

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.

[Quantify the Relationship between X and Y]

X	Y
→ 2	3
→ 4	5
→ 6	7
→ 8	9

X↑	Y↑
X↓	Y↑
X↑	Y↓
X↓	Y↓

Dataset

↓ ↑ Size of house

Price ↑ ↓

1200

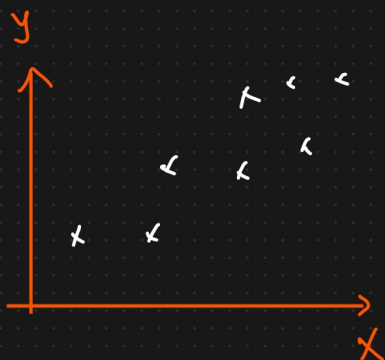
45 lakhs

1300

50 lakh

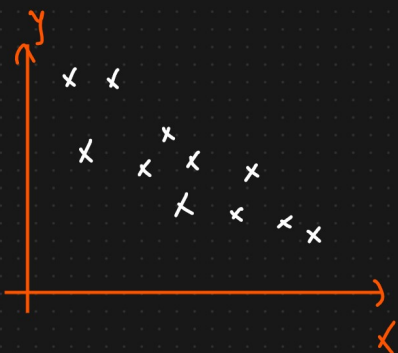
1500

75 lakh



X↑	Y↑
X↓	Y↓

⇒ +ve Covariance ⇒ +ve value



X↓	Y↑
X↑	Y↓

X	Y
7	10
6	12
5	14
4	16

⇒ -ve Covariance ⇒ -ve value

Covariance

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

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

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

$x_i \rightarrow$ Datapoint of random variable x

$\bar{x} \rightarrow$ Sample mean of x

$y_i \rightarrow$ Datapoints of random variable y

$\bar{y} \rightarrow$ Sample mean of y

Students

Hour studied (x)

Exam Score (y)

2

50

3

60

4

70

5

80

6

90

$x \uparrow \quad y \uparrow \Rightarrow +ve$
 $x \downarrow \quad y \downarrow$ (covariance)

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

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

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

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

$$Cov(x, y) = \underline{\underline{20}}$$

\Rightarrow The positive covariance indicates the no. of hours studied increased the Exam Score also.

$$\left\{ \begin{array}{c} x \\ 7 \\ 6 \\ 5 \end{array} \quad \begin{array}{c} y \\ 10 \\ 12 \\ 14 \end{array} \right\} \Rightarrow \underline{\underline{-ve}}$$

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

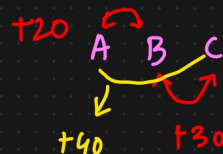
$$\begin{array}{cc} 0.96 & 0.98 \\ Cov(A, B) & Cov(B, C) \end{array}$$

$$\begin{array}{cc} -200 & -300 \\ +100 & +300 \\ 20 & 30 \\ Cov(A, B) & Cov(B, C) \end{array}$$

Advantages

$[-1 \text{ to } 1]$

Disadvantage



- ① Quantify the Relationship between X and Y

- ① Covariance does not have a specific limit value.

$$\text{Cov}(X, Y) \Rightarrow -\infty \text{ to } \infty$$

- ② Correlation
- Pearson Correlation Coefficient
 - Spearman Rank Correlation

- ① Pearson Correlation Coefficient $\Rightarrow [-1 \text{ to } 1]$

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

$$= \frac{20}{\sigma_x \cdot \sigma_y} \Rightarrow 0 \text{ to } 1$$

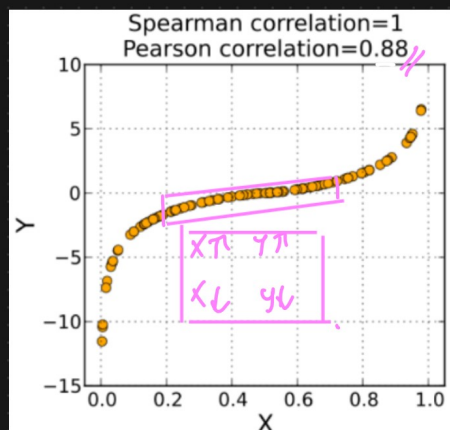
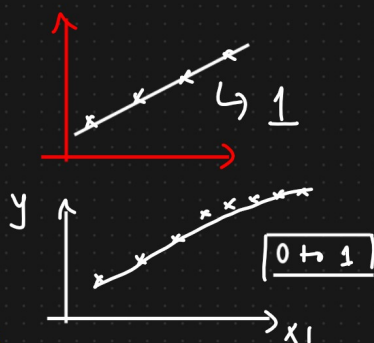
- ① The more the value towards $+1$ the more +ve correlated x & y is.

- ② The more the value towards -1 the more -ve correlated it is (x, y)

- ② Spearman Rank Correlation

Pearson Correlation

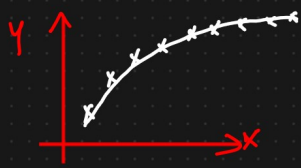
Correlation for
non linear data



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

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

Pearson Correlation
 $= 0.88$



x	y	R(x)	R(y)
1	2	2	1
3	4	3	2
5	6	4	3
7	8	5	5
9	7	6	4

$$r_s = \frac{\text{cov}(R(x), R(y))}{\sigma(R(x)) * \sigma(R(y))} \Leftarrow$$

Feature Selection

