

# 2nd Day of Statistics

## Agenda

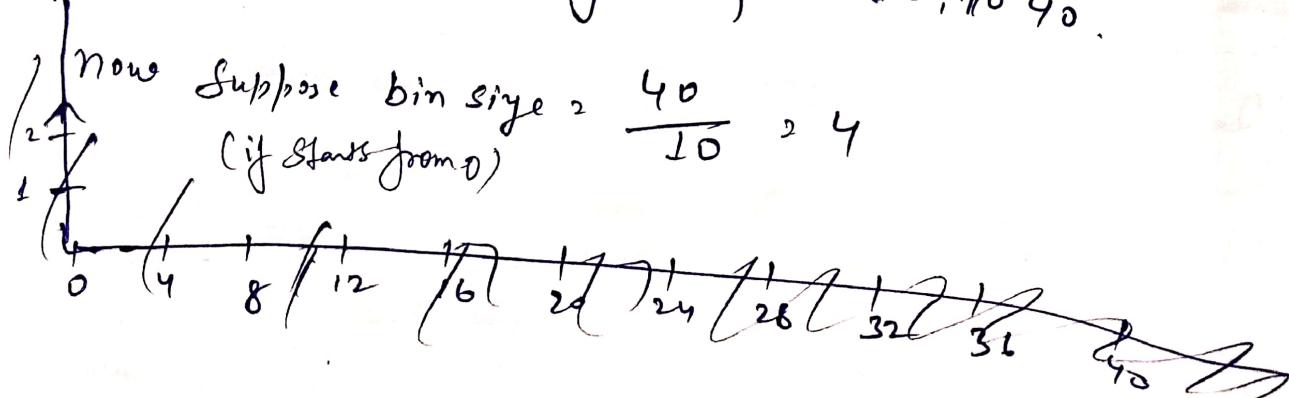
- ① Histogram
- ② Measure of Central Tendency
- ③ Measure of Dispersion
- ④ Percentile and Quartiles
- ⑤ 5 Number Summary (Box Plot)

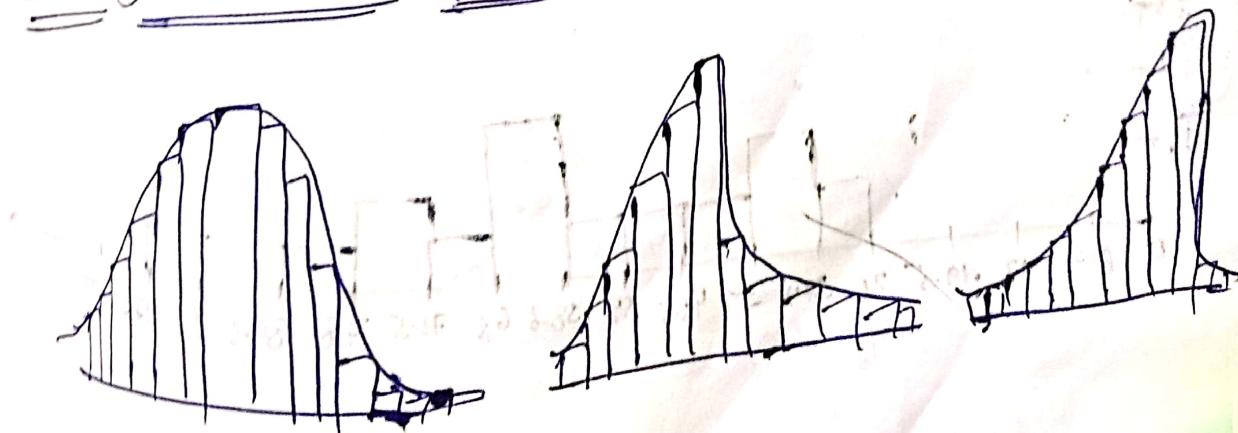
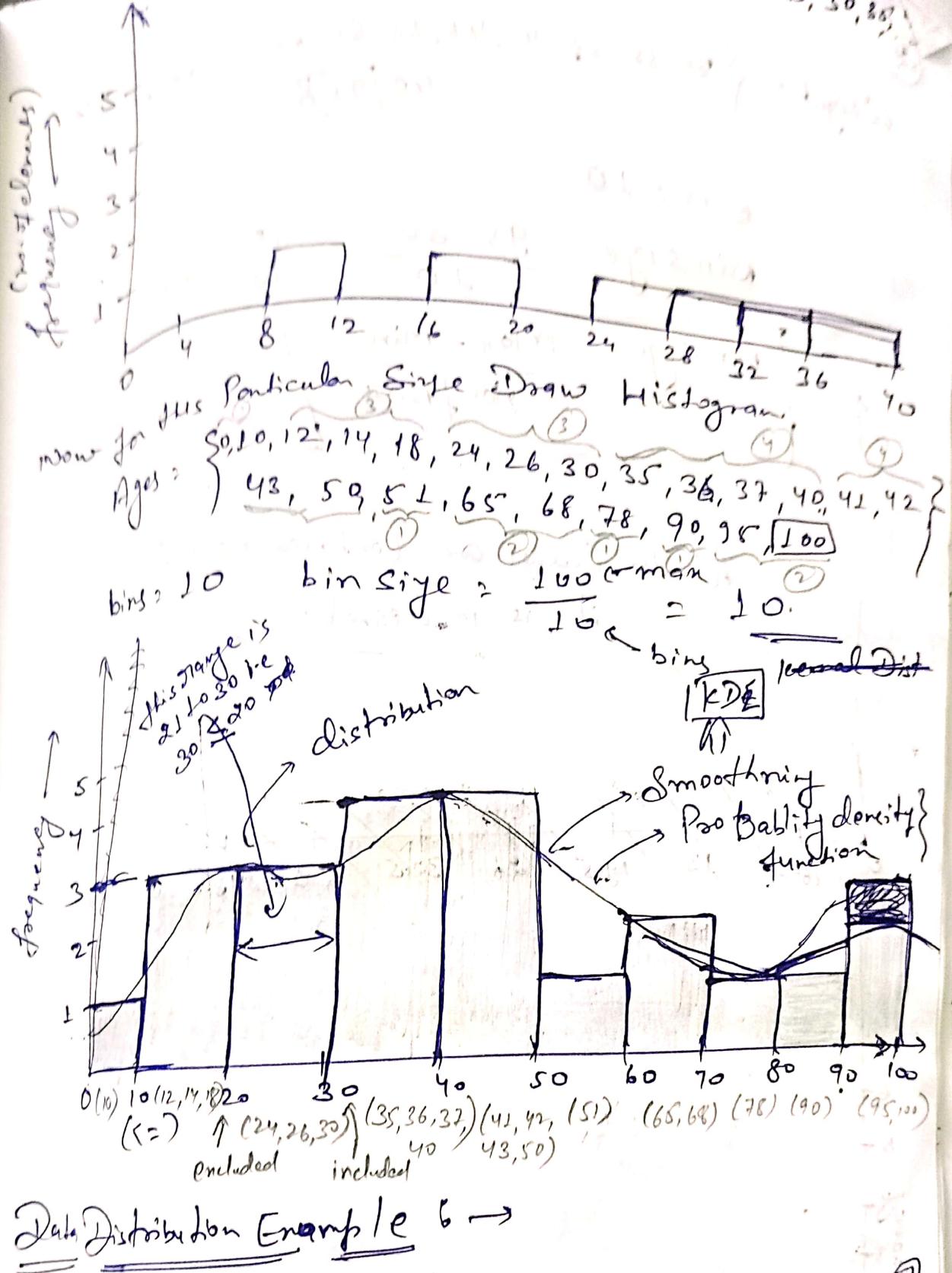
### ① Histogram →

Ages = {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 nos. (if not sorted)
  - ② Bins → No. of group. → [0, 10, 20, 25, 30, 35, 40]
  - ③ Bin Size :- Size of Bins      min      if starts from 0  
                                        max  
                                        bin size =  $\frac{40 - 10}{10} = \frac{30}{10} = 3$   
(divide it into 10 groups)
- 

But in Our Question Range is from 10, to 40.





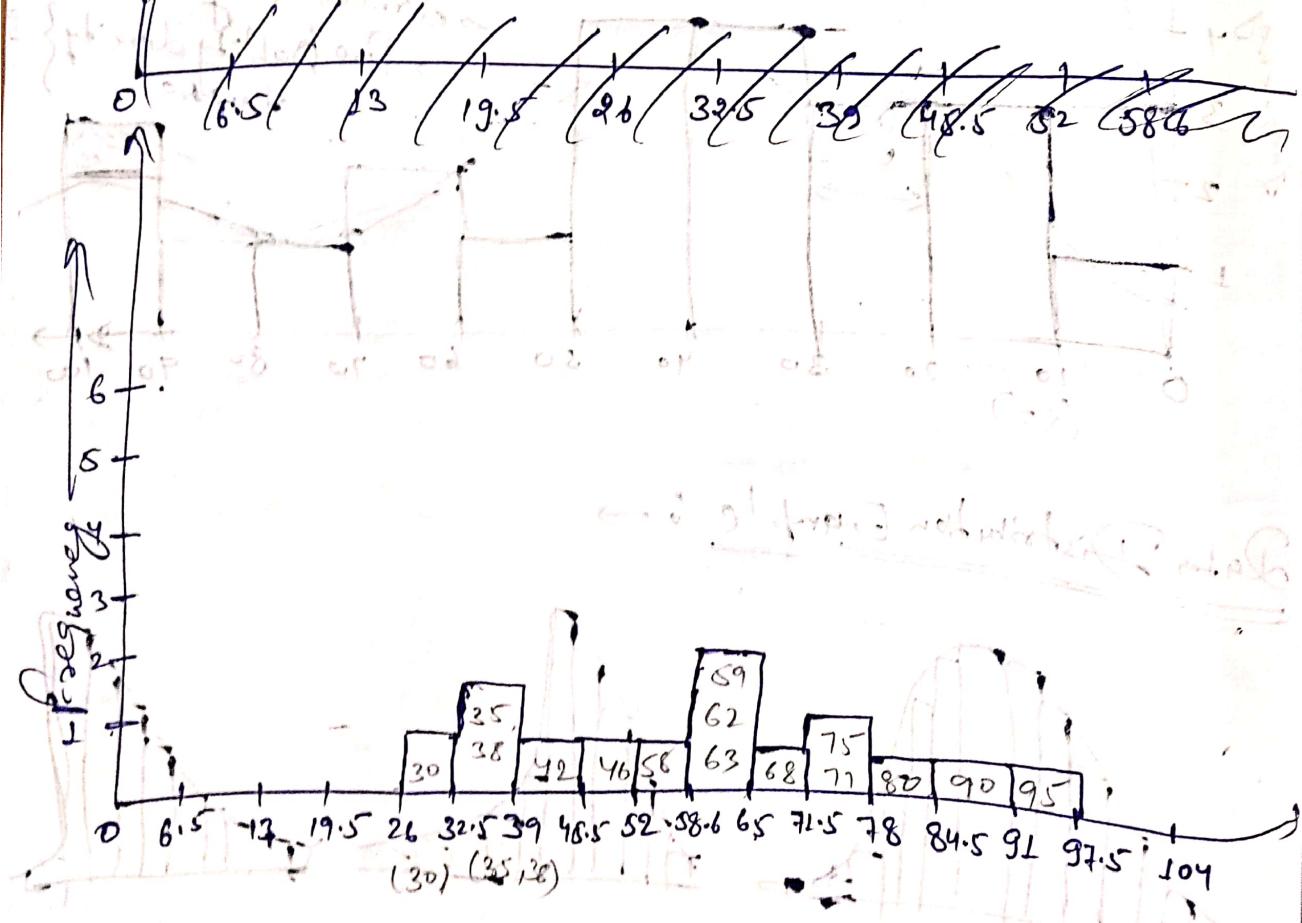
Q2  
Weight = { 30, 35, 38, 42, 46, 56, 59, 62, 63, 68, 78, 77, 80, 90, 95 } R

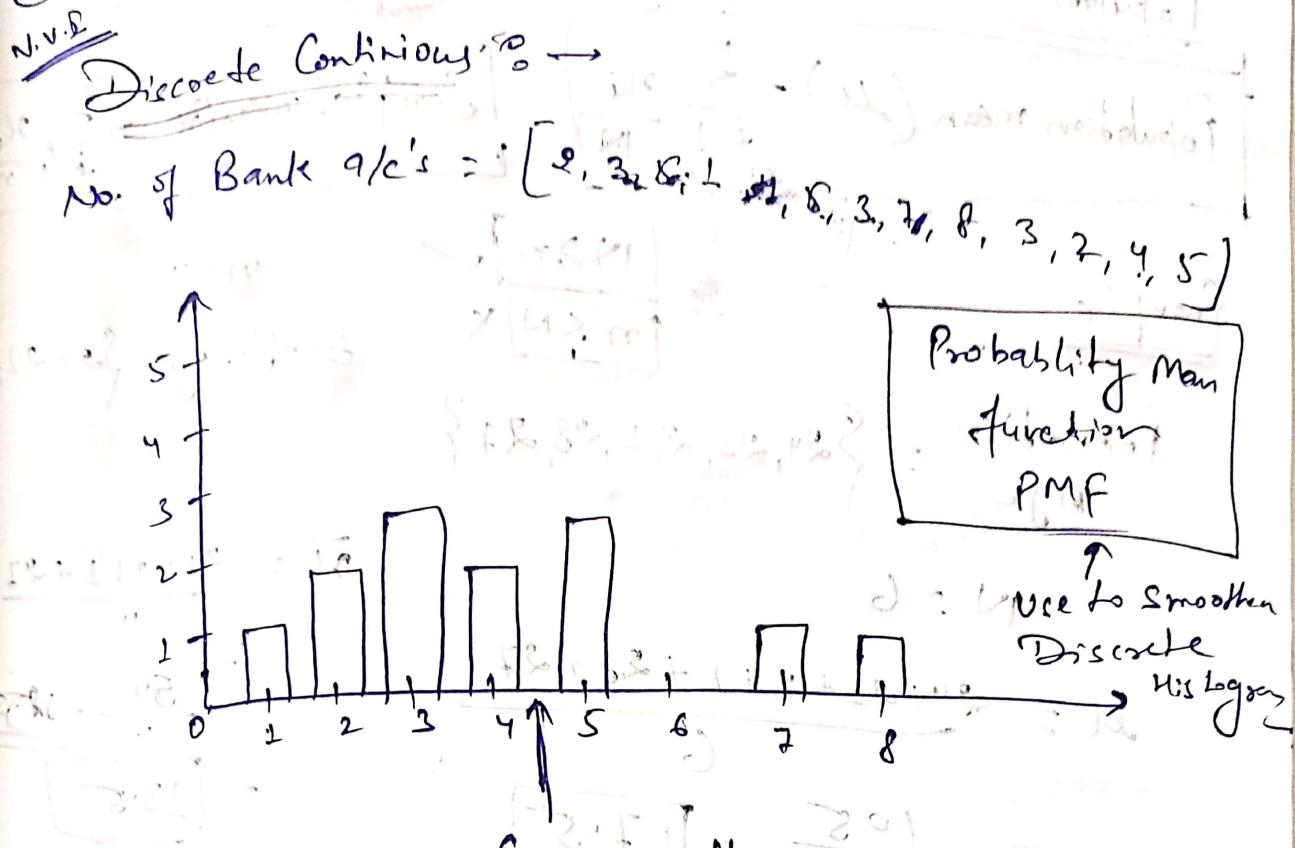
Bins = 10

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

$$\boxed{\text{Bin size} = \frac{\text{Max} - \text{Min}}{\text{no. of Bins}}}$$

All the values are continuous here. Since it is in decimal,





Graphs are there

bcz all the no's are discrete.

→ for continuous variable there is no decimal points.

PDF (Probability Density function) : → Use this method to

② Smoothen the continuous Histogram

PMF (Probability Mass Function) : → Use this method to

Smoothen Discrete Histogram.

for Discrete Variable

③ Measure of Central Tendency → (CT)

- ① Mean
- ② Median
- ③ Mode

A measure of CT is a single value that attempts to describe a set of data identifying the central position

④ Mean  $x = \{1, 2, 3, 4, 5\}$

$$\text{Avg} | \text{Mean} = \frac{1+2+3+4+5}{5} = 15/5 = 3$$

Population (N)

$N \geq n$

Sample (n)

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

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

Example

$$\begin{cases} N \geq n & \checkmark \\ n \leq N & \times \end{cases}$$

$n = 4$

Sample Age: {24, 23, 21, 27}

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

$$\text{Number of } N = 6$$

$$\bar{x} = \frac{24 + 23 + 21 + 27}{4}$$

$$\begin{aligned} \mu &= \frac{24 + 23 + 21 + 28 + 27}{6} \\ &= \frac{105}{6} \quad \boxed{17.5} \end{aligned}$$

$$\begin{aligned} \bar{x} &= \frac{54}{4} = 13.5 \\ &= \boxed{13.5} \end{aligned}$$

$$\begin{cases} \mu > \bar{x} \\ \mu < \bar{x} \end{cases}$$

Both is possible

Imp. note:  $\text{NaN} \Leftarrow \text{NULL Value}$

→ Practical Application ↗

Age	Salary	Family size
-	-	-
-	-	-
$\text{NaN}$	-	delete
-	-	-
-	$\text{NaN}$	-
-	$\text{NaN}$	$\text{NaN}$
-	$\text{NaN}$	-
$\text{NaN}$	-	-

Rather than deleting an entire row we can just  
the mean so that we  
won't lose any info.  
At last mean is  
Center of Data.

So whenever there is  $\text{NaN}$  value just



Eg:-

Age	Salary
24	45
28	50
29	NAN
MAIN	
31	75
36	80
NAN	NAN
80	200
24 + 28 + 29 + 31 + 36	

Mean Age =

(different from data)  
outlier

Mean Age =

$\frac{148}{5} = 29.6$

outliers (different from the other data)

Add here

now what will be the Avg. ?

$M_{Age} = \frac{148 + 80}{6}$

$M_{Salary} = \frac{45 + 50 + 60 + 75 + 80}{5} = 62$

$M_{Salary} = \frac{310 + 200}{6} = 85$

$= \frac{310}{5} = 62$

⇒ Because of outliers there is change in mean of both  $M_{Age}$  and  $M_{Salary}$ . Hence to prevent this we use 5th is called as Median.

ii) Median :-

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

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

$\bar{x} = 3 \xrightarrow{\text{Big change}} = \frac{115}{6} = 19.15$

Hence when there

is big change like this mean it is  
bit of outlier



## Steps to find out Median

- i) Sort the numbers.
- ii) Find the central no. If the no. of elements are even we find the Avg. of central elements
- iii) If the no. of elements are odd then we find the central ele.

Eg:-

i)  $\{1, 2, 3, 4, \boxed{5, 6}, 7, 8, 9, 10, 11, 12, 13\}$  ← Already sorted

No. of elements = 12 (even no.)

Hence find the Avg. of central ele. which is

i.e.  $\frac{5+6}{2} = \frac{11}{2} = 5.5 = \text{median}$

Hence Median = 5.5

mean =  $\frac{1+2+3+4+5+6+7+8+10+11+12+13}{12} = \frac{86}{12} = 7.16$

mean = 7.16

Hence whenever there is an outlier we use median.

and if there is no outlier we use mean

No outliers = mean  
Outliers = median

③ Mode → { Most frequent occurring elements }

Eg:- i)  $\{1, 2, 2, \boxed{3, 3, 3}, 4, 5, 2\}$

Mode = 3

ii)

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

Mode = 2, 3

## Practical Application of Mode

Eg:- Dataset { Categorical Variable }

Type of flower

Lily

Sunflower

Rose

NAN ← Rose

Rose

Sunflower

Rose

NAN ← Rose

Since Rose has Repeated man.  
no. of times i.e 3 times Hence  
we will replace NAN with  
Rose.

Hence Mode mostly we  
use with Categorical Var.

Q: If there are two values in Mode then from which  
value shall I replace NAN with.

Sol/  
Eg:- Mode : { Sunflower, Rose }

You can replace with anyone of them it's your  
choice.

Q: If NAN is occurred more than mode What to do?

Ans: Delete that particular column. If 90% of Data  
is having NAN value you have to delete that  
particular column.

③ Measure of Dispersion →

i) Variance ( $\sigma^2$ ) → Spread of Data.

ii) Standard Deviation ( $\sigma$ )

i) Variance ( $\sigma^2$ ): → Talks about spread of Data.

Population Variance ( $\sigma^2$ )

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

Sample Variance ( $s^2$ )

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

Distance  
from the mean  
from every data point(s)

why is it (n-1)  
Answer



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

$$\mu = 3$$

Q Suppose we have two Dist<sup>r</sup>

$$(i) \{1, 2, 3, 4, 5, 6, 7, 8, 9, 10\} \quad (ii) \{1, 2, 3, 4, 80, 60, 70\}$$

Find Variance

Find Variance

Now which Variance will be higher (i) or (ii) ?

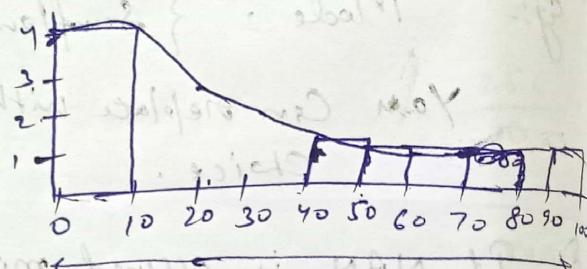
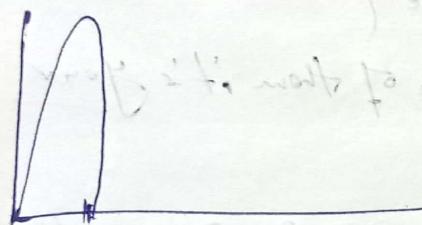
~~Second (ii) →~~

~~Variance of (ii)~~

~~Variance of (i)~~

~~Variance of (ii)~~

~~Variance of (ii)~~



Range is len,

Range is high

Hence  $\text{Variance of } (i) < \text{Variance of } (ii)$

Let's calculate:

$$\mu_i = \frac{1+2+3+4+5+6+7+8+9+10}{10}$$

$$\mu_{ii} = \frac{1+2+3+4+80+60+70+100}{8}$$

$$= \frac{55}{10} = 5.5$$

$$= \frac{290}{8} = 35$$

$$\sigma_i^2 = \frac{(1-5.5)^2 + (2-5.5)^2 + (3-5.5)^2 + (4-5.5)^2 + (5-5.5)^2 + (6-5.5)^2 + (7-5.5)^2 + (8-5.5)^2 + (9-5.5)^2 + (10-5.5)^2}{10}$$

$$= \frac{(4.5)^2 + (8.5)^2 + (10)^2 + (6.5)^2 + (1.5)^2 + (0.5)^2 + (-0.5)^2 + (-3.5)^2 + (-5.5)^2 + (-2.5)^2}{10}$$

$$20.25 + 12.25 + 6.25 + 2.25 + 0.25 + 0.25 \\ + 0.25 + 6.25 + 12.25 + 20.25$$

$$\frac{82.5}{10} = \underline{\underline{8.25}}$$

(ii)  $\sigma^2 = \frac{(1-35)^2 + (2-35)^2 + (3-35)^2 + (4-35)^2 + (50-35)^2 + (60-35)^2 + (70-35)^2 + (100-35)^2}{8}$

$$= \frac{(-34)^2 + (-33)^2 + (-32)^2 + (-31)^2 + (-25)^2 + (35)^2 + (65)^2}{8}$$

$$= 1156 + 1089 + 1024 + 961 + 225 + 625 \\ + 1225 + 4225$$

$$= \frac{10530}{8} = \underline{\underline{1316.25}}$$

~~(i)~~  $\sigma_{(i)}^2 < \sigma_{(ii)}^2$

which is  $8.25 < 1316.25$

### Example (ii)

①  $\{1, 2, 3, 4, 5\}$

$$n = 5$$

$$\sigma_{(i)}^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (4-3)^2 + (5-3)^2}{5} \\ = \frac{(-2)^2 + (-1)^2 + (0)^2 + (-1)^2 + (2)^2}{5} \\ = \frac{4+1+0+1+4}{5} = \frac{10}{5} = 2$$

②  $\{1, 2, 3, 4, 5, 6, 80\}$

$$n = 7$$

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

$$= \frac{(-13.4)^2 + \cancel{(-9)^2} + (-9)^2 + (-8)^2 + (-84)^2}{7}$$



$$\begin{aligned} & \frac{(13.4)^2 + (12.4)^2 + (11.4)^2 + (10.4)^2}{7} + \\ & + (8.8)^2 + (6.8)^2 \end{aligned}$$

$$\begin{aligned} & = 179.56 + 153.76 + 129.96 + 108. \\ & + 90.25 + 72.25 + 4303.36 \end{aligned}$$

$$\frac{7}{7} = 719.61$$

Now we can see:  $\sigma^2 = \frac{(x_1 - \bar{x})^2 + (x_2 - \bar{x})^2 + \dots + (x_n - \bar{x})^2}{n}$

$\sigma^2 < 719.61$

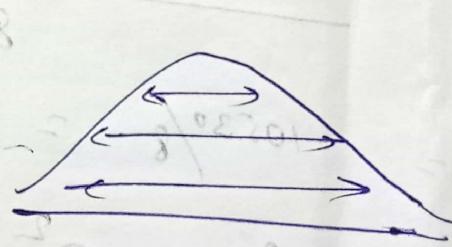
So can we conclude that as Variance keeps on increasing  
data set also increases.  
i.e. Variance ↑↑ Spread ↑↑



$$\sigma_{(i)}^2$$

Lower spread ↓

Lower Variance ↓



$$\sigma_{(ii)}^2$$

Higher spread (Width)

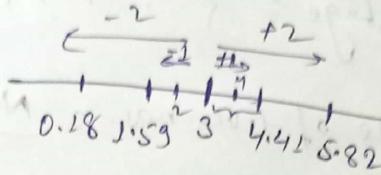
Higher Variance

ii) Standard Deviation ( $\sqrt{\sigma^2}$ ) →

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

$$\sigma^2 = \frac{2}{5} = 0.4$$

$$\sigma = \sqrt{0.4} = \underline{\underline{1.41}}$$



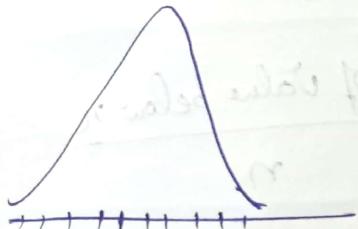
It tells that how many standard deviation away from the mean (each dist<sup>r</sup>) it falls.

Q: 4 is falling how many standard deviation away from the mean?

A: Within the  $\pm 1$  Range

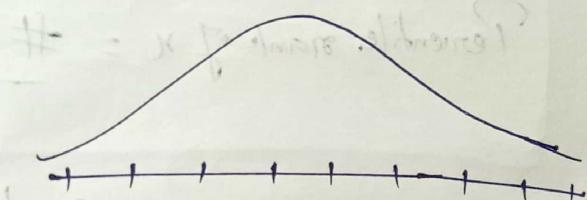
Similarly, 2 is falling  $-1$  away from the mean.

Ex:-



$$\sigma^2$$

$$\sqrt{\sigma^2}$$



$$\sigma^2 \text{ (Variance)}$$

$$\sqrt{\sigma^2} \text{ (S.D.)}$$

Variance can lie in both right or left. Hence we want to know which side if lie right side or left side for that we use Standard deviation.

#### Q) Percentile & Quantiles: →

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

$$\text{Percentage of Even nos.} = \frac{\text{No. of Even nos.}}{\text{Total no. of nos.}} = \frac{4}{8} = 50\%$$

⇒ Percentiles: → GATE, CAT, IELTS, SAT, GRE, JEE, NPTEL

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

99 Percentile = it means that the person has got better marks than 99% of entire student.

Index: - 0, 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18  
Dataset: → 2, 2, 3, 4, 5, 5, 5, 6, 7, 8, 8, 8, 8, 9, 9, 10, 11, 11, 11

What is the percentile rank of 10

Percentile rank of  $x = \frac{\# \text{ no. of value below } x}{n}$

$$= \frac{16}{20} = 0.8 = 80 \text{ percentile}$$

All the elements should be

in ascending order if not

just sort it.

Q: Can 2 people get same percentile?

Ans: YES.

$$\text{Percentile of 8} = \frac{9}{20} = 0.45 = 45 \text{ percentile}$$

This means that 45% of people are less than 8.  
For all 8 it will be some i.e. 45 percentile.



Percentile of 6  $\Rightarrow \frac{6}{20} = 0.35$ ,  $\underline{\underline{0.35 \text{ percentile}}}$

Percentile of 9  $\Rightarrow \frac{9}{20} = 0.45$ ,  $\underline{\underline{0.45 \text{ percentile}}}$

Q) What is the value that exists at 25 percentile?

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

for even 'n' multiply with  $(n+1)$  for odd 'n' multiply with  $n$

$$\frac{25}{100} \times 25 \rightarrow 5^{\text{th}} \text{ index.}$$

Hence O/P = 5  $\leftarrow$  5<sup>th</sup> index value is 5.

Q) What is value that exists at 95 percentile?

$$\text{Value} = \frac{95}{100} \times 25$$

$$\text{Value} = \frac{95}{100} \times 25 = 19.95 \text{ index}$$

$$\text{Hence O/P} \rightarrow 19.95 \text{ index}$$

⑤ 5 number Summary  $\rightarrow$

i) Minimum

ii) First Quartile (25 percentile) (Q1)

iii) Median

iv) Third Quartile (75 percentile) (Q3)

v) Maximum

All these we use to remove the outliers.

↓  
create Box Plot

Dataset 2:  $\{1, 2, 2, 2, 3, \boxed{3}, \boxed{3}, 4, 5, 5, 6, 6, 6, 6, 7, \boxed{8}, 8, 9, \boxed{27}\}$

How to know 27 is an outlier for that

We are creating a fence. Called

[lower fence]  $\longleftrightarrow$  [higher fence]  
they found out after many experiments

$$\text{Lower fence} = Q_1 - 1.5(\text{IQR}) \quad \text{where, IQR} = Q_3 - Q_1$$

(It specifies 1.5 times deviation from the mean)

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

Inter Quartile Range

$$Q_1 = \frac{25}{100} \times (n+1) = \frac{25}{100} \times 21 = 5.25^{\text{th}} \text{ index}$$

$$Q_1 = \frac{3+3}{2} = \boxed{3}.$$

$$Q_3 = \frac{75}{100} \times 21 = \frac{63}{4} = 15.75^{\text{th}} \text{ index}$$

$$Q_3 = \frac{7+8}{2} = 15/2 = \boxed{7.5}$$

Since we don't have 5.25 index so we will take the avg. of 5<sup>th</sup> & 6<sup>th</sup> index i.e.

Hence take the Avg. of 15<sup>th</sup> & 16<sup>th</sup> index i.e.

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

$$= 3 - 1.5(7.5 - 3)$$

$$= 3 - 1.5 \times 4.5$$

$$= \underline{\underline{-3.75}}$$

$$\begin{aligned}
 \text{Higher fence} &= 7.5 + 1.5(7.5 - 3) \\
 &= 7.5 + 1.5 \times 4.5 \\
 &= 7.5 + 6.75 \\
 &= \underline{\underline{14.25}}
 \end{aligned}$$

Hence all the values be in the range of  
 - 3.75 to 14.25

Hence  $d_7$  doesn't lie in this range hence it is an outlier.

Hence we can remove the outlier.

Middle element

Now,

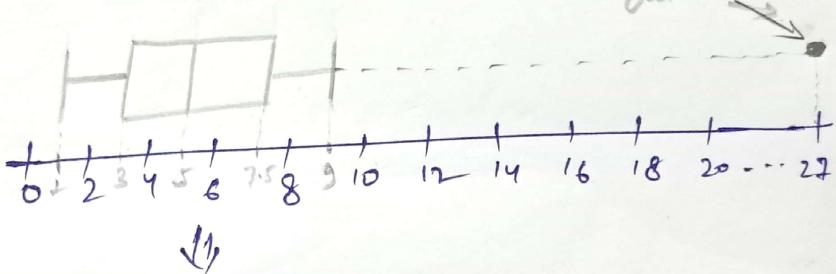
$$\{1, 2, 2, 2, 3, 3, 3, 4, 8, \boxed{5}, 5, 6, 6, 6, 6, 7, 8, 8, 9, \cancel{12}, \cancel{7}\}$$

remove outlier

Now find

- i) Minimum = 1
- ii)  $d_1 = 3$
- iii) Median = 5
- iv)  $d_3 = 7.5$
- v) Maximum = 9

Now create Box Plot



We create Box Plot to  
 Treat outliers.