① percentiles and Quantiles [

percentage :

EX: 1,2,3,4,5,6

% of numbers that are odd $= \dfrac{3}{6} = \dfrac{No.of \; odd \; numbers}{Total \; No.of \; numbers}$

$= \dfrac{1}{2} = 50\%$

percentiles :- A percentile is a value below which a certain percentage of data points lie.

$$x = \{2.3, 3, 4, 6, 6, 6, 7, 889, 9, 9, 10, 11, 12\}$$

① percentile Rank of $10 = \dfrac{\# \; of \; values \; below \; 10}{} \times 100$

$= \dfrac{\cancel{12}^{4}}{\cancel{15}_{8}} \times \cancel{100}^{20} = 80 \; percentile.$

80 percentile = 80% of the distribution fall below the value of 10.

② what value exists at 25th percentile ?

Value $= \dfrac{percentile}{100} * (n+1)$

$\dfrac{\cancel{25}^{1}}{\cancel{100}_{20 \; 4}} \times \cancel{16}^{4} = 4^{th} \; element$

Suppose if value is in decimal form in that case we actually needs to take average of the last No. present value at the value and next number in the sequence

$$X = \{2, 3, 3, \underset{\downarrow}{4}, \underset{\downarrow}{6}, 6, 6, 7, 8, 8, 9, 9, 10, 11, 12\}$$

when value = $4 \cdot 5 \Rightarrow \frac{4+6}{2} = 5$

## Quartiles :-

$Q_1 \longrightarrow 25$ percentile

$Q_2 \longrightarrow$ Median $\longrightarrow 50$ percentile.

$Q_3 \longrightarrow 75$ percentile.

## Five Number Summary :-

1) Minimum

2) First Quartile ($25$ percentile) $\rightarrow Q_1$

3) Median ($Q_2$)

4) Third Quartile ($75$ percentile) ($Q_3$)

5) Maximum.

## Removing the Outliers :-

$$X = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 29\}$$

[Lower Fence $\longleftrightarrow$ Higher Fence]

Lower fence $= Q_1 - 1 \cdot 5 \text{ (IQR)}$    Inter Quartile Range $= Q_3 - Q_1$

Higher Fence $= Q_3 + 1 \cdot 5 \text{ (IQR)}$

$$X = \{1, 2, 2, 2, 3, 3, 4, 5, 5, 5, 6, 6, 6, 6, 7, 8, 8, 9, 29\}$$

$Q_1 = 25^{st}$ percentile $= \dfrac{25}{100} \times (19+1)$

$= \dfrac{25}{100} \times 20 = 5^{th}$ element $= 3$

$Q_3 = 75$ percentile $= \dfrac{75}{100} \times (20)$

$= 15^{th}$ element $= 7$

$IQR = 7 - 3 = 4$

Lower Fence $= Q_1 - 1.5(IQR)$

$= 3 - 1.5(4)$

$= 3 - 6$

$= -3$

Higher Fence $= Q_3 + 1.5(IQR)$

$= 7 + 1.5(4)$

$= 13$

The values which are not in this range of values $[-3, 13]$

are considered as outliers

Box plot : [To visulalize outliers]

① min value = 1
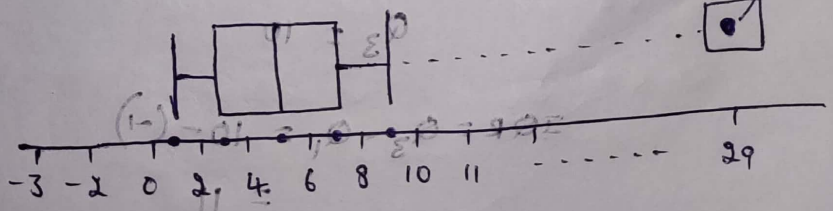
② $Q_1 = 3$

③ median $Q_2 = 5$

④ $Q_3 = 7$

⑤ Maximum = 9

outliers

-3  -2  0  2  4  6  8  10  11          29

## Right skewed



$$Q_3 - Q_2 > Q_2 - Q_1$$

mean > median > mode

## Left skewed



$$Q_2 - Q_1 < Q_3 - Q_2$$

mode > median > mean

---

Find the outliers and draw box plot ?

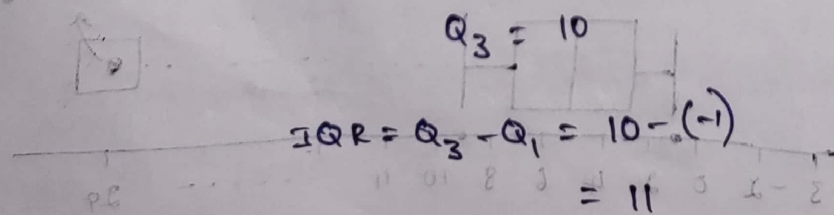$$y = \{-13, -12, -6, -5, 3, 4, 5, 6, 7, 7, 8, 10, 10, 11, 24, 55\}$$

$Q_1 = 25^{th}$ percentile $= \dfrac{25}{100} \times (16+1)$

$$= \dfrac{25}{100} \times 17 = 4.25$$

$$Q_1 = -\dfrac{2}{2}$$

$Q_3 = 75^{th}$ percentile $= \dfrac{75}{100} \times (16+1)$

$$= 12.75 \Rightarrow \left( \dfrac{10+10}{2} = 10 \right)$$

$$Q_3 = 10$$
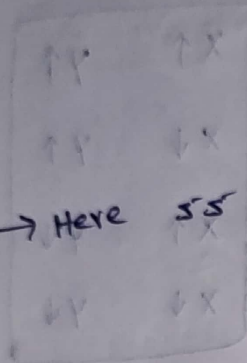
$$IQR = Q_3 - Q_1 = 10 - (-1)$$

$$= 11$$

Lower Fense $= -1 - 1.5(10+1)$

$= -1 - 1.5(11)$

$= -17.5$

Higher Fense $= 10 + 1.5(10+1)$

$= 10 + 16.5$

$= 26.5$

→ Here $55$ is an outlier

$z = \{1, 2, 4, 6, 7, 12, 18, 34, 77, 66, 108, 99, 14\}$

step : sort

$z = \{1, 2, 4, 6, 7, 12, 14, 18, 34, 66, 77, 99, 108\}$

$Q_1 = \frac{25}{100} \times 14 = 3.05$   $\left(\frac{4+6}{2} = 5\right) = 5$

$Q_3 = \frac{75}{100} \times 14 = 10.5$   $\left(\frac{66+77}{2}\right) = 71.5$

$IQR = 71.5 - 5 = 66.5$

$LF = 5 - 1.5(66.5)$   $HF = 71.5 + 1.5(66.5)$

$= -94.75$   $= 171.25$

→ There are no outliers in this dataset.

Covariance and correlation :-
→ one of the very important topic for data preprocessing, data analysis and feature selection

Let us consider two random variables :-

| size(x) | price(y) |
|---------|----------|
| 1200 sqm | 100 K $ |
| 1500 sqm | 200 K $ |
| 1800 sqm | 300 K $ |
| 2000 sqm | 400 K $ |

Relationship between size and price
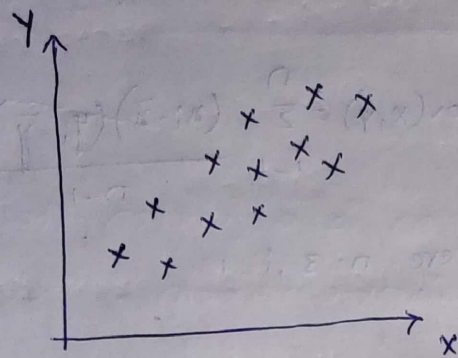
| | |
|---|---|
| X ↑ | Y ↑ |
| X ↓ | Y ↑ |
| X ↑ | Y ↓ |
| X ↓ | Y ↓ |

→ Covariance and correlation are two statistical concepts that are used to measure the relationship between two variables. Although they are different concepts and have different interpretations.

→ covariance :- measures how two variables vary together Specifically covariance measures how much two variables vary from their respective means at the same time

→ A positive covariance means that the two variables tend to increase or decrease together.

→ while a negative covariance means that one variable tends to increase while the other decreases.

covariance tells us about direction of relationship.

if x↑, y↑ = +ve direction

if x↑, y↓ = -ve direction
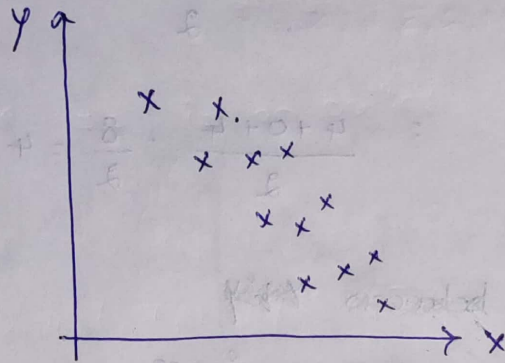
X↑   Y↑

X↓   Y↓

X↑ Y↓

X↓ Y↑

① Tells us about relationship between X & Y.

## Covariance :

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

→ covariance does not limit a specific value
(Co $\infty$ to $-\infty$)

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

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

$$\boxed{\text{Var}(x) \Leftarrow \text{Cov}(x,x)}$$

Cov(x,y)

| X↑ Y↑ |
| X↓ Y↓ |  ⇒ +ve Covariance

| X↑ Y↓ |
| X↓ Y↑ |  ⇒ −ve Covariance

Ex:-

| X | Y |
|---|---|
| 2 | 3 |
| 4 | 5 |
| 6 | 7 |

$\bar{x} = 4$   $\bar{y} = 5$

$$Cov(x,y) = \sum_{i=1}^{n} \frac{(x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

Here   $n = 3, i = 1$

$$= \frac{[(2-4)(3-5) + (4-4)(5-5) + (6-4)*(7-5)]}{2}$$

$$= \frac{4 + 0 + 4}{2} = \frac{8}{2} = 4 \quad \text{+ve covariance}$$

conclusion:

→ x and y are having a positive covariance.

Advantage

① Tells us about relationship between x & y.

disadvantage

① covariance does not have a specific limit value. $(-\infty$ to $\infty)$

## pearson correlation coefficient:-

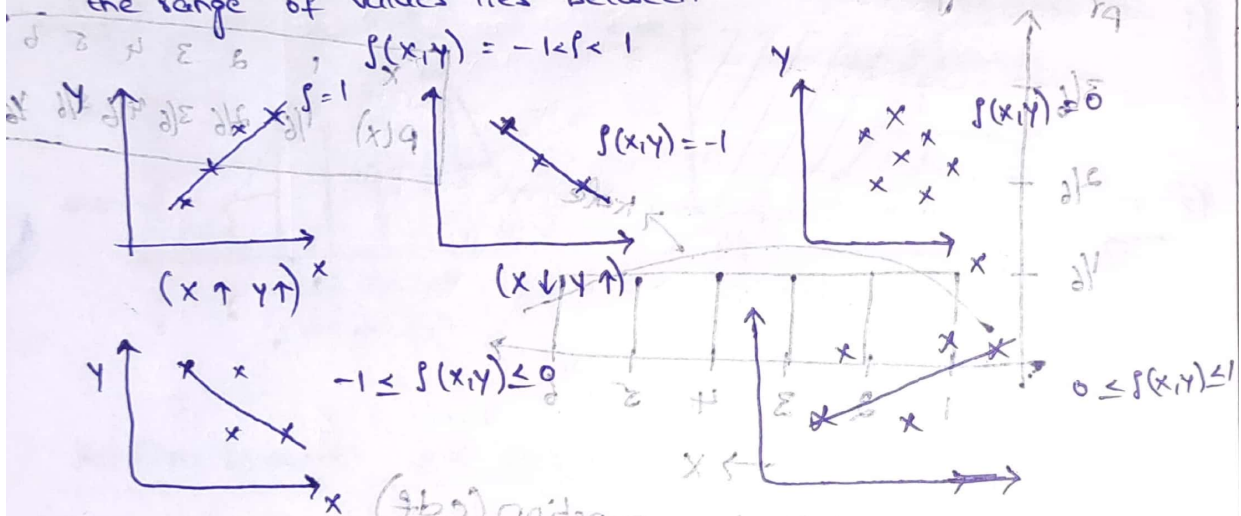$$f(x,y) = \frac{Cov(x,y)}{\sigma_x \cdot \sigma_y}$$

→ correlation measures the strength and direction of the linear relationship between two variables, without being affected by units of measurements.

→ The more the value towards +1 the more +ve correlated it is.

→ The more the value towards −1 the more +ve correlated it is.

→ when we try to find the pearson correlation coefficient the range of values lies between

$$f(x,y) = -1 < f < 1$$



$f = 1$
$(x \uparrow y \uparrow)$

$f(x,y) = -1$
$(x \downarrow y \uparrow)$

$f(x,y) = 0$

$-1 \leq f(x,y) \leq 0$

$0 \leq f(x,y) \leq 1$

## Spearman's correlation coefficient:

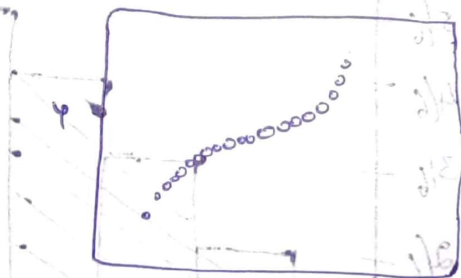→ spearman's rank correlation measures the strength and direction of association between two ranked variables.

$$r_s = \frac{Cov(R(x), R(y))}{\sigma_{R(x)} \cdot \sigma_{R(y)}}$$

→ spearman's correlation coefficient ranges between −1 and 1,

→ Here correlation 1 represents a perfect monotonic relationship between the two variables being correlated.

→ which means that as one variable increases the other variable also increases in a perfectly predictable manner.

→ This is observed when the two variables are perfectly ranked.

spearman correlation − 1
pearson correlation − 0.88



| x | y | R(x) | R(y) |
|---|---|---|---|
| 5 | 6 | 3 | 1 |
| 7 | 4 | 2 | 2 |
| 8 | 3 | 1 | 3 |
| 1 | 6 | 5 | 5 |
| 2 | 2 | 4 | 4 |