

Regression

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

LINES OF REGRESSION:

* The line of Regression of X on Y is

$$X - \bar{X} = r \frac{\sigma_x}{\sigma_y} (Y - \bar{Y})$$

(or)

$$X - \bar{X} = b_{xy} (Y - \bar{Y})$$

* The line of Regression of Y on X is

$$Y - \bar{Y} = r \frac{\sigma_y}{\sigma_x} (X - \bar{X})$$

(or)

$$Y - \bar{Y} = b_{yx} (X - \bar{X})$$

where b_{xy} , b_{yx} are regression coefficients.

Correlation coefficient $r_{xy} = \pm \sqrt{b_{xy} \times b_{yx}}$

Properties:

1. When b_{xy} and b_{yx} are +ve r is +ve
2. When b_{xy} and b_{yx} both are negative r is -ve
3. When $r = \pm 1$ the regression lines either parallel or coincide.
4. When $r = 0$, the regression lines are perpendicular to each other.

4. The regression lines always passes through (\bar{x}, \bar{y})

5. If one of the regression coefficient is greater than 1 and another one must be less than one.

6. The Geometric mean of the regression coefficient is correlation coefficient.

$$r_{xy} = \pm \sqrt{b_{yx} b_{xy}}$$

7. The arithmetic mean of the regression coefficient is greater than correlation coefficient

$$\text{i.e., } \frac{b_{yx} + b_{xy}}{2} > r_{xy}$$

8. The regression coefficients are independent of change of origin but not of scale.

$$\text{i.e., } b_{xy} = b_{uv} \text{ where } u = x - a, \text{ only.}$$

9. The angle between two lines of regression is

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

Proof: WKT $\theta = \tan^{-1} \left(\frac{m_2 - m_1}{1 + m_1 m_2} \right)$

Target equation of the straight line is

$$y - y_1 = \frac{-1}{m} (x - x_1)$$

Let the regression of y on x is

$$y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$-\frac{1}{m} = r \frac{\sigma_y}{\sigma_x}$$

$$\Rightarrow m_1 = \frac{-\sigma_x}{r\sigma_y}$$

Let the regression of x on y is

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$\Rightarrow y - \bar{y} = \frac{\sigma_y}{r\sigma_x} (x - \bar{x})$$

$$\Rightarrow -\frac{1}{m} = \frac{\sigma_y}{r\sigma_x}$$

$$\Rightarrow m_2 = -\frac{r\sigma_x}{\sigma_y}$$

$$\theta = \tan^{-1} \left(\frac{-\frac{r\sigma_x}{\sigma_y} + \frac{\sigma_x}{r\sigma_y}}{1 + \frac{\sigma_x}{r\sigma_y} \times r \frac{\sigma_x}{\sigma_y}} \right)$$

$$= \tan^{-1} \left(\frac{\frac{-r^2\sigma_x + \sigma_x}{r\sigma_y}}{\frac{\sigma_y^2 + \sigma_x^2}{\sigma_y^2}} \right)$$

$$= \tan^{-1} \frac{\sigma_x(1-r^2)\sigma_y}{r(\sigma_x^2 + \sigma_y^2)}$$

$$\therefore \theta = \tan^{-1} \left[\left(\frac{1-r^2}{r} \right) \frac{\sigma_x\sigma_y}{\sigma_x^2 + \sigma_y^2} \right]$$

* The two regression lines are $8x - 10y + 66 = 0$ and $40x - 18y - 214 = 0$. Find (\bar{x}, \bar{y}) . Given $v(x) = 9$, find $v(y)$, find r_{xy} .

Soln:

$$\text{Gn } 8x - 10y = -66$$

$$40x - 18y = 214$$

WKT the regression lines always pass thro (\bar{x}, \bar{y}) .

$$8\bar{x} - 10\bar{y} = -66$$

$$40\bar{x} - 18\bar{y} = 214$$

on Solving $\bar{x} = 13, \bar{y} = 17$

Assume $8x - 10y = -66$ as regression of x on y

$$8x = 10y - 66$$

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

Assume $40x - 18y = 214$ as regression of y on x

$$-18y = -40x + 214$$

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

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

$$\text{let } r = \pm \sqrt{b_{xy} \times b_{yx}}$$

$$= \pm \sqrt{\frac{10}{8} \times \frac{40}{18}}$$

$r \geq 1$ the assumption is wrong.

let $40x - 18y = 214$ as regression of x on y

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

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

$$\Rightarrow b_{xy} = \frac{18}{40} = \frac{9}{20}$$

let $8x - 10y = -66$ as regression of y on x

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

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

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

$$r = \pm \sqrt{b_{xy} \times b_{yx}} = \pm \sqrt{\frac{9}{20} \times \frac{8}{10}}$$

$$r = \pm 0.6$$

$$\boxed{r = 0.6}$$

Given $V(x) = 9$, find $V(y)$

$$\Rightarrow \sigma_x = 3$$

$$\text{let } b_{xy} = r \frac{\sigma_x}{\sigma_y}$$

$$\frac{18}{40} = (0.6) \frac{3}{\sigma_y}$$

$$\sigma_y = \frac{(0.6) \times 3 \times 40}{18}$$

$$\sigma_y = 4$$

$$\Rightarrow V(y) = 16$$

2. Out of two lines of regression $X+2Y=5$, $2X+3Y=8$
- (i) Find (\bar{X}, \bar{Y}) (ii) which one is the regression of X on Y
- (iii) If $v(X)=12$, find $v(Y)$.

Soln: Since the regression lines always pass through (\bar{x}, \bar{y}) .

$$\bar{X} + 2\bar{Y} = 5$$

$$2\bar{X} + 3\bar{Y} = 8$$

on solving $\bar{X} = 1, \bar{Y} = 2$

let us assume $X+2Y=5$ is the regression of Y on X .

$$2Y = -X + 5$$

$$Y = -\frac{1}{2}X + \frac{5}{2}$$

$$\Rightarrow b_{YX} = -\frac{1}{2}$$

let $2X+3Y=8$ as the regression of X on Y .

$$2X = -3Y + 8$$

$$X = -\frac{3}{2}Y + \frac{8}{2}$$

$$\Rightarrow b_{XY} = -\frac{3}{2}$$

$$\text{let } r = \pm \sqrt{b_{YX} \times b_{XY}} = \pm \sqrt{\left(-\frac{3}{2}\right)\left(-\frac{1}{2}\right)}$$

$$r = \pm \sqrt{\frac{3}{4}} = \pm 0.86$$

$$\Rightarrow r = -0.86$$

(ii) $2X+3Y$ is the regression of X on Y .

(ii) Given $v(x) = 12$

Wkt $b_{xy} = r \frac{\sigma_x}{\sigma_y}$

$$\therefore -\frac{3}{2} = -0.86 \times \frac{\sqrt{12}}{\sigma_y}$$

$$\sigma_y = \frac{0.86 \times 3.46}{1.5} = 1.98 \approx 2$$

$$\sigma_y = 2 \Rightarrow v(y) = 4.$$

* Obtain equations of lines of regression from the following data. Obtain the value of y when $x=71$. Also obtain the value of x when $y=70$.

x :	65	66	67	67	68	69	70	72
y :	67	68	65	68	72	72	69	71

Soln:

$$\bar{x} = \frac{\sum x}{n} = 68 = E(x)$$

$$E(xy) = 4695$$

$$\bar{y} = \frac{\sum y}{n} = 69 = E(y)$$

$$\sigma_x = \sqrt{E(x^2) - (E(x))^2} = \sqrt{4628.5 - (68)^2} = 2.12$$

$$\sigma_y = \sqrt{E(y^2) - (E(y))^2} = \sqrt{4766.5 - (69)^2} = 2.34$$

$$r_{xy} = \frac{E(xy) - E(x)E(y)}{\sigma_x \sigma_y} = 0.6014$$

Regression of X on Y is

$$X - \bar{X} = r \frac{\sigma_x}{\sigma_y} (Y - \bar{Y})$$

$$X - 68 = (0.604) \frac{2.12}{2.34} (Y - 69)$$

$$X - 68 = 0.5472 (Y - 69) = 0.5472Y - 37.75$$

$$X - 0.5472Y = 30.25 \rightarrow \textcircled{1}$$

Regression of Y on X is

$$Y - \bar{Y} = r \frac{\sigma_y}{\sigma_x} (X - \bar{X})$$

$$Y - 69 = (0.604) \frac{2.34}{2.12} (X - 68)$$

$$Y - 69 = 0.666X - 45.33$$

$$Y - 0.666X = 69 - 45.33$$

$$Y - 0.666X = 23.67 \rightarrow \textcircled{2}$$

To find X when $Y = 70$

put $Y = 70$ in $\textcircled{1}$

$$X - 0.5472(70) = 30.25$$

$$X - 38.29 = 30.25$$

$$\boxed{X = 68.54}$$

To find Y when $X = 71$

put $X = 71$ in $\textcircled{2}$

$$Y - 0.666(71) = 23.67$$

$$\Rightarrow Y = 70.956 \approx 71$$

* The regression equation of X on Y is $3Y - 5X + 108 = 0$. If the mean value of Y is 44 and the $V(X) = \left(\frac{9}{16}\right)^{\text{th}} V(Y)$ find the mean value of X . Also find the correlation coefficient

Soln:

Regression of X on Y is $3Y - 5X + 108 = 0$.

Given $\bar{Y} = 44$. Since the regression line passes through (\bar{X}, \bar{Y})

$$3\bar{Y} - 5\bar{X} + 108 = 0$$

$$3(44) - 5\bar{X} + 108 = 0$$

$$182 - 5\bar{X} + 108 = 0$$

$$240 - 5\bar{X} = 0$$

$$-5\bar{X} = -240$$

$$\bar{X} = 48.$$

Since $3Y - 5X + 108 = 0$ is the regression of X on Y

$$-5X = -3Y - 108$$

$$X = \frac{3}{5}Y + \frac{108}{5}$$

$$b_{xy} = \frac{3}{5}$$

$$\text{Wkt } b_{xy} = r \frac{\sigma_x}{\sigma_y}$$

$$V(X) = \frac{9}{16} V(Y)$$

$$\Rightarrow \sigma_x = \frac{3}{4} \sigma_y$$

$$\frac{3}{5} = r \frac{\frac{3}{4} \sigma_y}{\sigma_y}$$

$$\therefore \frac{3}{5} = r \frac{3}{4}$$

$$\Rightarrow r = \frac{4}{5} = 0.8$$

* The tangent angle between the lines of regression of Y on X and X on Y is 0.6 and $\sigma_x = \frac{1}{2}\sigma_y$. Find r_{xy} .

Soln:

$$\tan \theta = 0.6$$

$$\sigma_y = 2\sigma_x$$

$$\text{Hkt } \tan \theta = \frac{1-r^2}{r} \cdot \frac{\sigma_x \sigma_y}{\sigma_x^2 + \sigma_y^2}$$

$$0.6 = \left(\frac{1-r^2}{r}\right) \frac{\sigma_x \cdot 2\sigma_x}{\sigma_x^2 + 4\sigma_x^2}$$

$$0.6 = \frac{1-r^2}{r} \cdot \frac{2\cancel{\sigma_x^2}}{5\cancel{\sigma_x^2}}$$

$$0.6 = \left(\frac{1-r^2}{r}\right) \frac{2}{5}$$

$$(0.6)5r = 2(1-r^2)$$

$$3r = 2 - 2r^2$$

$$2r^2 + 3r - 2 = 0$$

$$r = 0.5, -2$$

$$\boxed{r = 0.5}$$

* If X and Y are jointly distributed with correlation

$r_{xy} = \frac{1}{2}$, $\sigma_x = 2$, $\sigma_y = 3$. Find the variance

$$2X - 4Y + 3.$$

Soln:

$$\text{Var}^v(2X - 4Y + 3) = \text{Var}^v(2X - 4Y) + \text{Var}^v(3)$$

$$= 4\text{Var}(X) + 16\text{Var}(Y) - 2(2)(4)\text{Cov}(X, Y)$$

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

$$\begin{aligned} \text{Cov}(x, y) &= r_{xy} \sigma_x \sigma_y \\ &= \frac{1}{2} \cdot 2 \cdot 3 \end{aligned}$$

$$\Rightarrow \text{Cov}(x, y) = 3.$$

$$\therefore V(2x - 4y + 3) = 4(4) + 16(9) - 16(3) = 112$$

* From the following data, find the most likely price at Madras corresponding to the price 70 at Bombay and that at Bombay corresponding to the price 68 at Madras.

	Madras	Bombay
Average Price	65	67
SD of Price	0.5	3.5

Standard difference b/w the price at Madras and Bombay is 3:1.

Soln: WKT $r = \frac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x\sigma_y}$

$$\begin{aligned} r &= \frac{(0.5)^2 + (3.5)^2 - (3.1)^2}{2(0.5)(3.5)} \\ &= \frac{0.25 + 12.25 - 9.61}{3.5} \end{aligned}$$

$$\Rightarrow r = 0.825$$

Regression of x on y is

$$x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$$

$$x - 65 = (0.825) \left(\frac{0.5}{3.5} \right) (y - 67)$$

$$x - 65 = (0.117) (y - 67)$$

When Bombay = 70 Madras = ?

$$x - 65 = (0.117) (70 - 67)$$

$$x - 65 = 0.351$$

$$\Rightarrow x = 65.351$$

Regression of y on x is

$$y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$$

$$y - 67 = (0.825) \left(\frac{3.5}{0.5} \right) (x - 65)$$

$$y - 67 = (5.775) (x - 65)$$

When $x = 68$ (Madras); Bombay = ?
 $y = 84.325$