

Correlation and Regression

↳ measure the degree of association between variables

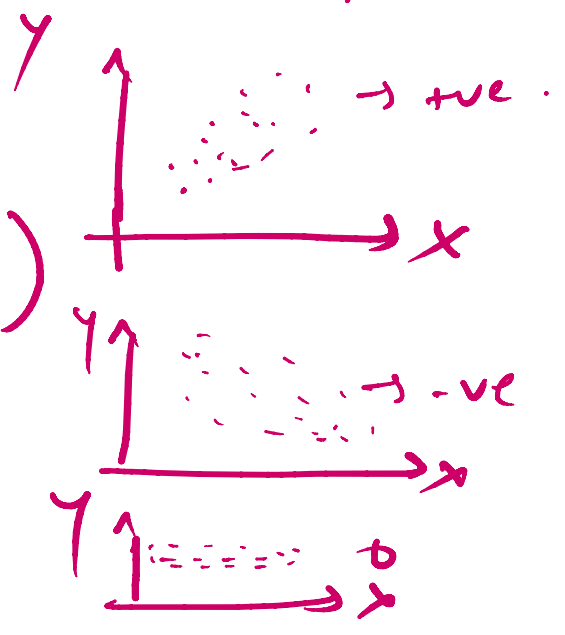
↳ bivariate data & relation between two variables

↑
relation

one indep → x
one dep → y

Graphical representation of this degree of association between x and y → Scatter Diagram.

- ① +ve correlation
- ② -ve correlation
- ③ uncorrelated (0)
(No relation between x & y).



— st —

Mathematically, the measure of degree of association between two variables x and y is correlation coefficient, denoted by

is correlation coefficient denoted by

$$r_{x,y} = \frac{\text{cov}(x,y)}{\sqrt{v(x)}\sqrt{v(y)}}$$

$$\text{Now, } \sqrt{\text{cov}(x,y)} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$$

$$\text{or, } \text{cov}(x,y) = \frac{1}{n} \sum_{i=1}^n x_i y_i - \bar{x} \cdot \bar{y}$$

$$\text{s.d, } \sigma_x = \sqrt{v(x)} = \sqrt{\frac{1}{n} \sum (x_i - \bar{x})^2} = \sqrt{\frac{1}{n} \sum x_i^2 - \bar{x}^2}$$

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

Properties: (i) $-1 \leq r_{x,y} \leq +1$

(ii) correlation coefficient do not depend on the change in origin and scale of variable.

ie if x and y are changed as $U = \frac{x - c}{d}$ and $V = \frac{y - c'}{d'}$

then $r_{x,y} = r_{u,v}$

Proof

(i) $-1 \leq r_{x,y} \leq +1$.

let x and y assumes the value (x_i, y_i) ,

let x and y assumes the values (x_1, y_1) , (x_2, y_2) , (x_3, y_3) , \dots , (x_n, y_n) , where

$\bar{x} = \frac{1}{n} \sum x_i$ and $\bar{y} = \frac{1}{n} \sum y_i$ the resp
mean of x and y
and σ_x and σ_y are the resp s.d.

Now let us write

$$u_i = \frac{x_i - \bar{x}}{\sigma_x}$$

$$\text{and } v_i = \frac{y_i - \bar{y}}{\sigma_y}$$

$$\begin{aligned} \sum_{i=1}^n u_i^2 &= \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma_x} \right)^2 \\ &= \frac{\sum (x_i - \bar{x})^2}{n \sigma_x^2} \end{aligned}$$

$$\begin{aligned} \text{and } \sum v_i^2 &= \sum \left(\frac{y_i - \bar{y}}{\sigma_y} \right)^2 \\ &= \frac{\sum (y_i - \bar{y})^2}{n \sigma_y^2} \end{aligned}$$

$$\sum u_i^2 = \frac{\sigma_x^2}{\sigma_x^2} = 1$$

$$\sum v_i^2 = \frac{\sigma_y^2}{\sigma_y^2} = 1$$

$$\begin{aligned} \text{and } \sum u_i v_i &= \sum_{i=1}^n \left(\frac{x_i - \bar{x}}{\sigma_x} \right) \left(\frac{y_i - \bar{y}}{\sigma_y} \right) \\ &= \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n \sigma_x \cdot \sigma_y} \\ &= \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y} \end{aligned}$$

$$\boxed{\sum u_i v_i = r_{x,y}} \quad \frac{\sum u_i v_i}{\sqrt{\sum u_i^2} \sqrt{\sum v_i^2}}$$

So, $\sum (u_i + v_i)^2 \geq 0$

$$\sum u_i^2 + \sum v_i^2 + 2 \sum u_i v_i \geq 0$$

$$1 + 1 + 2 \cdot r_{x,y} \geq 0$$

$$2 + 2 r_{x,y} \geq 0$$

$$2 (1 + r_{x,y}) \geq 0$$

$\Rightarrow 1 + r_{x,y} \geq 0$

$\therefore \boxed{r_{x,y} \geq -1} \quad \text{--- (1)}$

Again $\sum (u_i - v_i)^2 \geq 0$

$$\sum u_i^2 + \sum v_i^2 - 2 \sum u_i v_i \geq 0$$

$$1 + 1 - 2 r_{x,y} \geq 0$$

$$2 (1 - r_{x,y}) \geq 0$$

$\Rightarrow 1 - r_{x,y} > 0$

$\Rightarrow \boxed{r_{x,y} \leq 1} \quad \text{--- (2)}$

Comparing (1) and (2) $-1 \leq r_{x,y} \leq 1$ (Proved)

Proof (ii)

Let x and y be the variables
and $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be
the 'n' pair of observations.

Now let us change the origin and scale of
 x variable by 'a' and 'b' units

$$\text{then } u_i = \frac{x - a}{b}$$

Similarly if y variable is changed by 'c' units
of origin and 'd' units of scale then,

$$v_i = \frac{y - c}{d}$$

Here from the given observation of (x_i, y_i)
we have \bar{x}, \bar{y} as mean
and σ_x, σ_y as s.d of x and y resp

$$\therefore r_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y}$$

Since $u_i = \frac{x_i - a}{b}$

$$\Rightarrow x_i = a + b u_i$$

$$\therefore \bar{x} = a + b \bar{u}$$

$$\dots - 2 \dots (1 \dots - \bar{x})^2$$

Similarly $v_i = \frac{y_i - c}{d}$

$$\bar{y} = c + d \bar{v}$$

$$\begin{aligned} \sigma_x^2 &= \frac{1}{n} \sum (x_i - \bar{x})^2 \\ &= \frac{1}{n} \sum (\mu + b u_i - \mu - b \bar{u})^2 \\ &= \frac{1}{n} \cdot b^2 \sum (u_i - \bar{u})^2 \end{aligned}$$

$$\therefore \sigma_x^2 = b^2 \sigma_u^2$$

$$\sigma_x = |b| \sigma_u$$

Similarly $\sigma_y = |d| \sigma_v$

$$\begin{aligned} \text{Cov}(x, y) &= \frac{1}{n} \sum (x_i - \bar{x})(y_i - \bar{y}) \\ &= \frac{1}{n} \sum (\mu + b u_i - \mu - b \bar{u})(\mu + d v_i - \mu - d \bar{v}) \\ &= \frac{bd}{n} \sum (u_i - \bar{u})(v_i - \bar{v}) \\ &= b \cdot d \text{Cov}(u, v) \end{aligned}$$

\therefore Correlation coefficient between x, y is $r_{x,y} = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} = \frac{b \cdot d \cdot \text{Cov}(u,v)}{|b| \sigma_u \cdot |d| \sigma_v}$

$$r_{x,y} = \frac{bd}{|b||d|} \cdot r_{u,v}$$

if b and d are of same sign

$$\text{then } r_{x,y} = r_{u,v}$$

and if b and d are of opposite sign

$$\text{then } r_{x,y} = -r_{u,v}$$

$$\left. \begin{array}{l} \sum_{i=1}^{100} x_i = 280 \\ \sum_{i=1}^{100} y_i = 60 \end{array} \right\} \text{ calculate}$$

$$\begin{aligned} \sum_{i=1}^n x_i &= 280 & \sum_{i=1}^n y_i &= 60 \\ \sum x_i^2 &= 2384 & \sum y_i^2 &= 117 \\ \sum x_i y_i &= 438 \end{aligned}$$

Calculate $r_{x,y}$.

$$\bar{x} = \frac{1}{n} \sum x_i = \frac{1}{100} \times 280 = 2.8$$

$$\bar{y} = \frac{1}{n} \sum y_i = \frac{60}{100} = 0.6$$

$$\begin{aligned} \sigma_x &= \sqrt{\frac{1}{n} \sum x_i^2 - \bar{x}^2} = \sqrt{\frac{2384}{100} - 2.8^2} \\ &= \sqrt{16} \end{aligned}$$

$$\sigma_y = \sqrt{\frac{1}{n} \sum y_i^2 - \bar{y}^2} = \sqrt{0.81} = 0.9$$

$$\text{cov}(x,y) = \frac{1}{n} \sum xy - \bar{x}\bar{y}$$

$$\begin{aligned} &= \frac{438}{100} - (2.8 \times 0.6) = 4.38 - (1.68) \\ &= 2.70 \end{aligned}$$

$$\therefore r_{x,y} = \frac{\text{cov}(x,y)}{\sigma_x \sigma_y} = \frac{2.7}{4 \times 0.9} = \underline{\underline{0.75}}$$

x	65	63	67	69	68	62	70	66
y	68	66	68	65	69	66	68	65

	x	65	63	67	69	68	62	70	66
	y	68	66	68	65	69	66	68	65

Calculate r_{xy} .

x	y	x^2	y^2	xy
Σx	Σy	Σx^2	Σy^2	Σxy

↑ mean ↑ s.d ↑ cov. ∴ $r_{xy} = \text{ans.}$

Q While calculating the coefficient of correlation between two variables x and y , the following results were obtained:

$$n = 25 \quad \Sigma x = 125 \quad \Sigma y = 100$$

$$\Sigma x^2 = 650 \quad \Sigma y^2 = 460$$

$$\Sigma xy = 508$$

It was however later discovered at the time of checking that two pairs of obs (x, y) were copied (5, 14) and (8, 6) while the correct values were (8, 12) and (6, 8)

copied (6, 17) and (8, 12) were
 the correct values were (8, 12) and (6, 8)
 Susp. Determine the correct coefficient

of correlation.

From the question we have

$$\sum x = 125$$

We will subtract incorrect obs and add the correct obs, we get.

$$\text{Corrected } \sum x = 125 - 6 - 8 + 8 + 6$$

$$\therefore \text{corrected mean, } \bar{x} = \frac{1}{25} \times \frac{125}{5} = 5 \checkmark$$

Again $\sum y = 100$

after subtraction , we get-

$$\text{Corrected } \sum y = 100 - 14 - 6 + 12 + 8 = 100$$

$$\therefore \text{corrected mean } \bar{y} = \frac{1}{25} \times 100 = 4 \checkmark$$

given $\sum x^2 = 650$

$$\text{corrected } \sum x^2 = 650 - 6^2 - 8^2 + 8^2 + 6^2 = 650$$

$$\therefore \sigma_x = \sqrt{\frac{1}{n} \sum x^2 - \bar{x}^2} = \sqrt{\frac{650}{25} - 5^2} = \sqrt{26 - 25} = \sqrt{1} = 1$$

and $\sum y^2 = 460$

$$\text{correct } \sum y^2 = 460 - 14^2 - 6^2 + 12^2 + 8^2 = 436$$

$$\begin{aligned}\sigma_y &= \sqrt{\frac{1}{25} \times 436 - 4^2} \\ &= \sqrt{1.44} \\ \sigma_y &= 1.2\end{aligned}$$

Given

$$\sum xy = 508$$

$$\text{Correct } \sum xy = 508 - (6 \times 14) - (8 \times 6) + (8 \times 12) + (6 \times 8)$$

$$\therefore \sum xy = 520$$

$$\therefore \text{cov}(x, y) = \frac{1}{5} \times 520 - (5 \times 4)$$

$$= \frac{104}{5} - 20$$

$$= \frac{104 - 100}{5}$$

$$= \frac{4}{5}$$

$$\therefore r_{xy} = \frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = \frac{\frac{4}{5}}{1 \times 1.2} = \frac{4 \times 0.8}{5 \times 1.2} = \frac{8}{12} = \frac{2}{3} \quad (\text{ans})$$

— x —