# Topics

① Histogram
② Measure of central tendency
③ Measure of Dispersion
④ Percentiles & Quartiles
⑤ 5 Number Summary (Box plot)

## i) Histogram

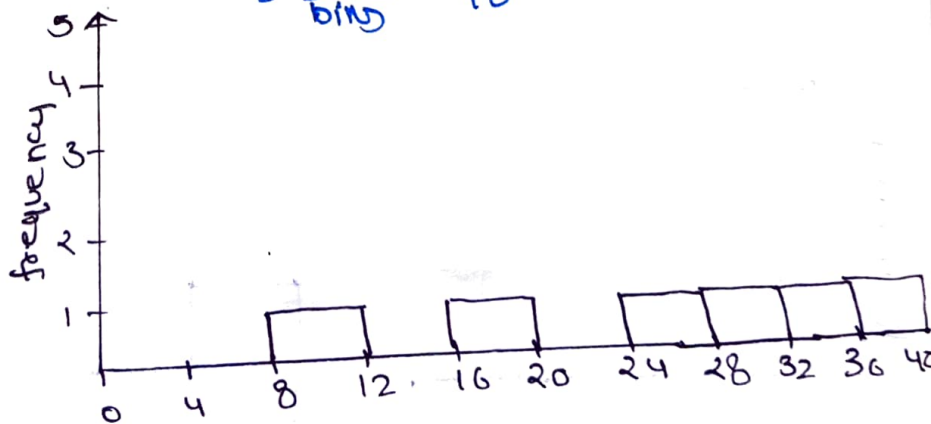Ages = {10, 12, 18, 24, 26, 30, 35, 36, 37, 38, 40, 41, 42, 43, 50, 51, 65, 68, 78, 90, 95, 100}

① Sort the no
② Bins → no. of groups
③ Bins size → size of Bins

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

min = 10        max = 40

bins = 10 → means 10 group b/w 0 to 40

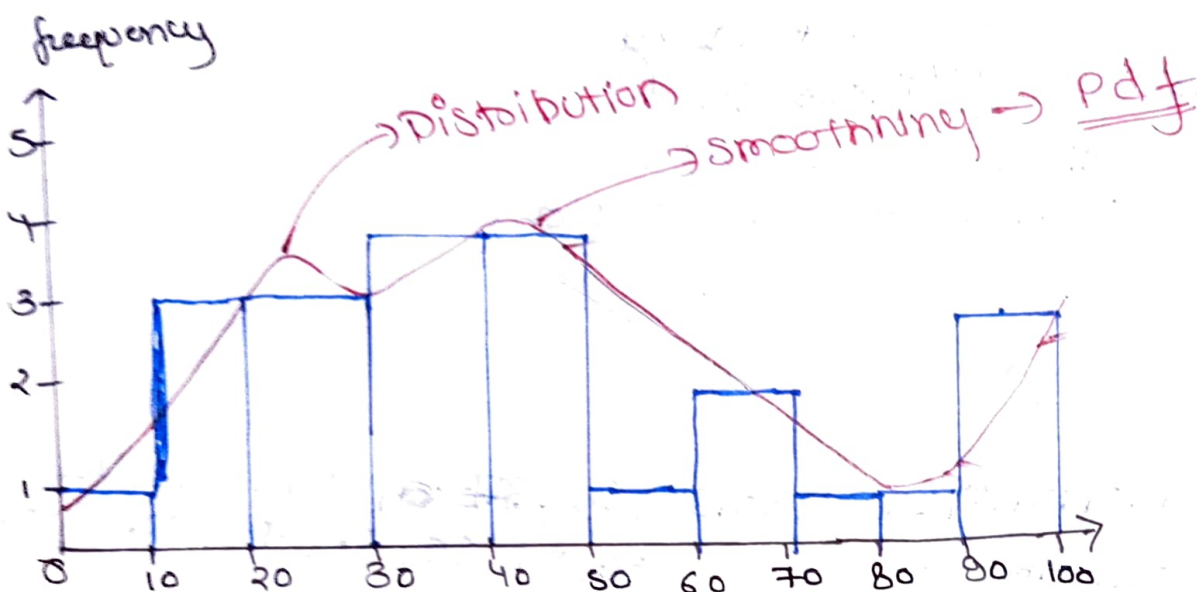$$= \frac{max}{bins} = \frac{40}{10} = 4$$

bing = 10          bin size = $\frac{100}{10}$ = 10 → 10 group b/w
                                                    0 to 100

frequency



→ Distribution          → smoothning → pdf

(histogram with axis 0, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100)

**eg** weight = { 30, 35, 38, 42, 46, 58, 59, 62, 63, 68, 75, 80,
                                                              77
                  90, 95 }

bin size = $\frac{65}{10}$ = 6.5

bins = 10

frequency



(x-axis: 30  36.5  43  49.5  56  62.5  69  75.5  82  88.5  95)

① Discreate & continous
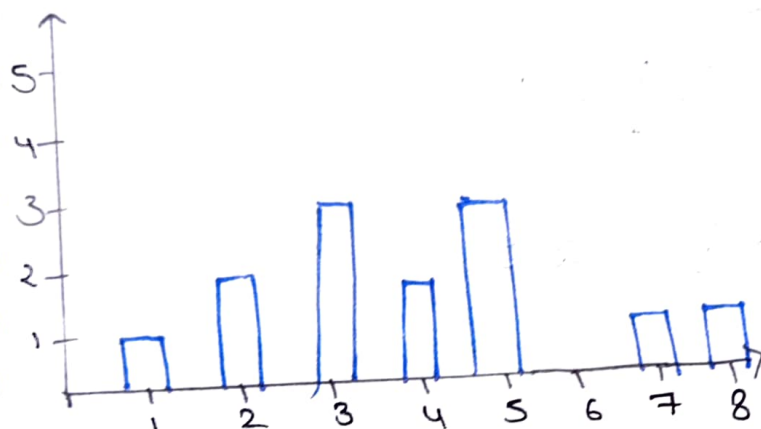
Ranks $= \{1, 2, 3, 4, 5, 6\}$

No. of Banks accounts $= [2, 3, 5, 1, 4, 5, 3, 7, 8 3, 2, 4, 5]$

Discrete values
↓
Smoothing
using
↓
Probability mass
function



Pdf = Probability density func → continous

pmf = Probability mass func → discreate

→ Measure of central Tendency

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

① Mean

eg $n = \{1, 2, 3, 4, 5\}$

→ specifying the central position

Avg/mean $= \dfrac{1+2+3+4+5}{5} = \dfrac{15}{5} = 3$

Population (N)

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

Sample (n)

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

$N \gg n$ → Population always greater than sample

ey

Population $= \{24, 23, 2, 1, 28, 27\}$    $\boxed{N=6}$
Age

$$\mu = \frac{24 + 23 + 2 + 1 + 28 + 27}{6}$$

$$= 17.6 \implies \boxed{\mu = 17.6}$$

Sample $= \{24, 2, 1, 27\}$    $\boxed{n = 4}$
Age

$$\bar{n} = \frac{24 + 2 + 1 + 27}{4} = 13.5$$

$$\boxed{\bar{n} = 13.5}$$

\* case
$\left. \begin{array}{l} \mu \geq \bar{n} \\ \bar{n} \geq \mu \end{array} \right\}$ these Situation can happen

→ Ponactical Application   (featocre Engy)

ey

| Age | salary | family size |
|---|---|---|
| — | — | — |
| — | — | — |
| NAN | — | — |
| — | NAN | — |
| — | — | NAN |
| — | — | — |
| — | — | — |
| NAN | — | — |

→ Donoping Row having NAN
    ↓
  loqq of Info
    ↓ instead
NAN are creplaced by Mear

## Eg

| Age | Salary |
|-----|--------|
| 24 | 45 |
| 28 | 50 |
| 29 | NAN |
| NAN | 60 |
| 31 | 75 |
| 36 | 80 |
| NAN | NAN |

Age Mean = 29.6    Salary mean = 62

new Data → 80    200    $\xrightarrow{\text{new mean}}$ Age mean = 38   Salary Mean = 85

→ Outliers

Outliers :- an observation that lies at abnormal distance from other values in a random sample from a population

## ② Median

$\{1, 2, 3, 4, 5\}$       $\{1, 2, 3, 4, 5, 100\}$
                                              ↑
$\bar{x} = 3$ ────→ $\bar{x} = 19.16$    outlier

→ Steps to find out median
  ① Sort the numbers
  ② find the central Number
    ① if no of Elements are even we find the avg. of central Elements
    ② add no. of Element ──→ find central Elements

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

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

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

Median $= 5$

\* no outlier $\longrightarrow$ Mean

with outlier $\longrightarrow$ Median

③ Mode $\rightarrow$ most fequent occuring element

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

mode $= 3$

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

mode $= 2, 3$

eg Types of flowers

Lily, Sunflower, Rose, NAN, Rose, Sunflower, Rose, NA

\* ~~mode~~

\* mode is mostly use with categorical variable

* Measure of Dispersion
  ① Variance $(\sigma^2)$ ← Spread of data
  ② Standard deviation $(\sigma)$

① Variance

Population Variance $(\sigma^2) = \sum\limits_{i=1}^{N} \dfrac{(x_i - \mu)^2}{N}$

Sample Variance $(S^2) = \sum\limits_{i=1}^{n} \dfrac{(x_i - \bar{x})^2}{n-1}$
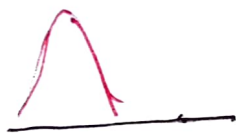
• $x_i - \mu \rightarrow$ deviation from mean

eg

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

$\mu = 3$

$\sigma^2 = 2$

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

$\mu = 14.4$

$\sigma^2 \neq 9.10$

- As variance keeps on inc spread inc

*

$\sigma^2$     <     $\sigma^2$

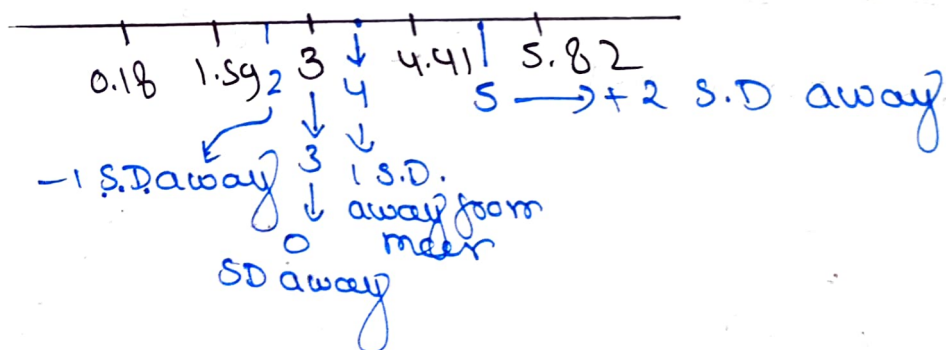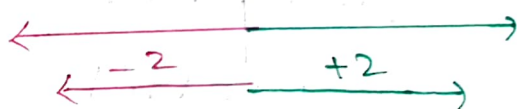② Standard deviation $\left(\sqrt{\sigma^2}\right)$

Ex $\quad 1, 2, 3, 4, 5$

$$\mu = \frac{1+2+3+4+5}{5} = \frac{15}{5} = 3$$

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

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

$$\boxed{\sigma^2 = 2} \qquad \boxed{\sigma = 1.41}$$



```
         ←            -2          →
         ←   -2          +2   →
              ←  -1   +1  →
```

```
  |    |    |  ↓    |    |
 0.18 1.59 2 3  4.41  5.82
          3  4    5 →+2 S.D away
          4
   -1 S.D away  3   1 S.D.
          ↓ away from
          O   mean
        SD away
```

→ How many standard deviation away a no falls from the mean.

# * Percentile & Quartiles

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

Percentage of even no. $= \dfrac{\text{no. of even no}}{\text{total no. of no}} = \dfrac{4}{8} = 0.5 = 50\%$

→ Percentile

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

eg

99 percentile → the person has got better marks than 99% of the entire students

eg

Dataset = 2, 2, 3, 4, 5, [5], 6, 7, 8, 8, 8, 8, 8, 9, 9, 10, 11, 11, 12

(0th index)    (5th index)

- what is the percentile rank of 10

$$\boxed{\text{Percentile Rank of } n = \dfrac{\text{no. of value below n}}{n}}$$

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

- Percentile rank of 8

$$= \dfrac{9}{20} = 45 \text{ percentile}$$

↓

45% of obs are below 8

* ## S number Summary

① Minimum
② first Quartile [25 percentile] [Q1]
③ Median
④ Third Quartile [75 percentile] [Q3]
⑤ Maxium

Remove the outlier
⇓
Box plot

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

note :- a small no can also be a outlier

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

$$\boxed{IQR = Q_3 - Q_1}$$
↓
Inter Quartile Range

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

$$[\text{Lower fence} \longleftrightarrow \text{Higher fence}]$$

$Q_1 = \dfrac{25}{100} * (n+1) = \dfrac{25}{100} \times 21 = 5.25 \Rightarrow Indx = 3$

$Q_3 = \dfrac{75}{100} * 21 = 15.75 = \dfrac{8+7}{2} = 7.5$

Lower fence $= 3 - 1.5(4.5) = -3.65$

Higher fence $= 7.5 - (1.5)(4.5) = 14.25$

Lower fence                    ←——————→  Higher fence
= -3.65                                   14.25

all values should lie b/w
—— this

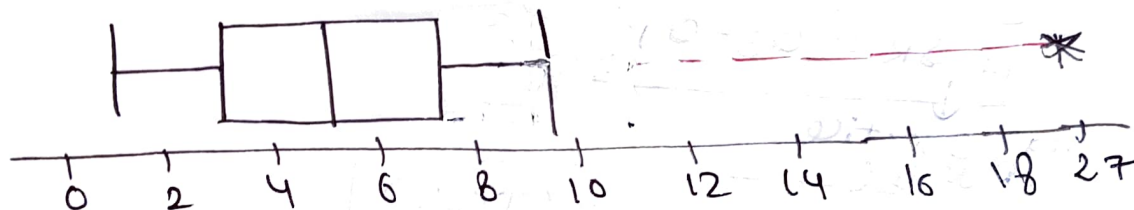→ 27 not lies b/w fence   so remove it

① $min = 1$          ③④ $q_3 = 7.5$
② $q_1 = 3$          ⑤ $max = 9$
③ $mean = 5$



Box plot
↓
To treat outliers