

Correlation & Regression

Covariance:

If X and Y are r.v's then the Covariance b/w them is

$$\text{Cov}(X, Y) = E(XY) - E(X) \cdot E(Y)$$

Properties:

1. If X and Y are independent, then $\text{Cov}(X, Y) = 0$.
2. $\text{Cov}(aX, bY) = ab \text{Cov}(X, Y)$
3. $\text{Cov}(X+a, Y+b) = \text{Cov}(X, Y)$
4. $V(X_1 + X_2) = V(X_1) + V(X_2) + 2\text{Cov}(X_1, X_2)$
5. $V(X_1 - X_2) = V(X_1) + V(X_2) - 2\text{Cov}(X_1, X_2)$

Correlation:

Karl Pearson's Coefficient of Correlation:

Let X & Y be given r.v. Then Karl Pearson's Coefficient of Correlation is denoted by r_{xy} (or) $r(X, Y)$.

$$r_{xy} = r(X, Y) = \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}} = \frac{\text{Cov}(X, Y)}{\sigma_X \sigma_Y}$$

$$\begin{aligned} \text{Cov}(X, Y) &= E(XY) - E(X)E(Y) \\ &= \frac{1}{n} \sum xy - \bar{x} \bar{y} \quad , \quad \bar{x} = \frac{\sum x}{n} \quad \bar{y} = \frac{\sum y}{n} \end{aligned}$$

'n' - no. of items.

$$\sigma_X^2 = \text{Var}(X) = \frac{1}{n} \sum X^2 - (\bar{x})^2$$

$$\sigma_Y^2 = \text{Var}(Y) = \frac{1}{n} \sum Y^2 - (\bar{y})^2$$

Note: Correlation Coefficient always lies b/w -1 to $+1$.

Problem 1: Calculate the Correlation Coefficient for the following height (in inches) of father (X) and their son Y .

X	65	66	67	67	68	69	70	72
Y	67	68	65	68	72	72	69	71

Soln

X	Y	XY	X^2	Y^2
-----	-----	------	-------	-------

65	67	4355	4225	4489
66	68	4488	4356	4624
67	65	4355	4489	4225
67	68	4556	4489	4624
68	72	4896	4624	5184
69	72	4968	4761	5184
70	69	4830	4900	4761
72	71	5112	5184	5041
<hr/>				
544	552	37560	37028	38132

$$\bar{x} = \frac{544}{8} = 68 \quad \bar{y} = \frac{552}{8} = 69 \quad \bar{x}\bar{y} = 68 \cdot 69 = 4692$$

$$\sigma_x = \sqrt{\frac{1}{n} \sum x^2 - (\bar{x})^2}$$

$$= \sqrt{\frac{1}{8} (37028) - (68)^2}$$

$$\sigma_x = 2.121$$

$$\sigma_y = \sqrt{\frac{1}{n} \sum y^2 - (\bar{y})^2}$$

$$= \sqrt{\frac{1}{8} (38132) - (69)^2}$$

$$\sigma_y = 2.345$$

$$\begin{aligned} \text{Cov}(x, y) &= \frac{1}{n} \sum xy - \bar{x}\bar{y} \\ &= \frac{1}{8} (37560) - 68 \cdot 69 \\ &= 4695 - 4692 \\ &= 3 \end{aligned}$$

The Correlation Coefficient of x & y is given by

$$r_{xy} = r(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \frac{3}{(2.121)(2.345)} = \frac{3}{4.973} = 0.6032$$

$$r_{xy} = 0.6032.$$

Rank Correlation:

If (x_i, y_i) $i=1, 2, \dots, n$ are the rank of the individuals in two characteristics A & B respectively, then the Rank Correlation Coefficient is given by

$$r = 1 - \frac{6}{n(n^2-1)} \sum_{i=1}^n d_i^2$$

$$d_i = x_i - y_i \quad n \rightarrow \text{no of items}$$

Karl Pearson's formula for rank Correlation Coefficient.

Problem Find the rank Correlation Coefficient from the following data.

Rank in X	1	2	3	4	5	6	7
Rank in Y	4	3	1	2	6	5	7

X	Y	$d_i = x_i - y_i$	d_i^2
1	4	-3	9
2	3	-1	1
3	1	2	4
4	2	2	4
5	6	-1	1
6	5	1	1
7	7	0	0
			<u>20</u>

$$n = 7$$

\therefore rank Correlation Coefficient is

$$\begin{aligned}
 r &= 1 - \frac{6 \sum d_i^2}{n(n^2-1)} \\
 &= 1 - \frac{6(20)}{7(49-1)} \\
 &= 1 - \frac{120}{7 \times 48} \\
 &= 1 - \frac{5}{14} = \frac{9}{14} = 0.6428
 \end{aligned}$$

Problem 10 Participants were ranked according to their performance in a musical test by 3 judges in the following order.

Rank by X:	1	6	5	10	3	2	4	9	7	8
" " Y:	3	5	8	4	7	10	2	1	6	9
Z:	6	4	9	8	1	2	3	10	5	7

Use rank correlation method & discuss which pair of Judges has the nearest approach to common likings of music.

Soln

$$r(x, y) = 1 - \frac{6 \sum d_1^2}{n(n^2-1)}$$

$$r(z, x) = 1 - \frac{6 \sum d_2^2}{n(n^2-1)}$$

$$P(Y,Z) = 1 - \frac{6 \sum d_2^2}{n(n^2-1)}$$

$$P(Z,X) = 1 - \frac{6 \sum d_3^2}{n(n^2-1)}$$

X	Y	Z	$d_1 = x - y$	d_1^2	$d_2 = y - z$	d_2^2	$d_3 = z - x$	d_3^2
1	3	6	-2	4	-3	9	5	25
6	5	4	1	1	1	1	-2	4
5	8	9	-3	9	-1	1	4	16
10	4	8	6	36	-4	16	-2	4
3	7	1	-4	16	6	36	-2	4
2	10	2	-8	64	8	64	0	0
4	2	3	2	4	-1	1	-1	1
9	1	10	8	64	-9	81	1	1
7	6	5	1	1	1	1	-2	4
8	9	7	-1	1	2	4	-1	1
				$\sum d_1^2 = 200$			$\sum d_2^2 = 214$	$\sum d_3^2 = 60$

Here $n = 10$.

$$P(X,Y) = 1 - \frac{6(200)}{10 \times \frac{99}{33}} = 1 - \frac{40}{33} = \frac{33-40}{33} = \frac{-7}{33} = -0.212$$

$$P(Y,Z) = 1 - \frac{6(214)}{10 \times \frac{99}{33}} = 1 - \frac{214}{165} = \frac{165-214}{165} = \frac{-49}{165} = -0.296$$

$$P(Z,X) = 1 - \frac{6(60)}{10 \times \frac{99}{33}} = 1 - \frac{4}{11} = \frac{7}{11} = 0.6363$$

Since $P(X,Y)$ & $P(Y,Z)$ are negative & $P(Z,X)$ is positive
 Z, X Judges has the nearest approach to common liking of music.