

# Covariance And Correlation

13 October 2023 21:13

Covariance and correlation are essential tools in various fields, such as statistics, data science, machine learning, and data analysis. They serve as useful measures for determining the relationship between two variables.

These concepts are particularly significant in artificial intelligence and machine learning, as they are frequently employed in linear regression and neural networks to model and predict the relationship between variables.

However, they have different properties and may be used in different contexts depending on the research question and the data being analysed.

Understanding the relationship between two variables is a critical aspect of data analysis. Two commonly used measures of relationships are correlation and covariance.

While both provide useful insights into the relationship between two variables, they have distinct properties and may be used in different contexts.

Correlation and covariance are sensitive to outliers, so checking for outliers is important before calculating these measures.

While correlation measures the linear relationship between two variables, it may not capture the full extent of the relationship if it is not linear. In cases where the relationship is not linear, other measures, such as nonparametric correlation coefficients or nonlinear regression, may be more appropriate.

Sometimes, a high correlation coefficient may not necessarily imply causality between the two variables. The correlation only measures the association between two variables, and other factors may affect the relationship between the two variables.

Covariance and correlation can be calculated using different methods, such as raw data, deviations from the mean, or data ranks. The choice of method can affect the resulting correlation or covariance coefficient.

## Example

Suppose we have two variables, X and Y, and we want to measure their relationship. We calculate the covariance and correlation coefficients and obtain the following results:

Covariance: 500  
Correlation: 0.8

At first glance, X and Y appear to have a strong, positive relationship. However, upon further inspection, we find one outlier in the data driving the results. After removing the outlier, we recalculate the covariance and correlation coefficients and obtain the following results:

Covariance: 200  
Correlation: 0.6

We can see that the correlation coefficient decreased, indicating a weaker relationship between X and Y, and the covariance decreased even more. This illustrates the importance of checking for outliers and the potential for spurious correlations when working with real-world data.

It's crucial to carefully examine the data and ensure that no outliers or other elements that might

skew the results are responsible for the findings.

## Correlation vs. Covariance

Measures of the relationship between two variables

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

correlation between  $X$  and  $Y$       standard deviation of  $X$       standard deviation of  $Y$       covariance normalized by standard deviation

@Danny Butvinik  
The AI Vanguard  
newsletter

### Correlation

- measures the **strength** and **direction** of the linear relationship between two variables.
- standardized values
- not affected by scale
- value range  $[-1 \dots 1]$

$$\text{corr}_{XY} = \rho_{XY} = \frac{E[(X - \mu_X)(Y - \mu_Y)]}{\sigma_X \sigma_Y}$$

Expected value operator  
Random variables  
means

### Covariance

- measures the **extent** to which two variables are linearly related or dependent,
- not standardized values
- affected by scale of the variable
- value range  $[-\infty \dots \infty]$

$$\text{cov}_{XY} = \sigma_{XY} = E[(X - \mu_X)(Y - \mu_Y)]$$

It doesn't have relationship

Covariance: One of the many topic Data preprocessing, quantity.  
 ID var → Dependent variable  
 Size → Price variable  
 1200 sqm → 1000 \$  
 1800 sqm → 2000 \$  
 1800 sqm → 3000 \$  
 Size ↔ Price or Relate  
 Quantifying my  
 ST P ↑ +ve  
 ST P ↓ -ve

Covariance (X, Y) =  $\frac{1}{n} \sum_{i=1}^n (x_i - \mu_x) \cdot (y_i - \mu_y)$   
 Size ↓ Y

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

⇒ Cov(x, x) = Var(x)  
 Covariance is +ve

Cov(x, y) =  $\frac{1}{n} \sum_{i=1}^n (x_i - \mu_x) \cdot (y_i - \mu_y)$   
 Random variables → ID & D  
 mean →

x ↑ y ↑ =  $\square$  +ve

x ↑ y ↓ = -ve

In covariance we find direction of relationship but don't know how much the +ve (-ve) value we use Pearson cc

08

282-084 • WK 41

2020

Thursday  
October

Covariance

Covariance

10	October 2020						
wk	M	T	W	T	F	S	S
40				1	2	3	4
41	5	6	7	8	9	10	11
42	12	13	14	15	16	17	18
43	19	20	21	22	23	24	25
44	26	27	28	29	30	31	

Pearson Correlation Coefficient

direction of relationship  
b/w x & y

+ve  
corr  
1 2  
 29 yr height 110 lbs  
 29 yr 110 lbs  
 -ve

$$\text{Pearson cc} = r(x, y) = \frac{\text{corr}(x, y)}{\sigma_x \sigma_y}$$

In covariance we don't know strength and direction of relation

Range b/w  $r(x, y) = [-1, 1]$

we divide by variance so it ranges

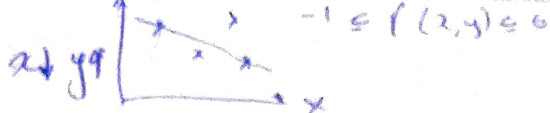
why?   
 the -1 to 1 Because it measures the linear relationship b/w two variables

-1 = Perfect negative linear relationship b/w 2 variables  
 It means 1 variable increases the other decreases in perfectly constant manner

NOTES

1 = Perfect +ve linear relationship between the two variables. This means 1 variable increases other also increases





$$-1 \leq r_{(x,y)} \leq 0$$

2020  
Friday  
October

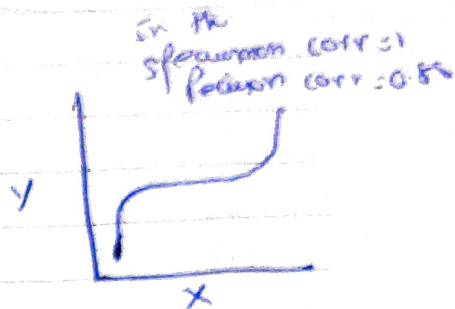
09

283-083 • WK 41

0 → no linear relationship

## SPEARMAN'S RANK CORRELATION

$$R_s = r_{r_{gx}, r_{gy}} = \frac{\text{cor}(r_{gx}, r_{gy})}{r_{gx} \cdot r_{gy}}$$



$$R = r_{\text{rank}}$$

First find rank

Steps

① Sort the data by the first column (x) create a new column and assign it the ranked values 1, 2, 3, ...

② Next Sort the data by the second column (y) and simultaneously assign it ranked values 1, 2, 3, ...

③ Create 5th column  $d_i$  to hold difference b/w two rank column (x and y)

④ Create one final column  $d_i^2$  to hold the values of column  $d_i$  squared.

# 10

284-082 • WK 41

2020

Saturday

October

Example

10	October 2020						
wk	M	T	W	T	F	S	S
40				1	2	3	4
41	5	6	7	8	9	10	11
42	12	13	14	15	16	17	18
43	19	20	21	22	23	24	25
44	26	27	28	29	30	31	

	$T_p(x)$	Hours of TV per week	score (x)	rank (y)	$d_i$	$d_i^2$
9.00						
10.00	86	0	1	1	0	0
	97	26	2	6	-4	16
11.00	99	28	3	8	-5	25
	100	27	4	7	-3	9
12.00	101	56	5	10	-5	25
	103	29	6	9	-3	9
1.00	106	7	7	3	4	16
	110	17	8	5	3	9
2.00	112	6	9	2	7	49
	113	12	10	4	9	36
3.00						

Spearman  
corr

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

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

11 Sunday

$$= r_s = \frac{-29}{15} = -0.193 \quad X \uparrow Y \downarrow$$

weak  
neg low

NOTES



November 2020

M	T	W	T	F	S	S		
1	2	3	4	5	6	7	8	
9	10	11	12	13	14	15		
16	17	18	19	20	21	22		
23	24	25	26	27	28	29		

2020

Monday

October

12

286-080 • WK 42

Advantage:- We can find relation with respect to Non linear Data.

Ques  
The Spearman correlation is less sensitive than the Pearson correlation to strong outliers that arise in tails of both samples.  
That is because Spearman's  $\rho$  limits the outlier to the value of ranks.

