

Covariance - Quantifying a relationship between variables. Direction of a relationship.

Size	Price
1200 sqm	₹ 1L
1800 sqm	₹ 2L
2500 sqm	₹ 3L

Quantify a relationship between Size and Price.

Size \uparrow Price \uparrow
Size \downarrow Price \downarrow

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

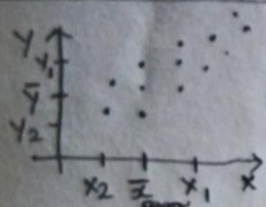
$$\text{Variance}(x), \text{var}(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 = \frac{1}{n} \sum_{i=1}^n (x - \bar{x}) * (x - \bar{x})$$

$$\text{Cov}(x, x) = \text{var}(x)$$

So, through covariance we will get a value.

Suppose $x \uparrow, y \uparrow$.

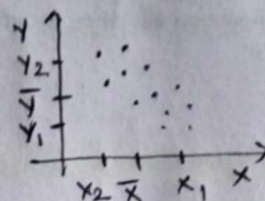
$$\begin{aligned} &= (x_1 - \bar{x}) * (y_1 - \bar{y}) \\ &= (+ve) * (+ve) \\ &= +ve \\ &= (y_2 - \bar{y}) * (x_2 - \bar{x}) \\ &= (-ve) * (-ve) \\ &= +ve \end{aligned}$$



If $x \uparrow, y \uparrow = \square$ +ve Covariance

$x \uparrow, y \downarrow = \square$ -ve Covariance

Suppose $x \uparrow, y \downarrow$



$$\begin{aligned} &= (x_1 - \bar{x}) * (y_1 - \bar{y}) \\ &= (+ve) * (-ve) \\ &= -ve \\ &= (x_2 - \bar{x}) * (y_2 - \bar{y}) \\ &= (-ve) * (+ve) \\ &= -ve \end{aligned}$$

So, always covariance is +ve when $x \uparrow, y \uparrow$.

So, always covariance is -ve when $x \uparrow, y \downarrow$.

Covariance find the direction of relationship. But we don't know exact value of strength.

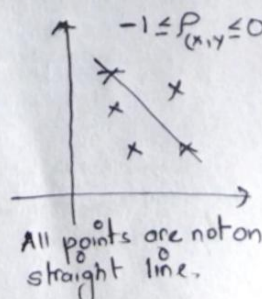
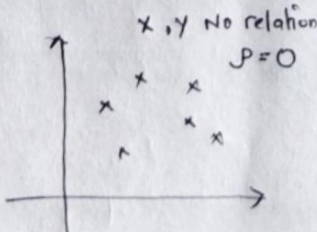
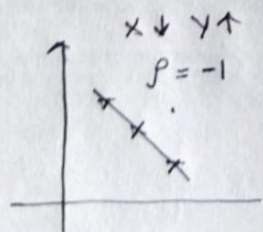
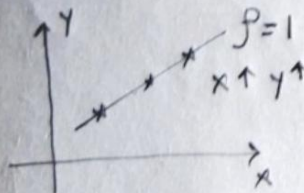
Pearson Correlation Coefficient

→ Strength of the relationship between variables and direction also of the relationship.

Range is $\rightarrow -1 \leq P(x, y) \leq 1$.

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

$\text{Cov}(x, y) \rightarrow$ covariance (x, y)
 $\sigma_x \rightarrow$ Standard deviation of x .
 $\sigma_y \rightarrow$ SD of y .

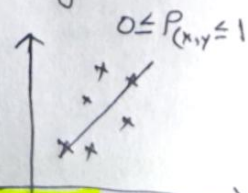


Suppose we have 3 variables a, b, c .

a and b have Pearson correlation coefficient of 1. In short they are same, then we will remove one of the variable/feature.

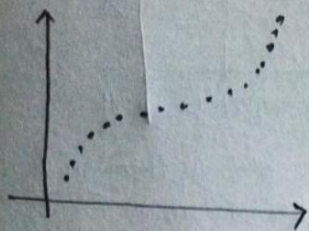
We also use **Spearman Rank Correlation** (Refer to Numerical Measure-Mean)

Pearson work only good for linear relationship. But for others we use Spearman Rank Correlation (Refer to page 4) coefficient.



Spearman's Rank Correlation coefficient

If the plot is non-linear relationship.



Spearman correlation = 1.

Pearson correlation = 0.88

$$\text{Pearson} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

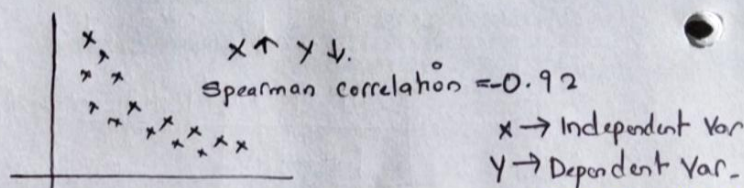
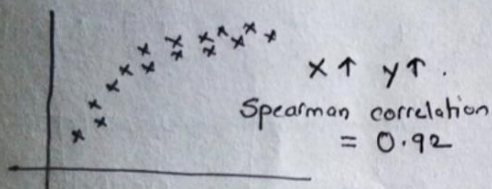
$$\text{Spearman} = \frac{\text{Cov}(\text{rank}_X, \text{rank}_Y)}{\sigma_{\text{rank}_X} \sigma_{\text{rank}_Y}}$$

- Pearson correlation assesses only linear relationship. whereas Spearman correlation assesses monotonic relationship (whether linear or not).

- Spearman correlation is from range -1 to 1.

- If the two variables are very similar, then Spearman correlation will be 1.

- If the two variables are dissimilar / fully opposed, then Spearman correlation will be -1.



Sign of Spearman correlation indicates the direction of association between X and Y

⊛ If Spearman correlation = 0, then there is no tendency of Y increase or decrease when X is increased.

Steps in the methods are - i) Sort the data and assigned the rank.

ii) Sort the second column and assign the rank.

iii) Create a difference column based on rank.

iv) Create final column based on difference by squaring.

IQ	Hours of TV	Sort →	IQ (x _i)	Hours of TV (y _i)	rank (x _i)	rank (y _i)	d _i	d _i ²
106	7		86	0	1	1	0	0
100	27		97	20	2	6	-4	16
86	2		99	28	3	8	-5	25
101	50		100	27	4	7	-3	9
99	28		101	50	5	10	-5	25
103	29		103	29	6	9	-3	9
97	20		106	7	7	3	4	16
113	12		110	17	8	5	3	9
112	6		112	6	9	2	7	49
110	17		113	12	10	4	6	36

$$\sum d_i^2 = 194.$$

$$n = 10.$$

$$P = 1 - \frac{6 \sum d_i^2}{n(n^2 - 1)}$$

Shortcut formula.

$$= 1 - \frac{6(194)}{10(10^2 - 1)}$$

$$= -0.175.$$

So, $P = -0.175$ which is very close to 0. So there is no correlation between IQ and watching TV for hours.