

Session - 39

Descriptive Stats-2

• Quantiles

↳ we divide our data (eg age column) into equal size bucket.

↳ They can be used to detect outliers.

• Types of Quantiles -

i) Quartile → Divide data into 4 equal parts

↳ 25th percentile — — — 100 per-

ii) Deciles → Divide data into 10 equal parts

iii) Percentiles → Divide data into 100 "

iv) Quintiles → " " " 5 " "

Things to remember -

→ Data should be sorted from low to high

→ 25th or any percentile may not be present in data

• Percentile -

↳ 75 per- means 75% people are behind you.

$$\text{formula} \rightarrow PL = \frac{P}{100} (N+1)$$

PL → desired percentile locn

P → Percentile rank

N → total n. of observations

2) calculate 75th percentile
78, 82, 84, 88, 91, 93, 94, 96, 98, 99

Ex: Sort data = 78, 82, 84, 88, 91, 93, 94, 96, 98, 99

$$PL = \frac{75}{100} (10+1) = 8.25$$

means $\frac{1}{2}w$ 829

$$\therefore 96 + 0.25(98 - 96) = 96.5$$

- Percentile of a value -

$$= \frac{X + 0.5Y}{2}$$

$x =$ no. of values below given values equal to " "

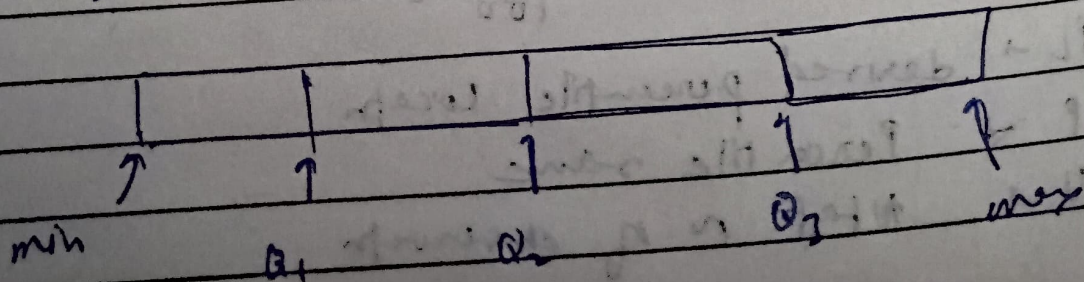
$x = \text{no. of values}$
 $y = \text{equal to " values}$

$N =$ total no. of values

4. $\Rightarrow \frac{3 + 0.5 \times 1}{10} = 0.85$
= 85 percentile

5 number Summary -

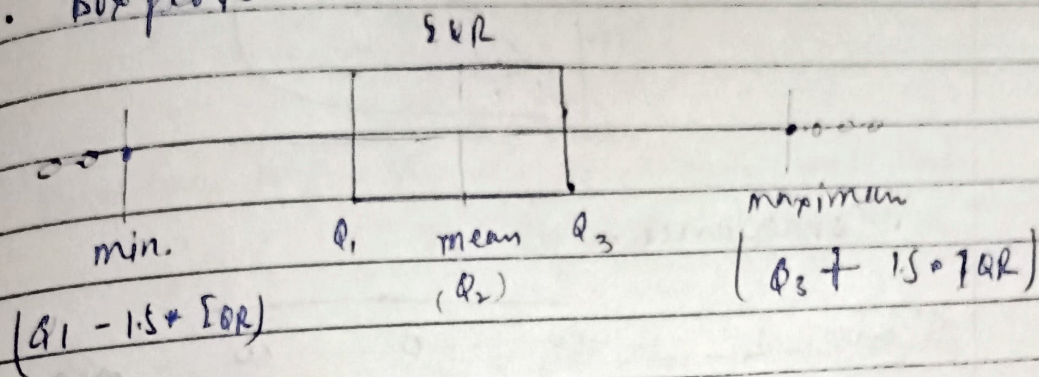
- 1) min. value
- 2) first quartile (Q_1)
- 3) median (Q_2)
- 4) Third quartile (Q_3)
- 5) maximum value



• Interquartile Range (IQR)

$$\rightarrow Q_3 - Q_1 \quad (\text{middle } 50\%)$$

• Box plots.

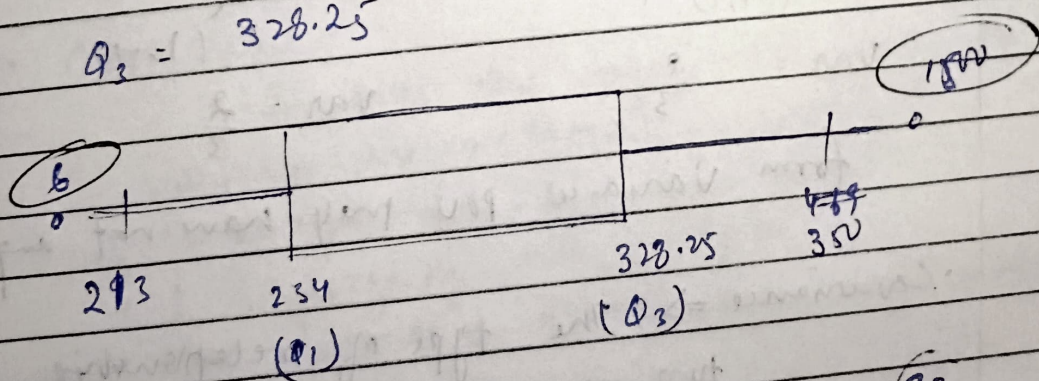


6	213	241	260	261	290	314	321	350	1500
1	2	3	4	5	6	7	8	9	10

$$Q_2 = \frac{50}{100} \times 11 = 5.5 = 285.5$$

$$Q_1 = \frac{25}{100} \times 11 = 2.75 = 213 + 0.75(241 - 213) = 234$$

$$Q_3 = 328.25$$

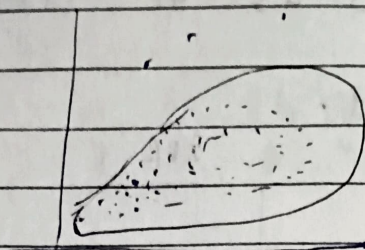


$$\text{min} = Q_1 - 1.5 \times \text{IQR} = 93$$

$$\text{max} = Q_3 + 1.5 \times \text{IQR} = 469$$

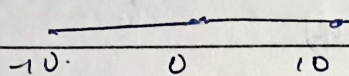
Since 93 is not present so we will stop at 213 (acc. to data)

• Scatter plot

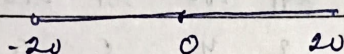


• Covariance -

Scenario-1 -



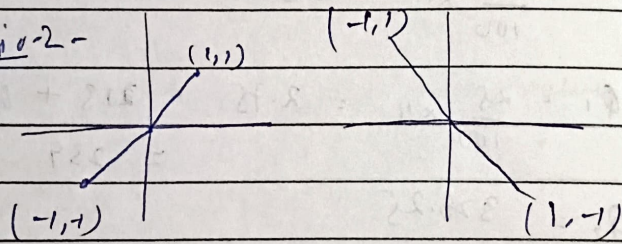
mean = 0



mean = 0

from mean POV they have not any diff.

Scenario-2 -

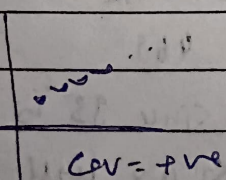


$$\text{var} = \frac{8}{3}$$

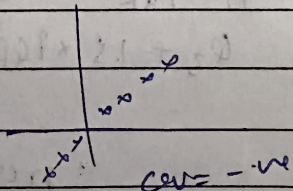
$$\text{var} = \frac{2}{3}$$

from Variance POV they have not any diff.

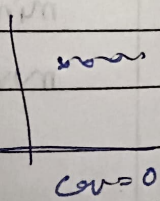
• Covariance \Rightarrow The type of relationship b/w two numerical column



one col \uparrow then
another col \uparrow



one col \uparrow then
another col \downarrow



No relation
ship

formula →

Population :-
$$\frac{\sum (X - \mu_x) (Y - \mu_y)}{N}$$

Sample :-
$$\frac{\sum (X - \bar{x}) (Y - \bar{y})}{n-1}$$

Disadvantage cov → It doesn't tell about the strength of relation

⇒ Covariance of a variable with itself is variance

$$\text{cov} = \frac{\sum (x_i - \bar{x}) (y_i - \bar{y})}{n-1} = \frac{\sum (x_i - \bar{x}) (x_i - \bar{x})}{n-1}$$

$$\Rightarrow \frac{\sum (x_i - \bar{x})^2}{n-1}$$

Correlation -

It tells the degree to which two variables are related.

- ranges $[-1, +1]$
- 1 ⇒ perfectly +ve relation
- more towards 1 means more +vely related

$$\text{Corr} = \frac{\text{cov}(x, y)}{\sigma_x * \sigma_y}$$

→ Correlation doesn't imply causation.