# Covariance And Correlation

**Covariance and correlation are two statistical measures used to determine the relationship between two variables. Both are used to understand how changes in one variable are associated with changes in another variable.**

## Covariance

**Definition**: Covariance is a measure of how much two random variables change together. If the variables tend to increase and decrease together, the covariance is positive. If one tends to increase when the other decreases, the covariance is negative.

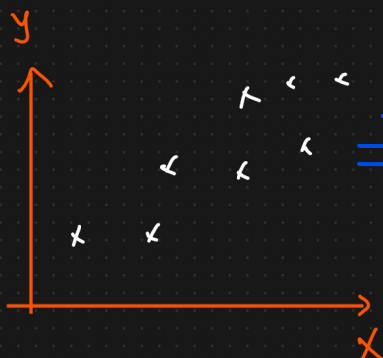to : [ Quantify the Relationship between X and Y ]

weather it is:

| X | Y |
|---|---|
| → 2 | 3 |
| → 4 | 5 |
| → 6 | 7 |
| → 8 | 9 |

X↑ Y↑
X↓ Y↑
X↑ Y↓
X↓ Y↓

Dataset
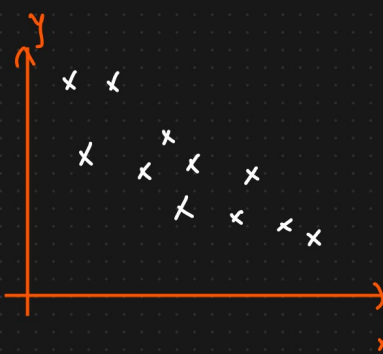
| ↓ ↑ Size of House | Price ↑↓ |
|---|---|
| 1200 | 45 lakhs |
| 1300 | 50 lakh |
| 1500 | 75 lakh |

it can be seen :

X↑ Y↑
X↓ Y↓

⟹ =) +ve Covariance =) +ve value

it can be seen that:

X↓ Y↑
X↑ Y↓

| X | Y |
|---|---|
| 7 | 10 |
| 6 | 12 |
| 5 | 14 |
| 4 | 16 |

=) −ve Covariance =) −ve value

**Covariance** of x to y :

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

thus , cov of x to x will be :

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

$$\boxed{Cov(x,x) = Var(x)}$$

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

- $x_i \rightarrow$ Datapoint of random variable $x$
- $\bar{x} \rightarrow$ Sample mean of $n$
- $y_i \rightarrow$ Datapoints of random variable $y$
- $\bar{y} \rightarrow$ Sample mean of $y$

thus the covariance of x to x is variance of x

e.g :

**Students**

| Hour Studied (X) | Exam Score (Y) |
|---|---|
| 2 | 50 |
| 3 | 60 |
| 4 | 70 |
| 5 | 80 |
| 6 | 90 |

as , we can see that :

$x\uparrow \quad y\uparrow$

$x\downarrow \quad y\downarrow$

$\Rightarrow$ +ve covariance

it is : $\downarrow$

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

proof :

① $\bar{x} = \dfrac{2+3+4+5+6}{5} = 4$ //

② $\bar{y} = \dfrac{50+60+70+80+90}{5} = 70$ //

$$Cov(x,y) = \frac{(2-4)(50-70) + (3-4)(60-70) + (4-4)(70-70) + (5-4)(80-70) + (6-4)(90-70)}{4}$$

$Cov(x,y) = 20$

$\Rightarrow$ The positive covariance indicates the no. of hours studied increased the exam score also.

$\begin{cases} x & y \\ 7 & 10 \\ 6 & 12 \\ 5 & 14 \end{cases} \Rightarrow$ cov is : $-ve$.

$x\uparrow \ y\downarrow$
$x\downarrow \ y\uparrow$

$0.96$
$Cov(A,B)$

$0.8$
$Cov(B,C)$

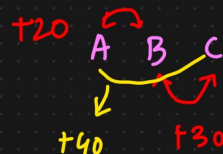$\begin{array}{c} -200 \\ +100 \\ 20 \\ Cov(A,B) \end{array}$

$\begin{array}{c} -300 \\ +300 \\ 30 \\ Cov(B,C) \end{array}$

Advantages:

$\boxed{-1 \text{ to } 1}$

Disadvantage

to -
① Quantify the Relationship between X and Y

① Covariance does not have a specific limit value.
$$Cov(X,Y) \Rightarrow -\infty \text{ to } \infty$$

+20  A  B  C
↓        ↺
+40    +30

② Correlation
- i. → Pearson Correlation Coefficient
- ii. → Spearman Rank Correlation

① Pearson Correlation Coefficient in: $\Rightarrow \boxed{-1 \text{ to } 1}$ always.

$$\int_{x,y} = \frac{Cov(x,y)}{\sigma_x \cdot \sigma_y} = \frac{20}{\sigma_x \cdot \sigma_y} \Rightarrow 0 \text{ to } 1$$
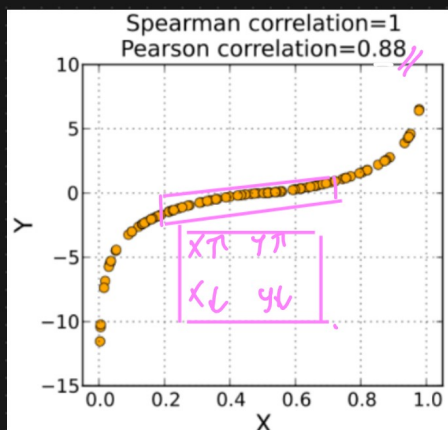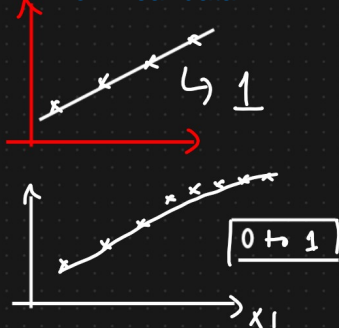
✓ The more the value towards +1 the more +ve correlated X & Y is.

✗ The more the value towards -1 the more -ve correlated it is (X,Y)

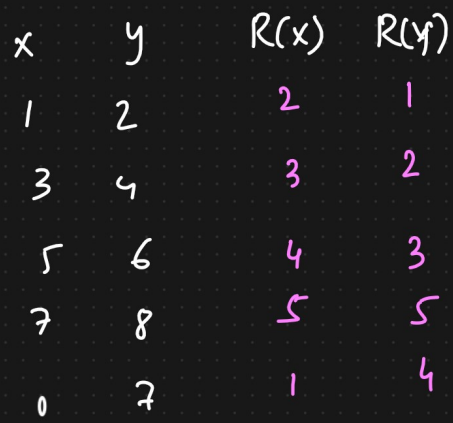① Spearman Rank Correlation

for non linear data:

for linear data:
↳ $\boxed{1}$

Pearson correlation not able to capture the correlation for non linear data

$\boxed{0 \text{ to } 1}$

Spearman correlation=1
Pearson correlation=0.88

X↑ Y↑
X↓ Y↓

A Spearman correlation of 1 results when the two variables being compared are monotonically related, even if their relationship is not linear. This means that all data points with greater x values than that of a given data point will have greater y values as well. In contrast, this does not give a perfect Pearson

relationship:
↓
⇒ X↑ Y↑
⇒ X↓ Y↓
↓
Pearson Correlation
= 0.88

$y$ (axis graph with curve)

| x | y | R(x) | R(y) |
|---|---|------|------|
| 1 | 2 | 2 | 1 |
| 3 | 4 | 3 | 2 |
| 5 | 6 | 4 | 3 |
| 7 | 8 | 5 | 5 |
| 0 | 7 | 1 | 4 |

$$Y_S = \frac{Cov\left(R_{(x)}, R_{(y)}\right)}{\sigma\left(R_{(x)}\right) * \sigma\left(R_{(y)}\right)}$$

⇐ R(x) means the rank of x , according to the values . eg. R(1) is 2 for x values.

these concepts are used for :

## Feature Selection

Size of ↑
House
=

+ve correlated

No. of Room ↑
=

Location ↑
=

No. of
People stays
in the house

~~No. of People stays in the house~~

No

Haunted

$\boxed{2=0}$ ↗

Price ↑↑ ↓↓
= //.

O/p
↑↓

-ve Correlation ↓
=

as there is no relation between this and the price is not possible

note :
t works well for non-linear but monotonic relationships.
If the relationship is non-monotonic, Spearman's rank correlation might not capture the relationship accurately.