

## # Day - 2

### Agenda

- ① Histograms
- ② measure of central tendency
- ③ measure of dispersion
- ④ Percentiles and quartiles
- ⑤ 5 number summary (Box plot)

### # Histogram

→ can have continuous variables as explicit

set ① Ages = { 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43, 50, 51, 55, 68, 78, 90, 95, 100 }

### Steps to create

→ sort the numbers

→ define no. of bins (no. of groups) (It depends on user's choice, they can take any value.)

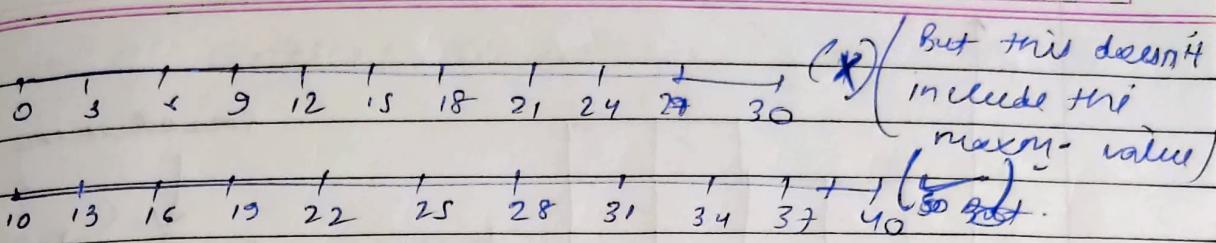
[ 10, 20, 25, 30, 35, 40 ]

$$\text{min} = 10$$

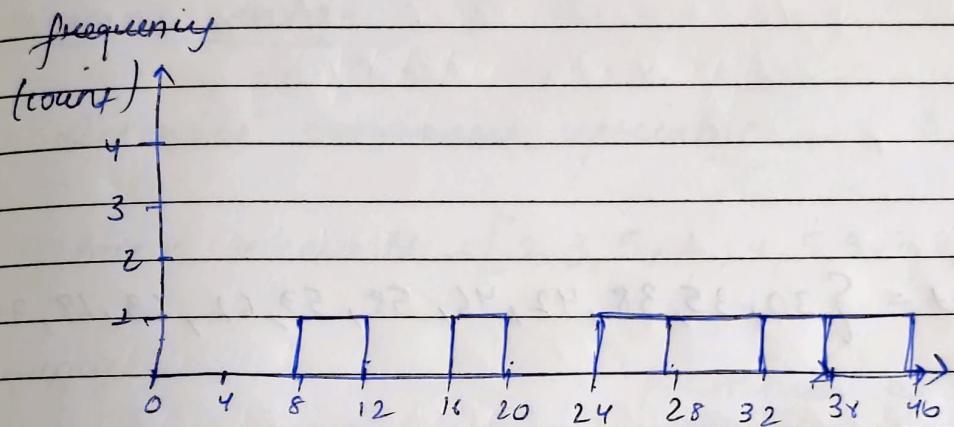
$$\text{max} = 40$$

$$\text{bins} = 10$$

$$\text{Bin size} = \frac{40 - 10}{10} = \frac{30}{10} = 3$$



$$\frac{40}{10} = 4 \quad \left( \text{Bin size} = \frac{\text{maxm. value}}{\text{no. of bins}} \right)$$

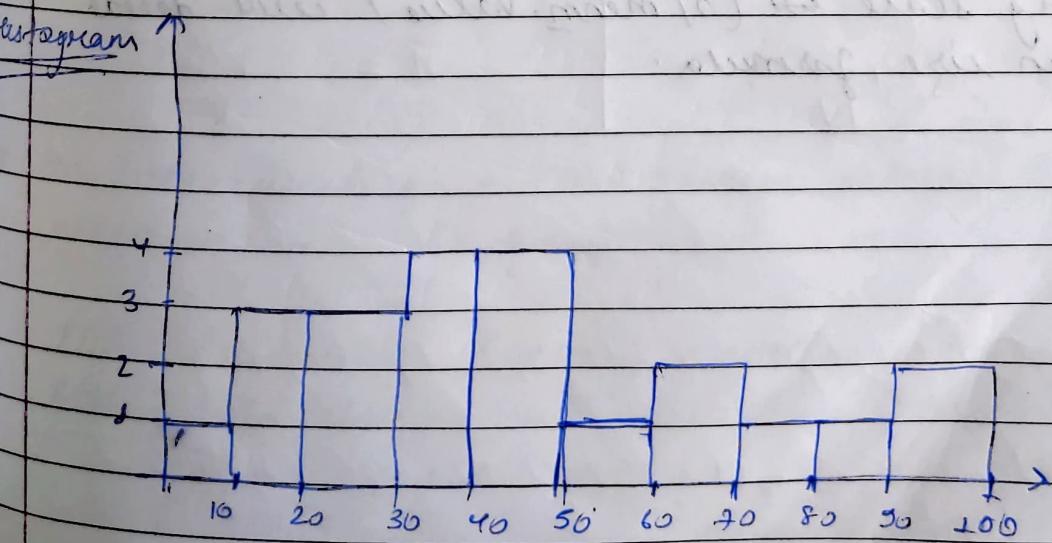


for set (A)

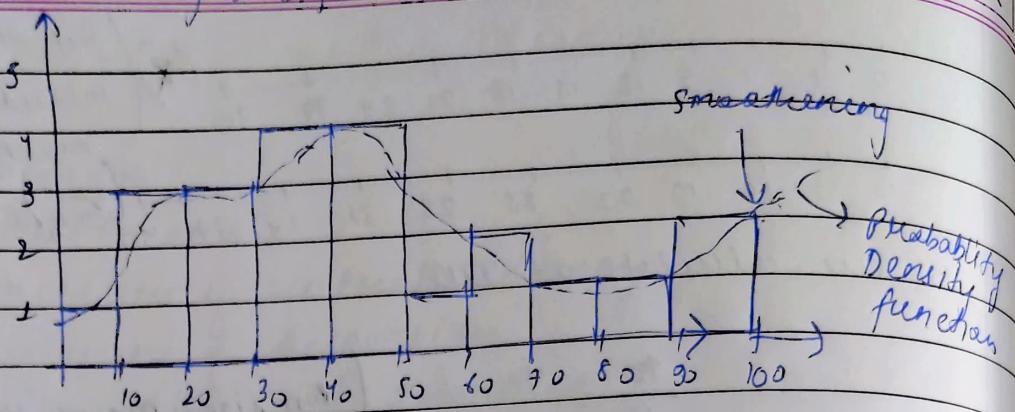
$$\text{maxm. value} = \frac{100}{10} = 10$$

(10)  
Bin size

If Bin size we want as 20,  
then no. of Bins will be

$$\frac{100}{20} = 5$$


## Histogram for continuous variable



set (B) weight = {30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 80, 92}

97

we want  $B_{min} = 10$

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

Here scale will start from 30

→ Starting scale w.r.t (or minm. value) will decide the bin size formula.

# Measures of central tendency

- (i) Mean
- (ii) Median
- (iii) Mode

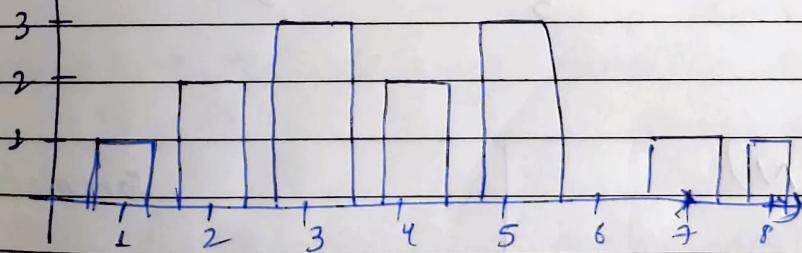
## Histogram contd.

For discrete continuous variable

No. of Bank accounts = [3, 3, 5, 1, 4, 5, 3, 7, 8, 3, 2, 14, 5]

frequency

(Histogram for discrete variables)



To smoothen we will use  $\rightarrow$  Probability Mass function (PMF)

In interview, questions about variable will be asked on basis of graph

PDF: Probability Density Function  $\rightarrow$  Continuous Variable

PMF: Probability Mass Function  $\rightarrow$  Discrete Variable

## #1 Measuring central tendency

- (1) Mean
- (2) Median
- (3) Mode

Defn: A measure of central tendency (CT) is a single value that attempts to describe a set of data identifying the central position.

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

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

Population ( $N$ )

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

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

sample ( $n$ )

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

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

Example:

$$\text{Population Age} = \{24, 23, 21, 1, 28, 27\}$$

$$\begin{aligned} \text{Population mean} (\bar{x}) &= \frac{24+23+21+1+28+27}{6} = \frac{105}{6} = 17.5 \end{aligned}$$

Sample Avg Age(n) = {24, 2, 1, 27}

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

$$= \frac{54}{4} =$$

$$\bar{x} = 13.5$$

$N > n$	✓
$n > N$	✗
$y > \bar{x}$	✓
$\bar{x} > y$	✓

## # Practical Application (Feature Engineering)

Features  $\rightarrow$  Age, Salary, Family size

Age	salary	Family size
—	—	—
—	—	—
NAN	—	—
—	—	—
—	NAN	—
—	—	NAN
—	NAN	—
NAN	—	—

| — If we drop any row with NAN value then there will be loss of info

| we can replace NAN with mean values of respective features

Age	salary
24	45
28	50
29	<u>NAN</u> ← 62
29.6 → <u>NAN</u>	60
31	75
36	80
29.1 → <u>NAN</u>	<u>NAN</u> ← 62

$$\text{Age mean} = 29.6$$

$$\text{Salary } n = 62$$

Suppose we are adding a new row with  
as

$$\underline{80} \quad \underline{200}$$

then mean for

$$\text{Age} = 38$$

$$\text{Salary} = 85$$

values 80, 200 are outliers since they  
are outside the range values

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

$$\{1, 2, 3, 4, 5, \underline{200}\} \bar{x} = 19.16$$

↑  
outliers

## Median

1) Sort the numbers

2) Find the central number

→ No. of elements → even → Find the average of central elements

→ No. of elements → Find the central element

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

$$\text{median} = \frac{5+6}{2} = 5.5 \quad , \quad \text{mean} = 25.1$$

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

$$\text{median} = 5 \quad , \quad \text{mean} = 25.6$$

\* If no outliers then use mean

If outliers are present then use median

# Mode

→ Most frequent occurring element

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

$$\text{mode} = 3$$

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

$$\text{mode} = 2, 3$$

(In older versions of Python  
it would return 2 & 3. But  
in newer versions, it  
returns an error.)

Types of flowers

Lily

Sunflowers

Rose

NAN ← Rose

Rose

Sunflowers

Rose

NAN ← Rose

Replacing data with

## # Measure of dispersion

(i) Variance ( $\sigma^2$ ) ← spread of data  
 (ii) standard deviation ( $\sigma$ )

### Variance

(i) Population variance ( $\sigma^2$ )  
 (ii) sample variance ( $s^2$ )

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

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

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

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

### # Population variance

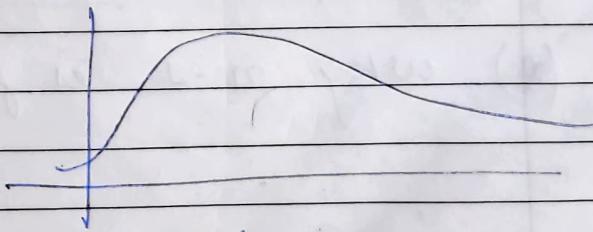
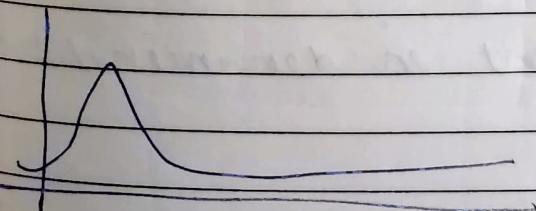
$$\{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\}$$

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

variance = less

variance = more

### Distribution A



spread over more area

{ 1, 2, 3, 4, 5 }

$$\bar{y} = 3$$

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

$$y - 14.4$$

$$\sigma^2 = \sum \frac{(x_i - \bar{y})^2}{N}$$

$$\begin{aligned}\sigma^2 &= \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5} \\ &= \frac{4+1+0+1+4}{5} = \frac{10}{5} = 2\end{aligned}$$

$$\sigma^2 = \frac{(1-14.4)^2 + (2-14.4)^2 + \dots}{7}$$

variance

variance ↑ spread ↑

$$(\sigma^2)$$

Assignment

Assignment :-

- (Q) Why  $n-1$  is present in denominator?

## standard deviation ( $\sqrt{\sigma^2}$ )

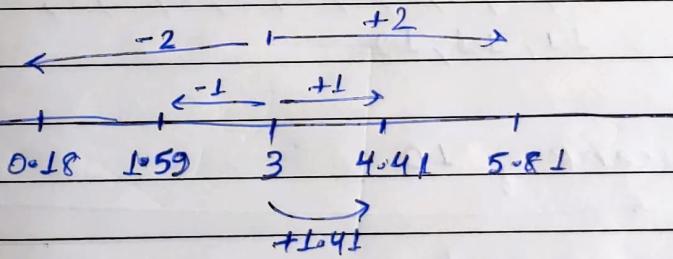
$\{ 1, 2, 3, 4, 5 \}$

$$\bar{x} = \frac{1+2+3+4+5}{5} = 3$$

$$\sigma^2 = \frac{10}{5} = 2$$

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

\* Find out how much away a number is present from the mean



\* where does 4 fall

Ans → 1 standard deviation towards right

variance  $\uparrow \Rightarrow$  Standard deviation  $\uparrow$

## # Percentiles and Quartiles

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

% age of even numbers = 50%

Defn : A percentile is a value below which certain percentage of observations lie.

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

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

(a) Percentile rank of 10

Percentile rank of  $n$  =  $\frac{\text{No. of value below } n}{n}$

$$= \frac{16}{20} = 80 \text{ percentile}$$

→ Dataset should be sorted while calculation  
 →

(a) what is the value that exists at 25 percentile

$$\text{value} = \frac{\text{Percentile}}{100} \times \frac{(n+1)}{2} \quad (\text{Even } n)$$

$$= \frac{25}{100} \times \frac{(19+1)}{2} = 5^{\text{th}} \text{ index}$$

5

$n \rightarrow$  Total no. of values  
in the list

elsewhere which is  
variations like

person has got  
of the entire

index 0 1 2 3 4 5 5 6 7 ...  
0 1 2 3 4 5

\* Index will start from 0

$$\text{value} = \frac{\text{Percentile}}{100} \times n \quad (n - \text{odd})$$

8, 8, 18, 9, 9, 10,

value at 95 percentile

$$\frac{95}{100}$$

below x

percentile

acceleration

## # 5 number summary

- (1) minimum
- (2) maximum
- (3) First quartile (25 percentile) (Q1)
- (4) Median
- (5) Third quartile (75 percentile) (Q3)

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

[ Lower fence  $\rightarrow$ , Higher fence ]

$$\text{Lower fence} = Q1 - 1.5(\text{IQR})$$

$$\text{Higher fence} = Q3 + 1.5(\text{IQR})$$

IQR

$$\text{IQR} = Q3 - Q1$$

$\downarrow$   
( Inter quartile Range )

$$Q1 = \frac{25}{100} \times (n+1)$$

$$= \frac{25}{100} \times 21 = 5.25 \rightarrow \text{Index}$$

SINCE 5.25 index doesn't exist to avg. of values

YUVRAJ

at index 586 will be taken,

$$\therefore \frac{3+3}{2} = 3$$

$$\therefore Q_1 = 3$$

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

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

$$IQR = Q_3 - Q_1 = 7.5 - 3.25 \\ 7.5 - 3 = 4.5$$

$$\text{Lower fence} = Q_1 - 1.5(IQR) \\ = 3 - 1.5(4.5) = -3.65$$

$$\text{Upper fence} = Q_3 + 1.5(IQR) \\ = 7.5 + 1.5(4.5) = 14.25$$

Minimum = 1

$Q_1 = 3$

Median = 5

$Q_3 = 7.5$

Maximum = 9 ( $-3.65$ ) Lower fence

Any value outside lower & upper fence is called outliers

If more outliers are present

