

[Day-2]

Page No.:

Date. / /

different types of statistics

Descriptive

- consist of organizing and summarizing

- includes

- Measure of central tendency
 - mean
 - Median
 - mode

② Measure of dispersion

- variance
- Standard deviation

college A has 1000 students.

e.g. Height = { 160, 132, 150, 140, 157 }

Descriptive - what is mean, median, mode of above data.

Inferential - taking sample predict the height of 1000 students.

Inferential

- collect sample data make conclusion of population data based on sample data.

- includes

- Hypothesis testing

① Measure of central tendency

i) Mean

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

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

e.g.

Age of students = {10, 12, 15, 17, 19}

$$\mu = \frac{10 + 12 + 15 + 17 + 19}{5} = 57.8$$

ii) Median:-

Suppose Age = {10, 12, 15, 17, 19}

then $\mu = 57.8$

If we add 87 in Age

then Age = {10, 12, 15, 17, 19, 87}

 $\mu = 87.5$

In above case 87 become outlier

to remove that we use median

∴ for even data set - arrange in asc

median = Average of two middle values

for odd data set

median = middle value.

iii) mode:

→ the data who had maximum frequency called mode

e.g. Age = { 10, 10, 12, 15, 19 }

mode = 10

② Measure of dispersion:

i) variance:

It tells about the spread of the data.

e.g. - Age = { 2, 2, 4, 4 } Age = { 1, 1, 5, 5 }

$\mu = 3$

$\mu = 3$

we can say about spread using above results.

but using variance we can tell.

Population (N)	Sample data (n)
$\sigma^2 = \sum_{i=1}^N \frac{(x_i - \mu)^2}{N}$	$s^2 = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}$
x_i = data points μ = mean N = size	x_i = data points \bar{x} = mean n = size

eg. Age = {2, 2, 4, 4}

$$\mu = \frac{2+2+4+4}{4}$$

$$\mu = 3$$

x_i	μ	$(x_i - \mu)^2$
2	3	1
2	3	1
4	3	1
4	3	1

$$N=4 \quad \sum (x_i - \mu)^2$$

$$= 4$$

$$\sigma^2 = \frac{4}{4}$$

$$\boxed{\sigma^2 = 1}$$

∴ data is less spread

Age = {1, 1, 5, 5}

$$\mu = \frac{1+1+5+5}{4}$$

$$\mu = 3$$

x_i	μ	$(x_i - \mu)^2$
1	3	4
1	3	4
5	3	4
5	3	4

$$N=4 \quad \sum (x_i - \mu)^2$$

$$= 16$$

$$\sigma^2 = \frac{16}{4}$$

$$\boxed{\sigma^2 = 4}$$

∴ data has more spread

Sample variance:

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

(n-1)

→ Bessel's correction

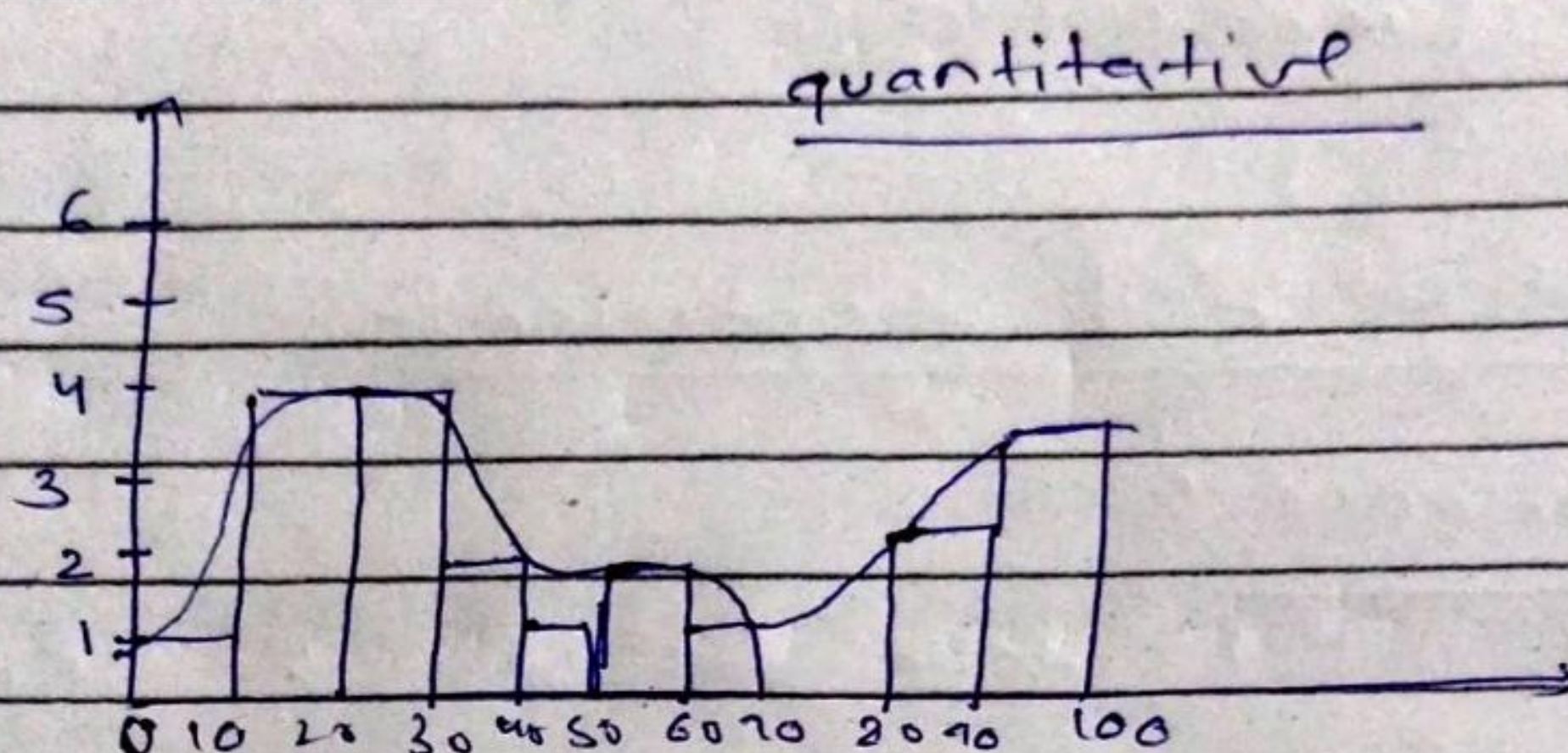
population variance:

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

Histogram: -

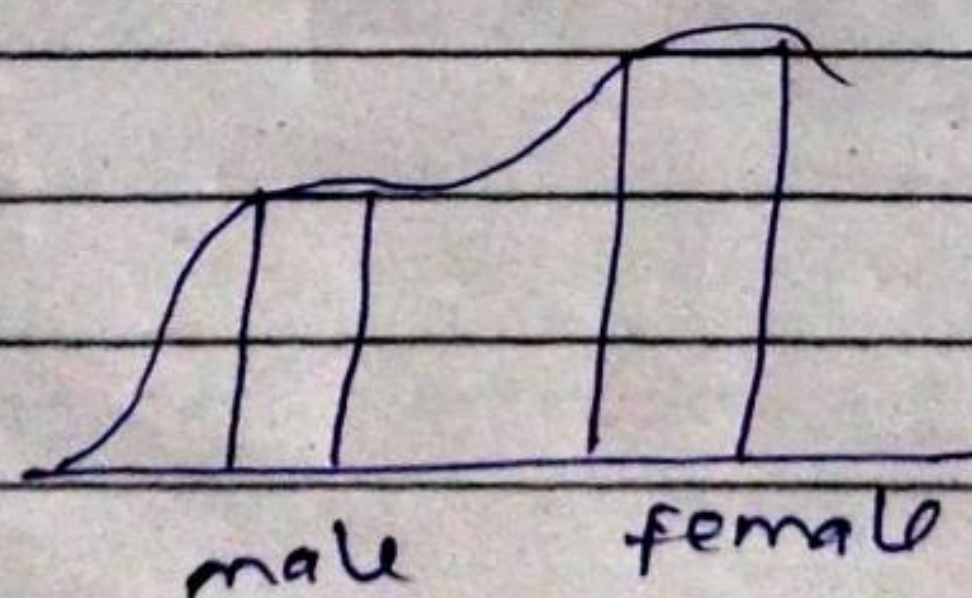
Data = { 10, 12, 13, 14, 20, 22, 24, 25, 26, 35, 38, 42, 55, 56, 68, 82, 84, 86, 92, 93, 94, 100 }

by default histogram bin = 10



pdf: - Smoother version of histogram, the technique is used kernel density function

qualitative :-



five number summary
 - to remove outlier

Here

- ① minimum
- ② first quartile
- ③ median
- ④ third quartile
- ⑤ maximum

using this we
 can create boxplot

It is used to remove outlier.

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

Here 27 is outlier.

to find this

[lower fence \leftrightarrow higher fence]

$Q_1 \rightarrow 25\%$ lower fence = $Q_1 - 1.5(IQR)$

$Q_3 \rightarrow 75\%$ higher fence = $Q_3 + 1.5(IQR)$

$$IQR = Q_3 - Q_1$$

$$Q_1 = \frac{25}{100} \times (19+1) = \frac{25 \times 20}{100} = 5^{\text{th}} \text{ index}$$

$$Q_1 = 3$$

$$Q_3 = \frac{75}{100} \times (19+1) = \frac{75 \times 20}{100} = 15^{\text{th}} \text{ index}$$

$$Q_3 = 7$$

$$\text{IQR} = 7 - 3 = 4$$

$$\begin{aligned} \text{lower fence} &= 3 - 1.5(4) \\ &= 3 - 6 \end{aligned}$$

$$\boxed{\text{lower fence} = -3}$$

$$\begin{aligned} \text{higher fence} &= 7 + 1.5(4) \\ &= 7 + 6 \\ &= 13 \end{aligned}$$

$$\text{data range } [-3, 13]$$

$\therefore 27$ is outlier in our data

$$\text{i.e. data} = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 27\}$$

$$\textcircled{1} \min = 1 \quad \textcircled{2} \text{median} = 5 \quad \textcircled{3} \max = 9$$

$$\textcircled{4} Q_1 = 3 \quad \textcircled{5} Q_3 = 7$$

box plot :

