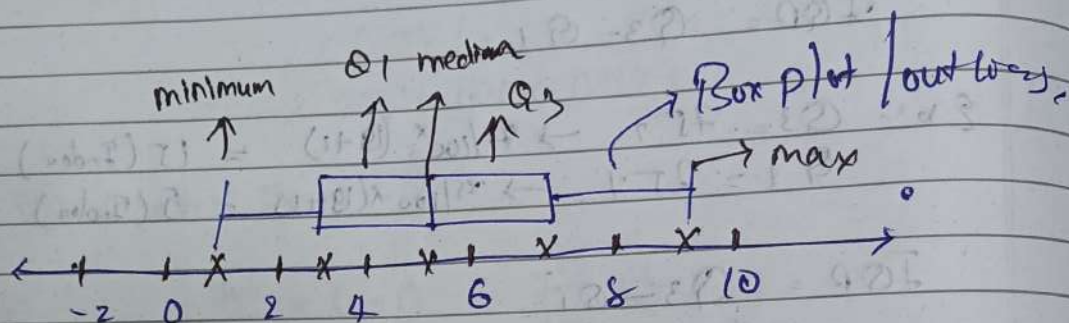
Median = 5

Q3 = 7

Max = 9

5 Number Summary

Box plot



\Rightarrow this we use for data visualization

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