

# ★ Statistics Day - 2

## Agenda

- ① Create histograms
- ② Measure of Centre Tendency
- ③ Measure of dispersion
- ④ Percentiles And Quartiles
- ⑤ 5 number summary (Box plot).

## ① Histogram

Age = {0, 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100}

- ① Sort the numbers.
- ② Bin - No. of groups
- ③ Bin Size - Size of bins.

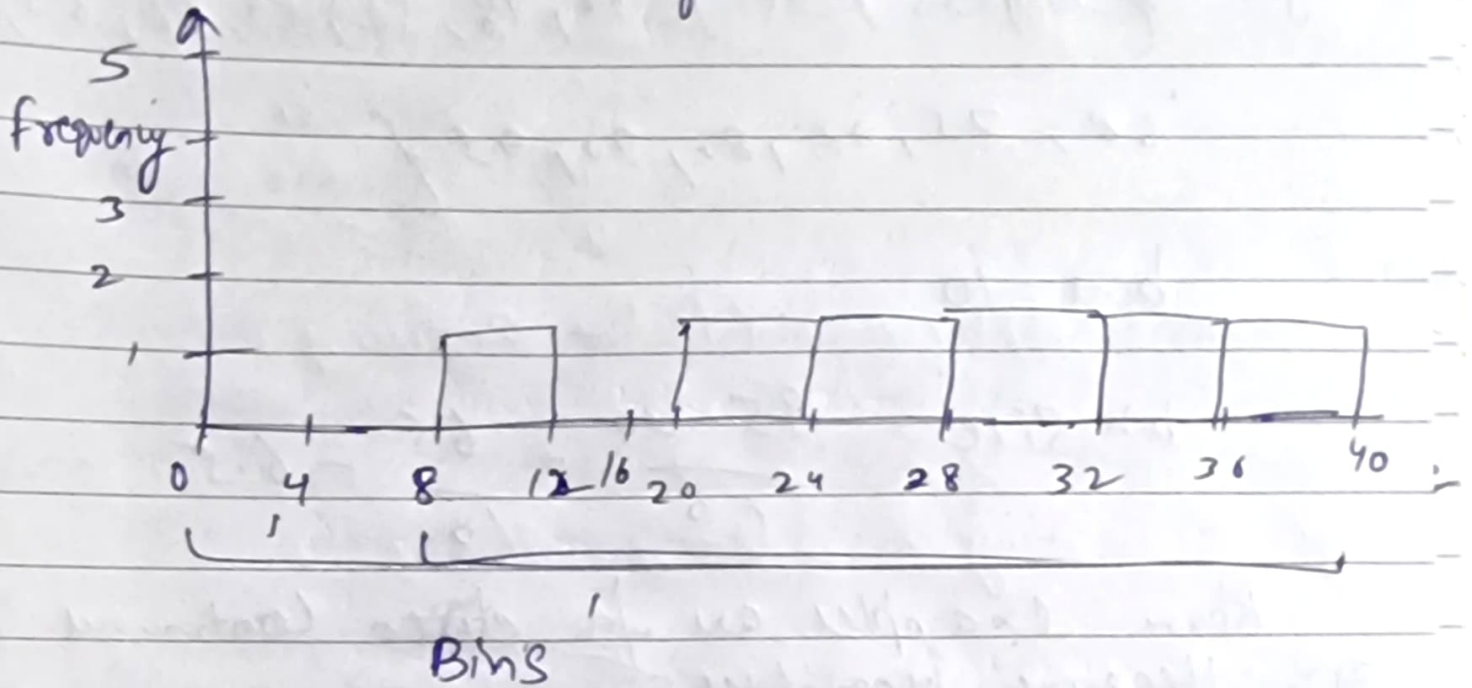
[10, 20, 25, 30, 35, 40] min = 10

$$\text{Bin Size} = \frac{\text{Max value}}{\text{Bin}}$$

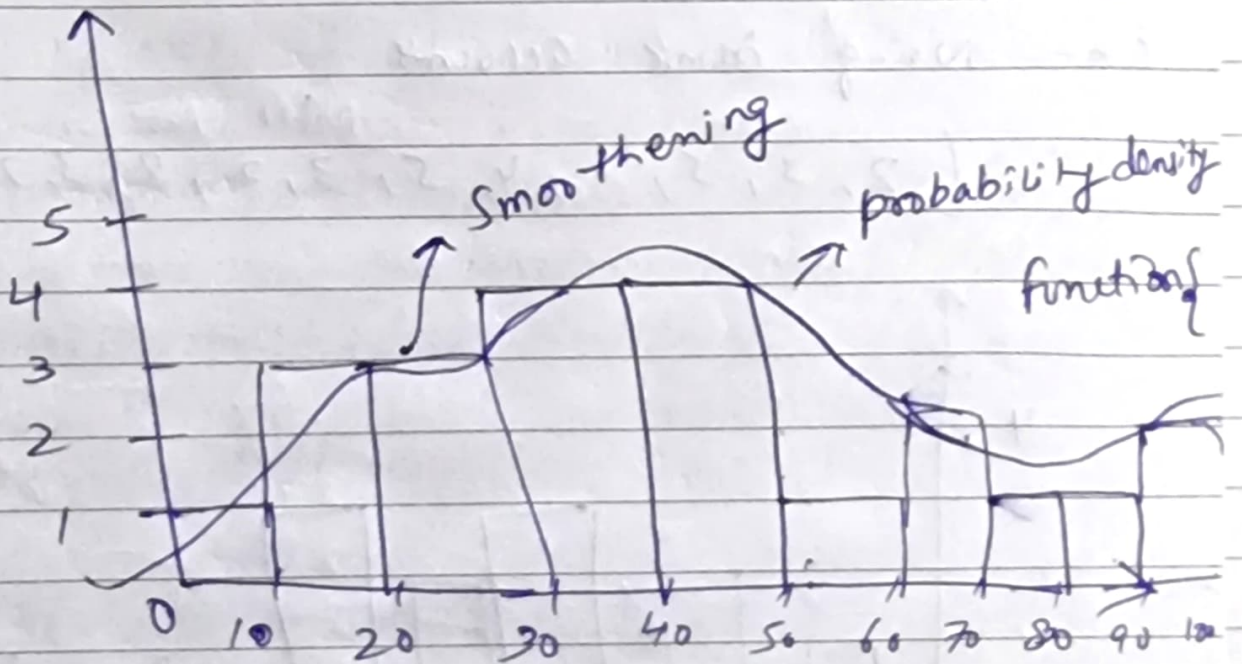
max = 40

$$= \frac{40}{10} = 4$$

Bin Size = 4 for above list then



If Bin Size =  $\frac{100}{10} = 10$ .





Outlier - w/c looks completely different outside the distribution.

### \* Median

$$\{1, 2, 3, 4, 5\} \text{ vs } \{1, 2, 3, 4, 5, \overset{\text{outlier}}{\downarrow} (100)\}$$
$$\bar{u} = 3 \qquad \bar{u} = 19.16$$

### \* Steps to find out median

- ① Sort the numbers
- ② Find the central number

↓  
[ if the no. of element are even, we find the average of central element ]

[ If the no. of element are odd, we find the central element. ]

Ex - Even  $\rightarrow \{1, 2, 3, 4, 5, 6, 7, 8, 100, 120\}$

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

odd  $\rightarrow \{0, 1, 2, 3, 4, 5, 6, 7, 8, 100, 120\}$

$$\text{median} = 5$$

Sample age:  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{24+28+1+28}{4}$

$\bar{x} = 13.5$

## Practical Application (Feature Engineering)

↓                      ↓                      ↓  
Age                      Salary                      Family size

[NaN]                      ← loss of info

NAN - Not a number

	Age	Salary
	24	45
	28	50
Age: 29.6	29	Nan
Salary: 62	Nan	60
	31	75
	36	80
	Nan	Nan



## ★ (2) Measure of Central Tendency

- (A) Mean
- (B) Median
- (C) Mode

Measure of Central tendency is the ~~measure~~ ~~of C.T.~~ Single value that attempts to describe a set of data identifying the central position.

$$\text{Mean} - X = \{1, 2, 3, 4, 5\}$$

$$\text{Average Mean} = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

(Center of all Elements)

Population (N)

Sample (n)

Population mean ( $\mu$ )

Sample mean ( $\bar{x}$ ) =  $\sum_{i=1}^n \frac{x_i}{n}$

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

Population Age =  $\{24, 23, 21, 28, 77\}$   
 $N = 6$

$$\text{Population Mean} = (\mu) = \frac{24+23+21+28+77}{6} = 17.5$$



Weight = { 30, 35, 38, 42, 48, 58, 59, 62, 63, 68, 75, 77, 80, 90, 95 }

bins = 10

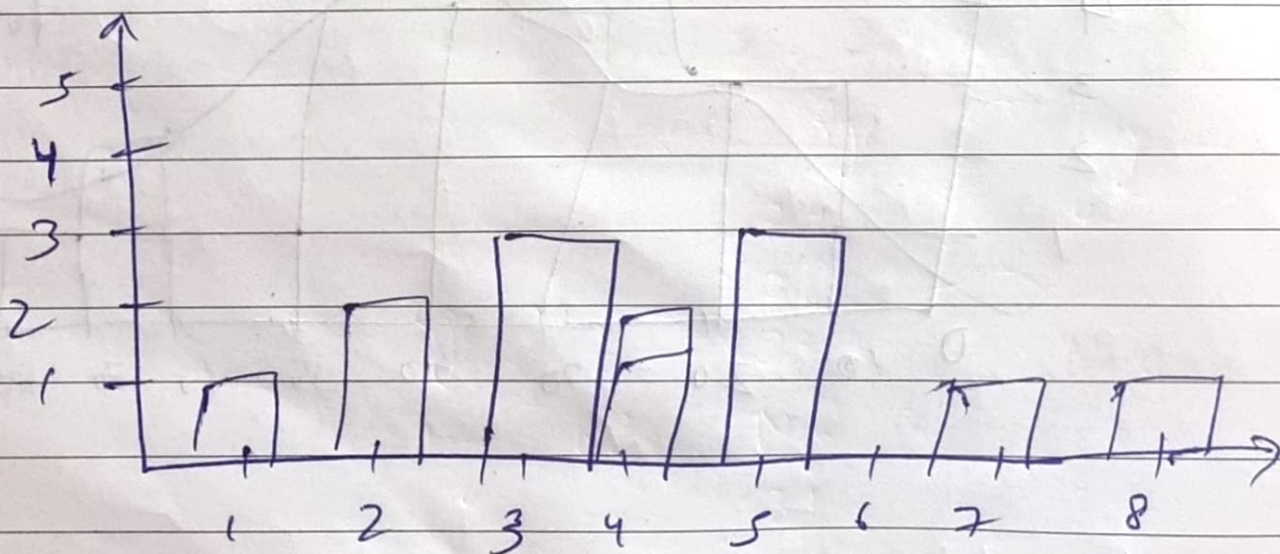
$$\text{bin size} = \frac{95 - 30}{10} = \frac{65}{10} = 6.5$$

Above examples are for ~~discrete~~ Continuous ~~Measure~~ Variable.

# Now - for Discrete Continuous

Ex - No. of Bank accounts

= [ 2, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 9, 5 ]



for smoother Probability mass function is used  
will discuss later.



No Outlier  $\rightarrow$  Mean  
Outliers - Median.

Mode - Most frequent occurring elements?

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

Mode = 3

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

[ 2, 3 ]

Practical aspects :

Data Set - Type of Flower

Lily  
Sunflower  
Nan  
Rose  
Sunflower  
Rose  
Nan

- In this we can replace Nan with Mode value i.e. Rose.



\* what is the value that consist of 25 percentile

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

$$= \frac{25}{100} \times 20 = 5^{\text{th}} \text{ index}$$

### ⑤ 5 number Summary

① Minimum

② First Quartile (25 percentile) ( $Q_1$ )

③ Median

④ Third Quartile (75 percentile) ( $Q_3$ )

⑤ Maximum

we remove  
the outliers  
(Box Plot)

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

[Lower Fence  $\longleftrightarrow$  Higher fence]

$\downarrow$   
 $Q_1$

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

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



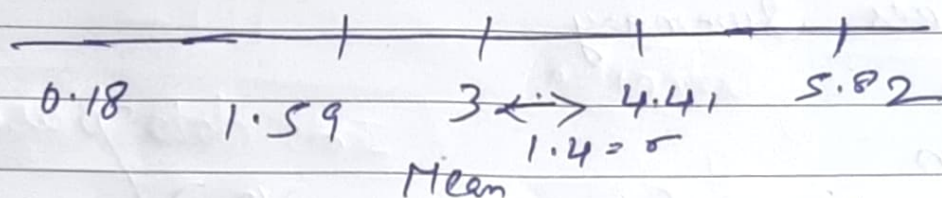
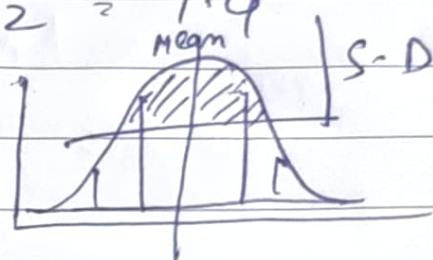
③ Standard deviation  $(\sqrt{\sigma^2}) = \sigma$

A Statistic that measure the dispersion of a dataset relative to its mean and is calculated by all the Square root of variance

$$\sigma^2 = 2$$

$$\sigma = 1.4$$

$$\sigma = \sqrt{2} = 1.4$$



④ Percentile & Quartile.

A percentile is a value below w/c a certain percentage of observation lie.

99 percentile - It means the person has got better marks than 99% of the entire students.

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

- What is the percentile rank of 10

$$\text{Percentile rank of } x = \frac{\text{No. of value below } x}{n} \times 100 = \frac{16}{20} \times 100 = 80$$



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

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

$$s^2 = \frac{(1-14.4)^2 + (2-14.4)^2 + (3-14.4)^2 + (4-14.4)^2 + (5-14.4)^2 + (6-14.4)^2 + (80 - (4 \cdot 4))^2}{7}$$

$$s^2 = 719.10$$

- Dof
- Bias Correction



### ③ Measure of Dispersion

① Variance ( $\sigma^2$ ) ← Spread of Data

② Standard deviation ( $\sigma$ ) ←

Variance

Population Variance ( $\sigma^2$ )

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

Sample Variance ( $S^2$ )

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

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

{1, 2, 3, 4, 50, 60, 70, 100}

{1, 2, 3, 4, 5}

$\mu = 3$

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

$\mu$

$\mu = 14.4$



$$Q_1 = \frac{25}{100} \times 21 = 5.25, \text{Index} = 3$$

IQR = (Inter Quartile Range)

$$= Q_3 - Q_1 = 7.5 - 2.5$$

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

① Minimum = 1

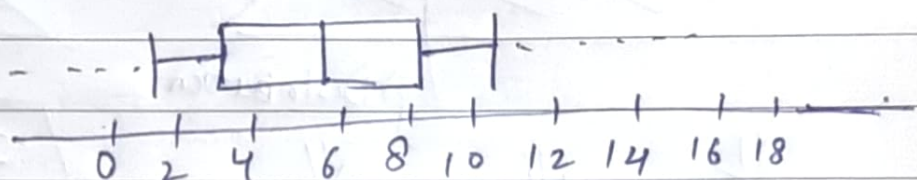
②  $Q_1 = 3$

Box Plot

③ Median = 5

④  $Q_3 = 7.5$

⑤ Max = 9



↓  
used to treat outliers

↓  
Number Summary

To Calculate  $Q_3$

$$Q_3 = \frac{75}{100} \times 21 = 15.75$$

$$\text{Index} = \frac{8+7}{2} = 7.5$$

$$Q_1 = \frac{25}{100} \times 21 = 5.25$$

$$\text{Index} = 3$$

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

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