# Agenda:

1. Histograms
2. Measure of centeral Tendency.
3. Measure of Dispersion
4. Percentiles and Quartiles
5. 5 Number Summary (Box plot)

## Histogram:

$$Ages = \{0, 10, 12, 14, 18, 24, 26, 30, 35, 36, 37, 40, 41, 42, 43,$$
$$50, 51, 65, 68, 78, 90, 95, 100\} \text{ (Continuous values)}$$

steps to follow:

① . Sort the numbers

② . Bins → No.of groups

③ . Bin Size → Size of Bins
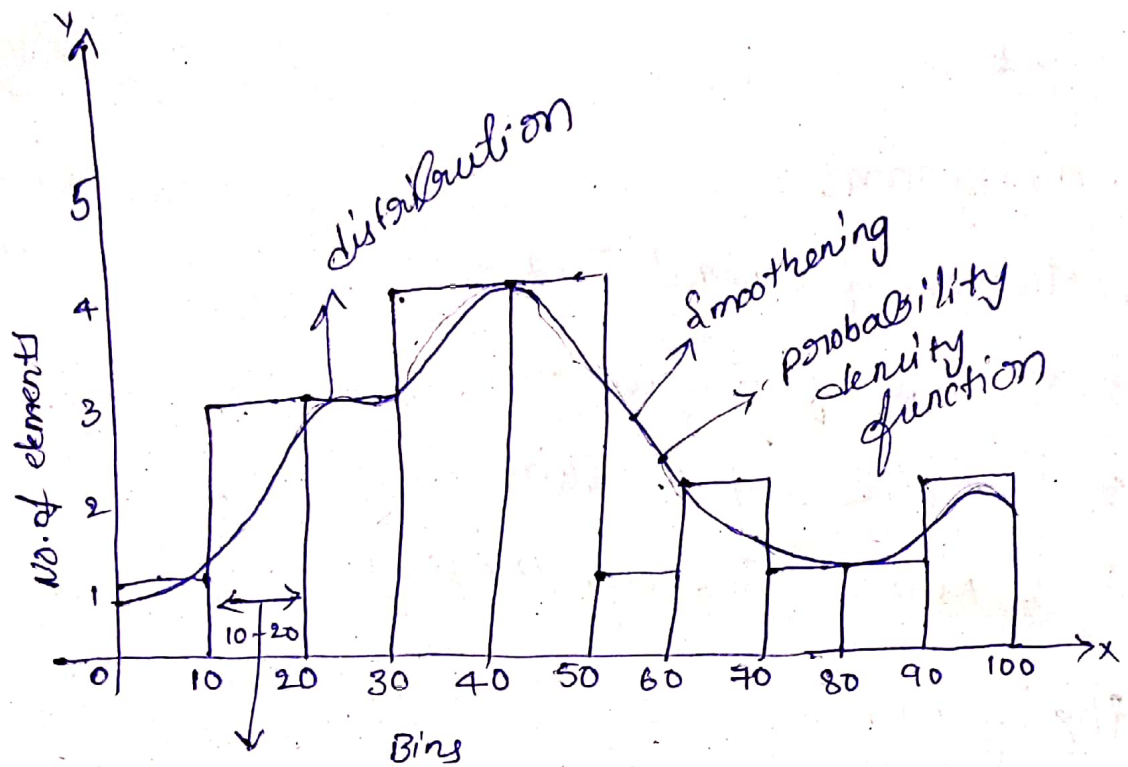
{ Bins are nothing
But groups

Binsize = groupsize }

→ the numbers are already sorted.

→ Bins = 10   (Here bins can taken by us

(Here I'm taking that how many no of groups(bins)
10 Bins)   we want, we can take)

→ Binsize = $\dfrac{max - min}{bins}$

$$= \dfrac{100 - 0}{10} = \dfrac{100}{10} = 10$$

So, Binsize = 10

distribution

smoothening

probability density function

No. of elements

10-20

10   20   30   40   50   60   70   80   90   100

Bins

between
10-20
we have
3 elements in the list,
Same for all.
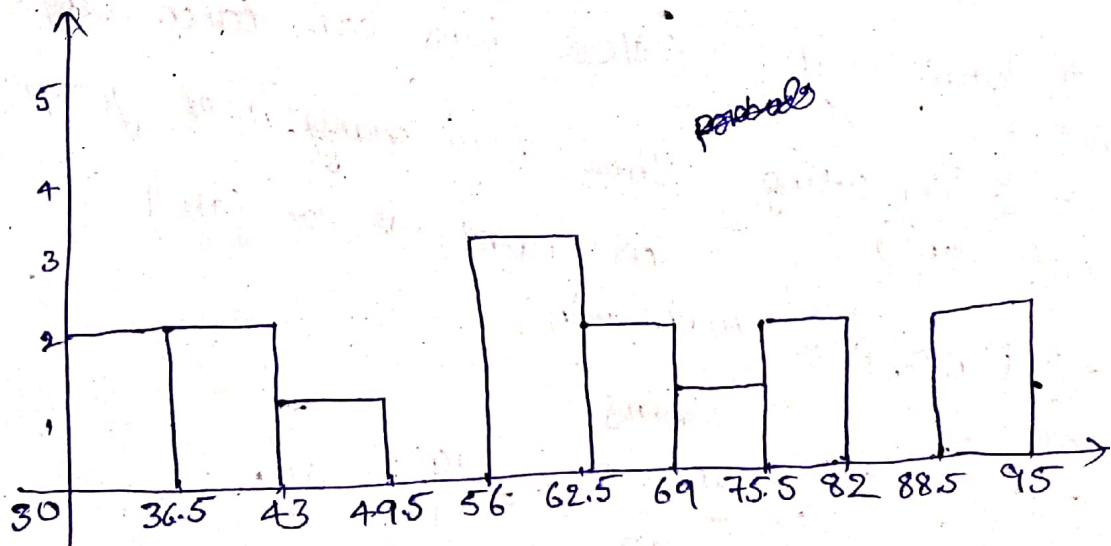
Ex-2

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

(continuous values)

bins = 10

$$binsize = \frac{95-30}{10} = \frac{65}{10} = 6.5$$



parabola

30   36.5   43   49.5   56   62.5   69   75.5   82   88.5   95

for Discrete value.

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



Probability
mass
function

Pdf = probability density function $\Big\} \rightarrow$ continuous

pmf = probability mass function $\Big\} \rightarrow$ discrete.

## * 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.

→ there are three methods to identify the central position

1). Mean
2) Median
3). Mode

Mean:

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

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

Population (N)     $N \geq n$     Sample (n)

$$\text{population mean } (\mu) = \left| \frac{\sum\limits_{i=1}^{N} x_i}{N} \right| \quad \text{sample mean}(\bar{x}) = \left| \frac{\sum\limits_{i=1}^{n} x_i}{n} \right|$$

population
age $= \{24, 23, 2, 1, 28, 27\}$

$\boxed{N = 6}$

sample age $= \{24, 2, 1, 27\}$

↳ Here 1 picked
4 values randomly
from population age.

$\boxed{n = 4}$

population mean $(\mu) = \dfrac{24 + 23 + 2 + 1 + 28 + 27}{6}$

$\boxed{\mu = 17.5}$

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

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

$\boxed{\mu \geq \bar{x}}$

~~$\bar{x} \geq \mu$~~

## Practical Application (Feature Engineering)

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

$(\mu)$
Age $= 29.6$
↓
$38$ wee for NAN value

Salary $= 62$
↓
85 for NAN app

If we remove total record due to NAN value then there is some loss of information.
So, by finding average (mean) we can give some values to NAN values.

# * Median:

In mean we have one problem. i.e., there may be a chance of occuring outliers, due to outliers the average value is changed.

for ex:

$\{1, 2, 3, 4, 5\}$ . $\{1, 2, 3, 4, 5, \overset{\text{outlier}}{\underset{\uparrow}{100}}\}$

$\Downarrow$            $\Downarrow$

for the mean $\bar{x} = 3$     for this $\bar{x} = 19.16$

so, we gone to median.

## Steps to find out median:

1. Sort the numbers

2. Find the central number.

$\Downarrow$

case ① : if the no. of elements are even we find the average of central elements

case ② : if the no. of elements are odd we find the central elements.

Ex:

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

median = 5

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

median $= \dfrac{5+6}{2} = 5.5$

**\* mode:** most frequent occuring element

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

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

mode = 3

mode = 2, 3

## practical application:

Dataset

Types of flower { categorical variable }

Lily

Sunflower

Rose

NAN → rose

Rose

Sunflower

Rose

NAN → Rose

Here rose is most frequently occuring, So, we replace NAN values with rose

## \* Measure of Dispersion

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

② standard deviation $(\sigma)$

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

## variance

population variance $(\sigma^2)$

$$\sigma^2 = \sum_{i=0}^{N} \frac{(x_i - \mu)^2}{N}$$

Sample variance $s$

$$s^2 = \sum_{i=0}^{n} \frac{(x_i - \bar{x})^2}{n-1}$$

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

$\mu = 3$

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

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

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

Ex:2

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

$\mu = 14 \cdot 4$

$$\sigma^2 = \frac{(1-14\cdot4)^2 + (2-14\cdot4)^2 + (3-14\cdot4)^2 + (4-14\cdot4)^2 + (5-14\cdot4)^2 + (6-14\cdot4)^2 + (80-14\cdot4)^2}{7}$$
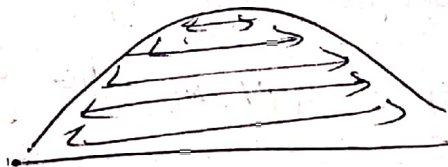
$$\boxed{\sigma^2 = 719 \cdot 10}$$

* Here we can observe that if variance increases
the spread of data is also increases

for example



when $\sigma^2 = 2$          when $\sigma^2 = 14 \cdot 2$

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

② standard deviation $\left(\sqrt{\sigma^2}\right) \Rightarrow \boxed{\sigma}$

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

$\mu = 3$

$\sigma^2 = 2$ ( previously we find out variance)

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

Here 1.41 is the standard deviation.

* Percentiles and Quantiles:

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

percentage of Even number $= \dfrac{no. \text{ of even numbers}}{Total \text{ no. of numbers}}$

$= \dfrac{4}{8} = 0.5 \Rightarrow 50\%$

percentiles:

Def: A percentile is a value below while a certain percentage of observation lie.

99 percentile means $\Rightarrow$ the pow person has got better marks than 99% of the entire students

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

*What is the percentile rank of 10?

$$\text{percentile rank of } x = \frac{\text{no. of values below } x}{n}$$

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

→ Here 16 values are present before 10 and the total values are 20.

* What is the value that exists at 25 percentile?

$$\text{value} = \frac{\text{percentile}}{100} * n+1$$

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

The value present at 5th index is 5.

* **5 number Summary:**

① Minimum

② First Quartile (25 percentile) (Q1)

③ median

④ Third Quartile (75 percentile) (Q3)

⑤ maximum

to
→ Remove the outliers
↓
and to plot
box plot

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

first to find the outliers, first we need to
find lower fence and higher fence.

$$[Lower\ Fence \longleftrightarrow Higher\ fence]$$

Lower Fence $= Q_1 - 1.5(IQR)$

Higher Fence $= Q_3 + 1.5(IQR)$

$$IQR = \overset{75}{Q_3} - \overset{25}{Q_1}$$

$$\Downarrow$$

Inter Quartile Range (IQR)

$Q_1 = \dfrac{25}{100} \times 21 = 5.25$ Index $= 3$

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

Lower Fence $= 3 - (1.5)(4.5)$

$\quad\quad = -3.65$

Higher Fence $= 7.5 + (1.5)(4.5)$

$\quad\quad = 14.25$

* In the given list, there is no values before
$-3.65$, So no need there is no outliers at

lower fence.

But in higher fence we have outlier which
is greater than 14.25.

The outlier is 27. So we have to remove
27 from the data set.

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

5 number summary used to plot the box plot here

① minimum = 1

② $Q_1 = 3$

③ median = 5

④ $Q_3 = 7.5$

⑤ maximum = 9

plot all these values

Boxplot
↓



outlier
↑

0  2  4  6  .8  10  12  14  16  18 --- 27