

AGENDA

- ① Histogram
- ② Measure of Central tendency
- ③ Measure of Dispersion
- ④ Percentiles & Quartiles
- ⑤ 5 number Summary (Box Plot)

22

* Histogram

⇒ Ages = {10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100}

* Sort the numbers

Bins: → No. of groups

Bin size: → Size of Bins

How to choose bins: It's totally your choice what you want to choose bin like 10, 5 or 20

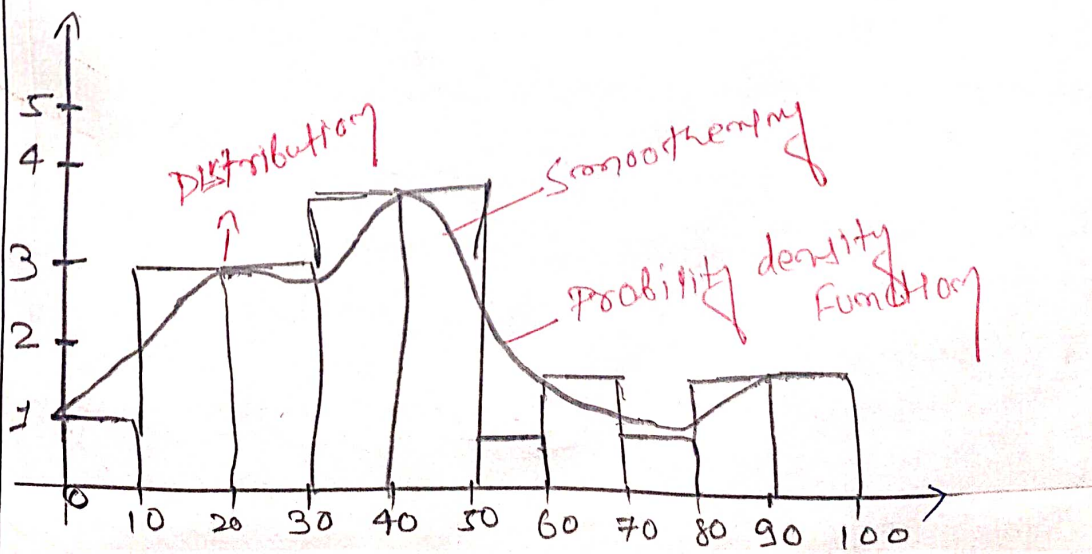
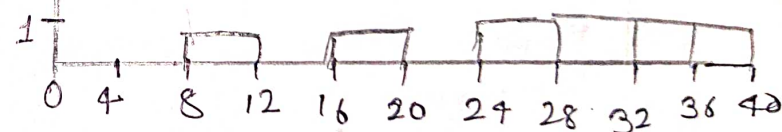
Frequency

5
4
3
2
1

[10, 20, 25, 30, 35, 40]

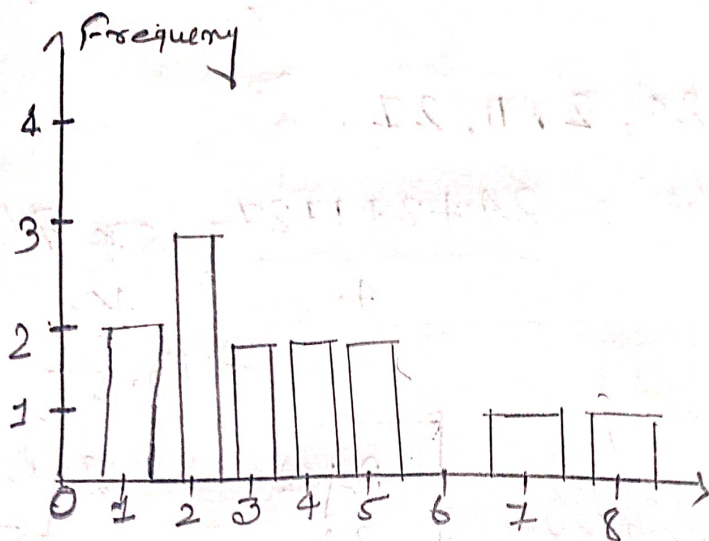
min = 10
Max = 40
bins = 10

bin = 10



① Discrete Continuous

Rank = $\{1, 2, 3, 5, 7, 8, 3, 2, 4, 5, 1, 2, 4\}$



2 Probability
Max function

② Measure of Central tendency

- ① Mean
- ② Median
- ③ Mode

A measure of central tendency is a single value that attempts to describe a set of data identify the central position.

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

$$\text{Average / Mean} = \frac{1+2+3+4+5}{5} = 3$$

Population (N) $N \gg n$

Sample (n)

Population mean (μ) =
$$\sum_{i=1}^N \frac{x_i}{N}$$

Sample mean (\bar{x}) \Rightarrow

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

Population age = { 24, 23, 2, 1, 28, 27 }

$$N = 6$$

$$\text{Population Mean } (\mu) = \frac{24 + 23 + 2 + 1 + 28 + 27}{6} = 17.5$$

Sample age = { 24, 2, 1, 27 }

$$\text{Sample mean } (\bar{x}) = \frac{24 + 2 + 1 + 27}{4} = \frac{54}{4} = 13.5$$

$$\bar{x} = 13.5$$

Practical implementation

Age	Salary	Family Size
—	—	—
—	NaN	—
NaN	—	—
—	—	—
—	—	NaN

(Replace NaN with mean, Median, Mode According to situation)

Age	Salary
24	45
28	50
29	NaN
31	60
36	75
NaN	80
NaN	NaN
Mean of Age = $\frac{24 + 28 + 29 + 31 + 36}{5} = 29.6$	
Mean of Salary = $\frac{45 + 50 + 60 + 70 + 80}{5} = 62$	

Median

{ 1, 2, 3, 4, 5 }

$$\bar{x} = 3$$

{ 1, 2, 3, 4, 5, 100 }

$$\bar{x} = 19.16$$

Steps to find out the Medians:-

Sort the no.

Find Central Number

→ if the no. of element are even find avg. of two central elements

→ if the no. of elements are odd we simply choose the Central Number.

Eg:- For even

{ 1, 2, 3, 4, 5, 6, 7, 8, 100, 120 }

$$\text{Mean} = 25.6$$

$$\text{Median} = \frac{5+6}{2} = \frac{11}{2} = 5.5$$

* When there is no outliers \rightarrow Mean

* When there is outliers \rightarrow Median

* Mode

\rightarrow Most Frequent occurring element

Eg:- { 1, 2, 2, 2, 3, 3, 3, 3, 4, 5 } or { 1, 2, 2, 2, 3, 3, 3, 4, 5 }

Mode = 3

Mode = 2, 3

Practical Application

Types of flowers

Lily
Sunflower

{ Replace NAN with Rose }

Rose

NAN

Rose

Sunflower

Rose

\Rightarrow use Mode with categorical variable

Measure of Dispersion

* Variance (σ^2) \rightarrow Spread of Data

* Standard deviation (σ)

Variance

Population variance (σ^2)

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

Sample variance (s^2)

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

Population Variance (σ^2)

$\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$

$\{1, 2, 3, 4, 5\}$

$$\mu = 3$$

$$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5}$$

$$= \frac{4 + 1 + 0 + 1 + 4}{5} = \frac{10}{5} = 2$$

$$\boxed{\sigma^2 = 2}$$

Sample Variance (s^2)

$\{1, 2, 3, 4, 50, 60, 70, 80\}$

$\{1, 2, 3, 4, 5, 6, 80\}$

$$\mu = 14.4$$

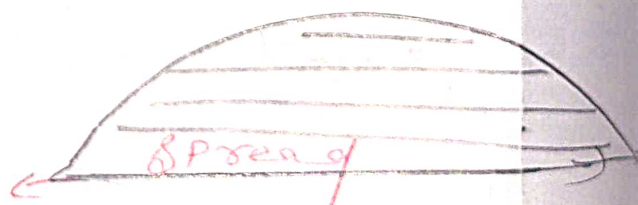
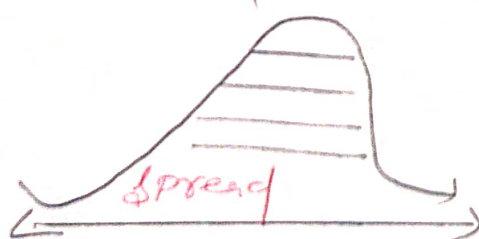
$$s^2 = \frac{(1-14.4)^2 + (2-14.4)^2 + \dots + (80-14.4)^2}{n}$$

$$\boxed{s^2 = 719.10}$$

Variance \uparrow



Spread \uparrow



Standard deviation ($\sqrt{\sigma^2}$)

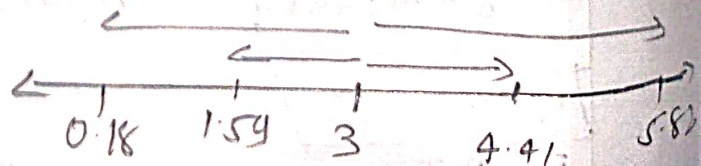
$\{1, 2, 3, 4, 5\}$

$$\mu = 3$$

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

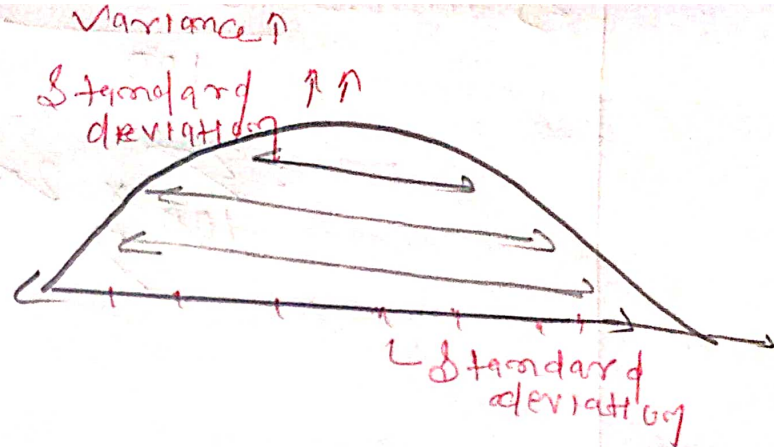
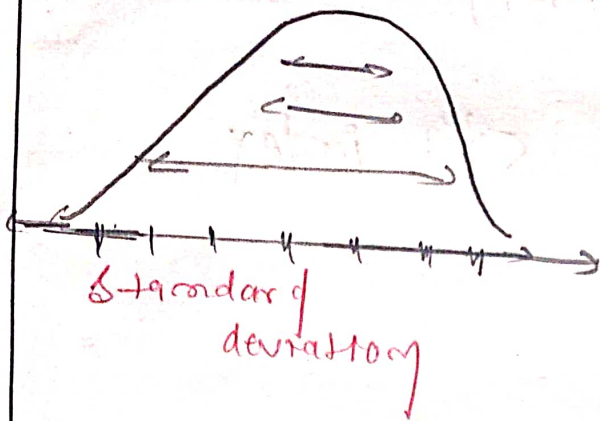
$$\sigma^2 = 2$$

$$\sigma = \sqrt{2} = 1.41$$



$$\mu + \sigma = 3 + 1.41 = 4.41$$

$$\mu - \sigma = 3 - 1.41 = 1.59$$



④ Percentile & Quartiles

Percentiles: - $\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$

Percentile of even no. $\Rightarrow \frac{\text{No. of 70 even no.}}{\text{Total no}} = \frac{6}{12} = 50\%$

Definition

\Rightarrow A Percentile is a value below which a certain percentage of observation lie

99 Percentile \Rightarrow It means the person has got better marks than 99% of the entire students.

* Rank for same multiple value will be same

Dataset $\{2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12\}$

Q. What is the Percentile Rank of 10?

$n = 20$

$$\text{Percentage Rank} = \frac{\text{No. of value below } x}{n} = \frac{16}{20} = 80\%$$

② What is the value that exist at 25% Percentile

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

$$= \frac{25}{100} \times 26 = 5^{\text{th}} \text{ index}$$

⑤

5 number Summary

* Minimum

* First Quartile (Q_1) 25%

Box Plot

* Median

* Third Quartile (Q_3) 75%

\Rightarrow Remove the outliers

* Maximum

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

outliers
7

How to Remove outliers

$$\text{Lower Fence} = Q_1 - 1.5 \times (IQR) \quad \left. \begin{array}{l} IQR = Q_3 - Q_1 \\ \text{Inter Quartile Range} \end{array} \right\}$$

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

Standard deviation

$$Q_1 = \frac{25}{100} \times (n+1)$$

$$= \frac{25}{100} \times (26+1) \Rightarrow 5.25 \quad \boxed{\text{index} = 3}$$

$$Q_3 = \frac{75}{100} \times (n+1)$$

$$= \frac{75}{100} \times (26+1) = 18.75 \quad \boxed{\text{index} = \frac{8+7}{2} = 7.5}$$

$$\text{Lower Fence} = 3 - (1.5)(4.5) \Rightarrow -3.65$$

$$\text{Higher Fence} = 7.5 + (1.5)(4.5) = 14.25$$

$$\boxed{IQR = 7.5 - 3 = 4.5}$$

Box Plot

5 No. Summary

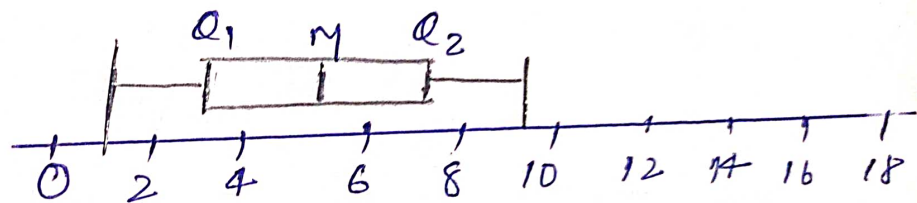
① Minimum - 1

② Q_1 - 3

③ Median - 5

④ Q_3 - 7.5

⑤ Maximum - 9



To treat outliers