

CORRELATION AND REGRESSION

Correlation:

If the change in one variable affects a change in the other variable, the variables are said to be correlated.

Two types of correlations are Positive correlation, Negative correlation

Positive Correlation

If the two variables deviate in the same direction

Eg: Height and Weight of a group of persons, Income and Expenditure.

Negative Correlation

If the two variables constantly deviate in opposite directions.

Eg: Price and Demand of a commodity, the correlation between volume and pressure of a perfect gas.

Covariance

If X and Y are random variables, then covariance between X and Y is defined as

$Cov(X, Y) = E(XY) - E(X).E(Y)$. If X and Y are independent then $Cov(X, Y)$

Karl- Pearson's coefficient of correlation

Correlation coefficient between two random variables X and Y , usually denoted by

$$r(X, Y) = \frac{COV(X, Y)}{\sigma_X \cdot \sigma_Y}$$

Where $COV(X, Y) = \frac{1}{n} \sum XY - \bar{X} \bar{Y}$

$$\sigma_X = \sqrt{\frac{1}{n} \sum X^2 - \bar{X}^2}, \bar{X} = \frac{\sum X}{n}$$

$$\sigma_Y = \sqrt{\frac{1}{n} \sum Y^2 - \bar{Y}^2}, \bar{Y} = \frac{\sum Y}{n}$$

Note

1. Correlation coefficient may also be denoted by $\rho(X, Y)$ or ρ_{XY}
2. If $\rho(X, Y) = 0$, We say that X and Y are uncorrelated.

Note

Types of correlation based on 'r'

| Value of 'r' | Correlation is said to be |
|--------------|---------------------------|
| $r = 1$ | Perfect and positive |
| $0 < r < 1$ | Positive |
| $-1 < r < 0$ | Negative |
| $r = 0$ | Uncorrelated |

Problems under Correlation

1. Calculate the correlation coefficient for the following heights (in inches) of fathers **X** and their sons **Y**.

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

Solution:

| X | Y | XY | X² | Y² |
|---------------------|------------------|---------------------|----------------------|-------------------------|
| 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 |
| $\sum (X)$ = 544 | $\sum (Y) = 552$ | $\sum (XY) = 37560$ | $\sum (X^2) = 37028$ | $\sum (Y^2)$ = 38132 |

$$\bar{X} = \frac{544}{8} = 68, \bar{Y} = \frac{552}{8} = 69$$

$$\bar{X} \bar{Y} = 68 \times 69 = 4692$$

$$\sigma_X = \sqrt{\frac{1}{n} \sum X^2 - \bar{X}^2} = 2.121$$

$$\sigma_Y = \sqrt{\frac{1}{n} \sum Y^2 - \bar{Y}^2} = 2.345$$

$$r(X, Y) = \frac{COV(X, Y)}{\sigma_X \cdot \sigma_Y} = \frac{\frac{1}{n} \sum XY - \bar{X}\bar{Y}}{\sigma_X \cdot \sigma_Y} = \frac{\frac{1}{8} 37560 - 4692}{2.121 \times 2.345} = 0.6031$$

2. Find the correlation coefficient between industrial production and export using the following data:

| | | | | | | | |
|----------------------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| Production(X) | 55 | 56 | 58 | 59 | 60 | 60 | 62 |
| EXPORT(Y) | 35 | 38 | 37 | 39 | 44 | 43 | 44 |

Solution:

| U | V | $U = X - 58$ | $V = Y - 40$ | UV | U ² | V ² |
|----|----|--------------|--------------|----|----------------|----------------|
| 55 | 35 | -3 | -5 | 15 | 9 | 25 |
| 56 | 38 | -2 | -2 | 4 | 4 | 4 |
| 58 | 37 | 0 | -3 | 0 | 0 | 9 |
| 59 | 39 | 1 | -1 | -1 | 1 | 1 |
| 60 | 44 | 2 | 4 | 8 | 4 | 16 |
| 60 | 43 | 2 | 3 | 6 | 4 | 9 |
| 62 | 44 | 4 | 4 | 16 | 16 | 16 |

| | | | | | | |
|--|--|---------------|---------------|--------------------|---------------------|---------------------|
| | | $\sum(U) = 4$ | $\sum(V) = 0$ | $\sum(UV)$ = 48 | $\sum(U^2)$ = 38 | $\sum(V^2)$ = 80 |
|--|--|---------------|---------------|--------------------|---------------------|---------------------|

Now $\bar{U} = \frac{\sum U}{n} = \frac{4}{7} = 0.5714$

$$\bar{V} = \frac{\sum V}{n} = 0$$

$$Cov(U, V) = \frac{\sum UV}{n} - \bar{U}\bar{V} = 6.857$$

$$\sigma_X = \sqrt{\frac{1}{n} \sum U^2 - \bar{U}^2} = 2.2588$$

$$\sigma_Y = \sqrt{\frac{1}{n} \sum V^2 - \bar{V}^2} = 3.38$$

$$r(U, V) = \frac{Cov(U, V)}{\sigma_U \cdot \sigma_V} = \frac{\frac{1}{n} \sum UV - \bar{U}\bar{V}}{\sigma_U \cdot \sigma_V} = 0.898$$

3. Two R.V.'S X and Y have joint p.d.f of $f(x, y) = \begin{cases} \frac{xy}{96} & 0 < x < 4, 1 < y < 5 \\ 0 & \text{elsewhere} \end{cases}$.

Find (i) $E(X)$ (ii) $E(Y)$ (iii) $E(XY)$ (iv) $E(2X + 3Y)$ (v) $Var(X)$ (vi) $Var(Y)$ (vii)

$Cov(X, Y)$

Solution:

(i) $E(X) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x f(x, y) dx dy$

$$\begin{aligned}
 &= \int_1^5 \int_0^4 x \frac{xy}{96} dx dy \\
 &= \frac{1}{96} \int_1^5 \int_0^4 x^2 y dx dy \\
 &= \frac{1}{96} \int_1^5 y \left(\frac{x^3}{3} \right)_0^4 dy \\
 &= \frac{64}{288} \int_1^5 y dy = \frac{2}{9} \left[\frac{y^2}{2} \right]_1^5 = \frac{24}{9}
 \end{aligned}$$

(ii) $E(Y) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y f(x,y) dx dy$

$$\begin{aligned}
 &= \int_1^5 \int_0^4 y \frac{xy}{96} dx dy \\
 &= \frac{1}{96} \int_1^5 \int_0^4 xy^2 dx dy \\
 &= \frac{1}{96} \int_1^5 y^2 \left(\frac{x^2}{2} \right)_0^4 dy \\
 &= \frac{1}{96(2)} \int_1^5 y^2 (4^2 - 0) dy
 \end{aligned}$$

$$\begin{aligned}
 &= \frac{16}{192} \left[\frac{y^3}{3} \right]_1^5 = \frac{124}{36} \\
 \Rightarrow E(Y) &= \frac{31}{9}
 \end{aligned}$$

(iii) $E(XY) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} xy f(x,y) dx dy$

$$= \int_1^5 \int_0^4 xy \left(\frac{xy}{96} \right) dx dy$$

$$\begin{aligned}
 &= \frac{1}{96} \int_1^5 \int_0^4 x^2 y^2 dx dy \\
 &= \frac{1}{96} \int_1^5 y^2 \left(\frac{x^3}{3} \right)_0^4 dy \\
 &= \frac{1}{96(3)} \int_1^5 y^2 (4^3 - 0) dy \\
 &= \frac{64}{288} \int_1^5 y^2 dy = \frac{2}{9} \left[\frac{y^3}{3} \right]_1^5 = \frac{248}{27}
 \end{aligned}$$

$$\Rightarrow E(XY) = \frac{248}{27}$$

$$(iv) E[2X + 3Y] = 2E[X] + 3E[Y]$$

$$\begin{aligned}
 &= 2\left(\frac{8}{3}\right) + 3\left(\frac{31}{9}\right) \\
 &= \frac{16+31}{3} = \frac{47}{3}
 \end{aligned}$$

$$(v) Var(X) = E(X^2) - [E(X)]^2$$

$$\text{Now, } E(X^2) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x^2 f(x, y) dx dy$$

$$= \int_1^5 \int_0^4 x^2 \left(\frac{xy}{96} \right) dx dy$$

$$= \frac{1}{96} \int_1^5 \int_0^4 x^3 y dx dy$$

$$= \frac{1}{96} \int_1^5 y \left(\frac{x^4}{4} \right)_0^4 dy$$

$$= \frac{1}{96(4)} \int_1^5 y (4^4 - 0) dy$$

$$= \frac{256}{384} \int_1^5 y dy = \frac{2}{3} \left[\frac{y^2}{2} \right]_1^5 = \frac{24}{3} = 8$$

$$\Rightarrow E(X^2) = 8$$

$$\text{Var}(X) = E(X^2) - [E(X)]^2 = 8 - \left(\frac{8}{3}\right)^2 = \frac{8}{9}$$

$$\sigma_X^2 = \frac{8}{9} \Rightarrow \sigma_X = \frac{\sqrt{8}}{3}$$

$$E(Y^2) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} y^2 f(x, y) dx dy$$

$$= \int_1^5 \int_0^4 y^2 \left(\frac{xy}{96}\right) dx dy$$

$$= \frac{1}{96} \int_1^5 \int_0^4 xy^3 dx dy$$

$$= \frac{1}{96} \int_1^5 y^3 \left(\frac{x^2}{2}\right)_0^4 dy$$

$$= \frac{1}{96(2)} \int_1^5 y^3 (4^2 - 0) dy$$

$$= \frac{16}{192} \int_1^5 y^3 dy$$

$$= \frac{1}{12} \left[\frac{y^4}{4} \right]_1^5 = \frac{624}{48} = 13$$

$$\Rightarrow E(Y^2) = 13$$

$$(vi) \text{Var}(Y) = E(Y^2) - [E(Y)]^2$$

$$= 13 - \left(\frac{31}{9}\right)^2$$

$$= \frac{92}{81}$$

$$\sigma_Y^2 = \frac{92}{81} \Rightarrow \sigma_Y = \frac{\sqrt{92}}{9}$$

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

$$= \frac{248}{27} - \frac{248}{27} = 0$$

4. If the independent random variables X and Y have the variances 36 and 16 respectively, find the correlation coefficient between $X + Y$ and $X - Y$

Solution:

Given that $\text{Var}(X) = 36$, $\text{Var}(Y) = 16$. Since X and Y are independent,

$$E(XY) = E(X) \cdot E(Y)$$

Let $U = X + Y$ and $V = X - Y$

$$\text{Var}(U) = \text{Var}(X + Y)$$

$$= 1^2 \text{Var}(X) + 1^2 \text{Var}(Y)$$

$$(\because \text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y))$$

$$= 36 + 16 = 52$$

$$\Rightarrow \sigma_U = \sqrt{52}$$

$$\text{Var}(V) = \text{Var}(X - Y)$$

$$= 1^2 \text{Var}(X) + (-1)^2 \text{Var}(Y)$$

$$(\because \text{Var}(aX + bY) = a^2 \text{Var}(X) + b^2 \text{Var}(Y))$$

$$= 36 + 16 = 52$$

$$\Rightarrow \sigma_V = \sqrt{52}$$

$$\text{Cov}(U, V) = E(UV) - E(U) \cdot E(V) \quad \dots (1)$$

$$E(UV) = E[(X + Y)(X - Y)]$$

$$= E[X^2 - Y^2]$$

$$= E(X^2) - E(Y^2) \quad \dots (2)$$

$$E(U) = E(X + Y) = E(X) + E(Y) \quad (3)$$

$$E(V) = E(X - Y) = E(X) - E(Y) \quad \dots (4)$$

Substituting (2), (3), (4) in (1) we get

$$\text{Cov}(U, V) = E(X^2) - E(Y^2) - [E(X) + E(Y)][E(X) - E(Y)]$$

$$= E(X^2) - E(Y^2) - [E(X)]^2 + [E(Y)]^2 - E(X)E(Y) + E(X)E(Y)$$

$$= \{E(X^2) - [E(X)]^2\} - \{E(Y^2) - [E(Y)]^2\}$$

$$= \text{Var}(X) - \text{Var}(Y)$$

$$\text{Cov}(U, V) = 36 - 16 = 20$$

$$\text{Hence } \rho(U, V) = \frac{\text{Cov}(U, V)}{\sigma_U \sigma_V} = \frac{20}{\sqrt{52} \sqrt{52}} = \frac{20}{52} = \frac{5}{13}$$

Regression

Regression is a mathematical measure of the average relationship between two or more variables in terms of the original limits of the data.

Lines of Regression

If the variables in a bivariate distribution are related we will find that the points in the scattered diagram will cluster around some curve called the curve of regression. If the curve is a straight line, it is called the line of regression and there is said to be linear regression between the variables, otherwise regression is said to be curvilinear.

The line of regression of Y on X is given by

$$y - \bar{y} = r \cdot \frac{\sigma_Y}{\sigma_X} (x - \bar{x})$$

Where r is the correlation coefficient, σ_Y and σ_X are standard deviations.

The line of regression of X on Y is given by

$$x - \bar{x} = r \cdot \frac{\sigma_X}{\sigma_Y} (y - \bar{y})$$

Note

Both the lines of regression passes through the mean (\bar{x}, \bar{y})

Angle between two lines of regression

If the equations of lines of regression of Y on X and X on Y are

$$y - \bar{y} = r \cdot \frac{\sigma_Y}{\sigma_X} (x - \bar{x})$$

$$x - \bar{x} = r \cdot \frac{\sigma_Y}{\sigma_X} (y - \bar{y})$$

The angle θ between the two lines of regression is given by

$$\tan\theta = \frac{1 - r^2}{r} \left(\frac{\sigma_Y \sigma_X}{\sigma_Y^2 + \sigma_X^2} \right)$$

Problems on Regression

1. From the following data, find (i) the two regression equations (ii) The coefficient of correlation between the marks in Economics and Statistics (iii) The most likely marks in statistics when marks in Economics are 30.

| | | | | | | | | | | |
|---------------------|----|----|----|----|----|----|----|----|----|----|
| Marks in Economics | 25 | 28 | 35 | 32 | 31 | 36 | 29 | 38 | 34 | 32 |
| Marks in Statistics | 43 | 46 | 49 | 41 | 36 | 32 | 31 | 30 | 33 | 39 |

Solution:

| X | Y | $X - \bar{X}$ | $Y - \bar{Y}$ | $(X - \bar{X})^2$ | $(Y - \bar{Y})^2$ | $(X - \bar{X})(Y - \bar{Y})$ |
|-----|-----|---------------|---------------|-------------------|-------------------|------------------------------|
| 25 | 43 | -7 | 5 | 49 | 25 | -35 |
| 28 | 46 | -4 | 8 | 16 | 64 | -32 |
| 35 | 49 | 3 | 11 | 9 | 121 | 33 |
| 32 | 41 | 0 | 3 | 0 | 9 | 0 |
| 31 | 36 | -1 | -2 | 1 | 4 | 2 |
| 36 | 32 | 4 | -6 | 16 | 36 | -24 |
| 29 | 31 | -3 | -7 | 9 | 49 | 21 |
| 38 | 30 | 6 | -8 | 36 | 64 | -48 |
| 34 | 33 | 2 | -5 | 4 | 25 | -10 |
| 32 | 39 | 0 | 1 | 0 | 1 | 0 |
| 320 | 380 | 0 | 0 | 140 | 398 | -93 |

Now $\bar{X} = \frac{\sum X}{n} = \frac{320}{10} = 32$

$\bar{Y} = \frac{\sum Y}{n} = \frac{380}{10} = 38$

Coefficient of regression of Y on X is $b_{YX} = \frac{\sum(X-\bar{X})(Y-\bar{Y})}{\sum(X-\bar{X})^2}$

$= -\frac{93}{140} = -0.6643$

$$\begin{aligned} \text{Coefficient of regression of X on Y is } b_{XY} &= \frac{\Sigma(x-\bar{x})(y-\bar{y})}{\Sigma(y-\bar{y})^2} \\ &= -\frac{93}{398} = -0.2337 \end{aligned}$$

Equation of the line of regression of X on Y is $x - \bar{x} = b_{XY}(y - \bar{y})$

$$\Rightarrow x - 32 = -0.2337(y - 38)$$

$$\Rightarrow x = -0.2337y + 0.2337 \times 38 + 32$$

$$\Rightarrow x = -0.2337y + 40.8806$$

Equation of the line of regression of Y on X is $y - \bar{y} = b_{YX}(x - \bar{x})$

$$\Rightarrow y - 38 = -0.6642(x - 32)$$

$$\Rightarrow y = -0.6642x + 0.6642 \times 32 + 38$$

$$\Rightarrow y = -0.6642x + 59.2576$$

Correlation of coefficient $r^2 = b_{YX} \times b_{XY}$

$$= -0.6643 \times (-0.2337)$$

$$= 0.1552$$

$$r = \pm\sqrt{0.1552}$$

$$= \pm 0.394$$

Now we have to find the most likely marks in statistics (Y) when marks in Economics (X) are 30. We use the line of regression of Y on X .

$$\Rightarrow y = -0.6642x + 59.2576$$

Put $x = 30$ we get

$$\Rightarrow y = -0.6642 \times 30 + 59.2576$$

$$\Rightarrow y = 39.3286$$

2. The two lines of regression are $8x - 10y + 66 = 0$, $40x - 18y - 214 = 0$. The variance of X is 9. Find (i) the mean value of X and Y (ii) correlation coefficient between X and Y .

Solution:

Since both the lines of regression passes through the mean values \bar{x} and \bar{y} , the point (\bar{x}, \bar{y}) must satisfy the two given regression lines.

$$(1) \times 5 \Rightarrow 40\bar{x} - 50\bar{y} = -330$$

$$(2) \Rightarrow 40\bar{x} - 18\bar{y} = 214$$

Subtracting (1) - (2) we get

$$\Rightarrow 32\bar{y} = 544$$

$$\Rightarrow \bar{y} = 17$$

Sub $\bar{y} = 17$ in (1) we get,

$$(1) \Rightarrow 8\bar{x} - 10\bar{y} = -66$$

$$\Rightarrow 8\bar{x} - 10 \times 17 = -66$$

$$\Rightarrow 8\bar{x} = -66 + 170$$

$$\Rightarrow \bar{x} = 13$$

Hence the mean value are given by $\bar{x} = 13$ and $\bar{y} = 17$

(ii) Let us suppose that equation (A) is the equation of line of regression of Y on X and (B) is the equation of the line regression of X on Y, we get after rewriting (A) and (B)

$$\Rightarrow 10y = 8x + 66$$

$$\Rightarrow y = \frac{8}{10}x + \frac{66}{10}$$

$$\Rightarrow b_{YX} = \frac{8}{10}$$

$$\Rightarrow 40x = 18y + 214$$

$$\Rightarrow x = \frac{18}{40}y + \frac{214}{40}$$

$$\Rightarrow b_{XY} = \frac{18}{40}$$

Correlation of coefficient $r^2 = b_{YX} \times b_{XY}$

$$= \frac{8}{10} \times \frac{18}{40} = \frac{9}{25}$$

$$\Rightarrow r = \pm \frac{3}{5}$$

$$= \pm 0.6$$

Since both the regression coefficients are positive, r must be positive.

Hence $r = 0.6$

Important note:

If we take equation (A) as the line of regression of X on Y we get,

$$\Rightarrow 8x = 10y - 66$$

$$\Rightarrow x = \frac{10}{8}y - \frac{66}{8}$$

$$\Rightarrow b_{XY} = \frac{10}{8}$$

$$\Rightarrow 18y = 40x - 214$$

$$\Rightarrow y = \frac{40}{18}x - \frac{214}{18}$$

$$\Rightarrow b_{YX} = \frac{40}{18}$$

Correlation of coefficient $r^2 = b_{YX} \times b_{XY}$

$$= \frac{10}{8} \times \frac{40}{8} = \frac{25}{9}$$

$$r = 2.78$$

But r^2 should always lies between 0 and 1. Hence our assumption that line (A) is line of regression of X on Y and the line (B) is line of regression of Y on X is wrong.

3. The two lines of regression are $4x - 5y + 33 = 0, 20x - 9y - 107 = 0$. The variance of X is 25. Find (i) the mean value of X and Y (ii) correlation coefficient between X and Y.

Solution:

Since both the lines of regression passes through the mean values \bar{x} and \bar{y} , the point (\bar{x}, \bar{y}) must satisfy the two given regression lines.

$$(1) \Rightarrow 20\bar{x} - 9\bar{y} = 107$$

$$(2) \times 5 \Rightarrow 20\bar{x} - 25\bar{y} = -165$$

Subtracting (1) - (2) we get

$$\Rightarrow 16\bar{y} = 272$$

$$\Rightarrow \bar{y} = 17$$

Sub $\bar{y} = 17$ in (1) we get,

$$(2) \Rightarrow 4\bar{x} - 5\bar{y} = -33$$

$$\Rightarrow 4\bar{x} - 5 \times 17 = -33$$

$$\Rightarrow 4\bar{x} = -33 + 85$$

$$\Rightarrow \bar{x} = 13$$

Hence the mean value as given by $\bar{x} = 13$ and $\bar{y} = 17$

(ii) Let us suppose that equation (A) is the equation of line of regression of Y on X and (B) is the equation of the line regression of X on Y, we get after rewriting (A) and (B)

$$\Rightarrow 5y = 4x + 33$$

$$\Rightarrow y = \frac{4}{5}x + \frac{33}{5}$$

$$\Rightarrow b_{YX} = \frac{4}{5}$$

$$\Rightarrow 20x = 9y + 107$$

$$\Rightarrow x = \frac{9}{20}y + \frac{107}{20}$$

$$\Rightarrow b_{XY} = \frac{9}{20}$$

Correlation of coefficient $r^2 = b_{YX} \times b_{XY}$

$$= \frac{4}{5} \times \frac{9}{20} = \frac{3}{5}$$

$$r = \pm 0.6$$

4. Can $Y = 5 + 2.8X$ and $X = 3 - 0.5Y$ be the estimated regression equations of Y on X and X on Y respectively? Explain your answer.

Solution:

Given,

$$\Rightarrow X = 3 - 0.5Y$$

$$\Rightarrow b_{XY} = -0.5$$

$$\Rightarrow Y = 5 + 2.8X$$

$$\Rightarrow b_{YX} = 2.8$$

Correlation of coefficient $r^2 = b_{YX} \times b_{XY}$

$$= 2.8 \times (-0.5) = -1.4$$

$$r = \sqrt{-1.4} \text{ which is imaginary quantity.}$$

Here r cannot be imaginary.

Hence the given lines are not estimated as regression equations.

