# Correlation Analysis -

Dr. Uday Kashid

- Karl pearson's coefficient of correlation [or] product moment coefficient of correlation :- (r)

① $\quad r = \dfrac{cov(x,y)}{\sigma_x \cdot \sigma_y} = \dfrac{\Sigma(x-\bar{x})(y-\bar{y})}{N \cdot \sigma_x \cdot \sigma_y}$

where $\bar{x} = \dfrac{\Sigma X}{N}$, $\bar{y} = \dfrac{\Sigma Y}{N}$

$\sigma_x^2 = $ variance of $x = \dfrac{\Sigma(x-\bar{x})^2}{N}$

$\sigma_y^2 = $ variance of $y = \dfrac{\Sigma(y-\bar{y})^2}{N}$

and $\sigma_x = \sqrt{\sigma_x^2} = \sqrt{V(x)} = S.D$ of $x$

$\sigma_y = \sqrt{\sigma_y^2} = \sqrt{V(y)} = S.D$ of $Y$.

**or**

② $\quad r = \dfrac{\Sigma(x-\bar{x})(y-\bar{y})}{\sqrt{\Sigma(x-\bar{x})^2}\ \sqrt{\Sigma(y-\bar{y})^2}}$

[Use when $\bar{x}$ & $\bar{y}$ are in integer form]

**or** Direct method :-

③ $\quad r = \dfrac{N\Sigma xy - \Sigma x \cdot \Sigma y}{\sqrt{N\Sigma(x^2) - [\Sigma(x)]^2}\ \sqrt{N\Sigma(y^2) - (\Sigma y)^2}}$

**or**

④ $\quad r = \dfrac{\Sigma(xy) - N(\bar{x})(\bar{y})}{\sqrt{\Sigma(x^2) - N(\bar{x})^2} \cdot \sqrt{\Sigma(y^2) - N(\bar{y})^2}}$

**or**

⑤ Assumed mean Formula :

IF $A$ & $B$ are assumed means of $x$ & $y$ respectively then

$r = \dfrac{\Sigma(x-A)(y-B) - \dfrac{\Sigma(x-A) \cdot \Sigma(y-B)}{N}}{\sqrt{\Sigma[(x-A)^2] - \dfrac{[\Sigma(x-A)]^2}{N}}\ \sqrt{\Sigma[(y-B)^2] - \dfrac{[\Sigma(y-B)]^2}{N}}}$

Note: ① $-1 \le r \le 1$

② If $r = 1$ (perfect +ve correlation) ③ If $r = -1$ (perfect (-ve) correlation)

(11)

| X | 57 | 42 | 38 | 42 | 45 | 42 | 44 | 40 | 46 | 44 | 43 | 40 |
|---|----|----|----|----|----|----|----|----|----|----|----|----|
| Y | 10 | 26 | 41 | 29 | 27 | 27 | 19 | 18 | 19 | 31 | 29 | 33 |

Ans
$r = -0.74$

Direct calculating method of $r$

$$r = \frac{\Sigma xy - N(\bar{x})(\bar{y})}{N\sqrt{\frac{\Sigma x^2}{N} - (\bar{x})^2}\sqrt{\frac{\Sigma y^2}{N} - (\bar{y})^2}}$$

**Spearman's Rank correlation coefficient :. (R)**

$$R = 1 - \left[\frac{6\Sigma di^2}{n^3 - n}\right]$$

where $di = (R_1 - R_2) = [\text{difference ben}$
Ranks of $i^{th}$ item$]$

$$-1 \le R \le 1$$

Find Rank correlation coeff.

① Industry : 

| | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Rank (profit) : | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
| Rank (capital): | 13 | 16 | 14 | 15 | 10 | 12 | 4 | 11 | 5 | 9 | 8 | 3 | 1 | 6 | 7 | 2 |

(For repeated values)

Ans $R = -0.8176$

② 

| X | 85 | 74 | 85 | 50 | 65 | 78 | 74 | 60 | 74 | 90 |
|---|----|----|----|----|----|----|----|----|----|----|
| Y | 78 | 91 | 78 | 58 | 60 | 72 | 80 | 55 | 68 | 70 |

$R = 0.45$

Growth of employment in lakhs in india ben 1988 - 1995

③ 

| year | 1988 | 1989 | 1990 | 1991 | 1992 | 1993 | 1994 | 1995 |
|------|------|------|------|------|------|------|------|------|
| Public sector | 98 | 101 | 104 | 10.7 | 113 | 120 | 125 | 128 |
| private sector | 65 | 65 | 67 | 68 | 68 | 69 | 68 | 68 |

Note: For repeated items $Rank = \frac{i^{th} + j^{th}}{2}$ for Two repeated items

If $m_1, m_2, \cdots$ are repeated nos of items, Then spearman

Rank correlation coeff $= R = 1 - \left\{\frac{6\left[\Sigma di^2 + \frac{1}{12}(m_1^3 - m_1) + \frac{1}{12}(m_2^3 - m_2) + \cdots\right]}{n^3 - n}\right\}$

Del-21
MCQ: ④ Rank correlation coefficient of the following data is $[R = 1]$ Ans

| x | 23 | 25 | 27 | 29 | 31 | 33 |
|---|----|----|----|----|----|----|
| y | 43 | 45 | 47 | 49 | 51 | 53 |

Ans:- $R = 1 - \frac{6\Sigma(R_1 - R_2)^2}{N^3 - N} = 1 - 0$

For Not repeated ⑤

| x | 35 | 38 | 43 | 30 | 54 | 68 | 70 | 92 | 44 | 56 |
|---|----|----|----|----|----|----|----|----|----|----|
| Y | 51 | 37 | 48 | 62 | 93 | 73 | 56 | 72 | 70 | 92 |

Find Rank correlation coeff
$R = 0.59$
$N = 10$

For not repeated ⑥

| X | 105 | 110 | 112 | 108 | 111 | 116 | 120 | 104 | 115 | 125 |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Y | 39 | 41 | 45 | 38 | 48 | 58 | 60 | 35 | 54 | 69 |

$R = 0.9636$

For repeated ⑦

| x | 98 | 101 | 104 | 107 | 113 | 120 | 125 | 128 |
|---|----|-----|-----|-----|-----|-----|-----|-----|
| y | 65 | 65 | 67 | 68 | 68 | 69 | 68 | 68 |

$R = 0.98$

⑧ If $r = 0.4$, $cov(x,y) = 1.6$, $\sigma_y^2 = 25$, find $\sigma_x$. (Ans $\Rightarrow 0.8$)

Scanned with CamScanner

# REGRESSION

Regression line y on x is given by

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

Uday Kashid.

But $r \frac{\sigma_y}{\sigma_x}$ is called slope g line & called as

Regression coeff. y on x & denoted by $b_{yx}$.

$$\therefore \quad b_{yx} = r \frac{\sigma_y}{\sigma_x}$$

Hence $\quad y - \bar{y} = b_{yx} (x - \bar{x})$

similarly Regression line $\bar{x}$ on y .is

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

$$x - \bar{x} = b_{xy} (y - \bar{y}) \quad \text{where } b_{xy} = r \frac{\sigma_x}{\sigma_y}$$

① $\quad b_{yx} \cdot b_{xy} = r \frac{\sigma_y}{\sigma_x} \cdot r \frac{\sigma_x}{\sigma_y}$

$$b_{yx} \cdot b_{xy} = r^2$$

$$\therefore \quad r = \sqrt{b_{yx} \cdot b_{xy}}$$

⑪ $b_{yx} = r \frac{\sigma_y}{\sigma_x} = \frac{cov(x,y)}{\sigma_x \cdot \sigma_y} \frac{\sigma_y}{\sigma_x} = \frac{cov(x,y)}{\sigma_x^2}$

$$\Rightarrow b_{yx} = \frac{\Sigma(x-\bar{x})(y-\bar{y})}{\cancel{N} \, \Sigma(x-\bar{x})^2} = \frac{\Sigma(x-\bar{x})(y-\bar{y})}{\Sigma(x-\bar{x})^2}$$

∴ similarly

$$b_{xy} = r \frac{\sigma_x}{\sigma_y} = \frac{\Sigma(x-\bar{x})(y-\bar{y})}{\Sigma(y-\bar{y})^2}$$

Q⑬ show that Arithmetic mean of coeff. g regression is greater than or equal to coeff. g correlation.

→ for $\sigma_x$ & $\sigma_y$ we can write

$$(\sigma_x - \sigma_y)^2 \geq 0$$

$$\sigma_x^2 + \sigma_y^2 - 2\sigma_x \sigma_y \geq 0$$

$$\therefore \quad \sigma_x^2 + \sigma_y^2 \geq 2 \sigma_x \sigma_y$$

$$\frac{\sigma_x^2 + \sigma_y^2}{2 \sigma_x \sigma_y} \geq 1$$

$$\frac{1}{2} \left[ \frac{\sigma_x}{\sigma_y} + \frac{\sigma_y}{\sigma_x} \right] \geq 1 \qquad \text{multiply by } r$$

$$\frac{1}{2} \left[ r \frac{\sigma_x}{\sigma_y} + r \frac{\sigma_y}{\sigma_x} \right] \geq r$$

$$\frac{1}{2} \left[ b_{xy} + b_{yx} \right] \geq r \qquad \rightarrow \text{proved.}$$

Q 1) A panel of two judges Mrs. Madhuri Dixit-Nene & Dipika
Mrs. KaJOL graded dramatic performances by independently
Anushka Rashmika Mandana
awarding marks as follows.

| Performance NO | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| Marks By Madhuri Dipika | 36 | 32 | 34 | 31 | 32 | 32 | 34 |
| Marks by Anushka | 35 | 33 | 31 | 30 | 34 | 32 | 36 |

The eighth performance however which Anushka could not
Judge
attend, got 38 marks by Judge Dipika. If Judge Anushka
had also been present, how many marks would she be
expected to have awarded to the eight performance.

Method I:
$\bar{X} = \frac{\Sigma x}{N} = 33$, $\bar{Y} = \frac{\Sigma Y}{N} = 33$, $b_{yx} = \frac{\Sigma(x-\bar{x})(y-\bar{y})}{\Sigma(x-\bar{x})^2} = 0.72$

$x - \bar{y} = b_{yx}(x-\bar{x})$ ⟹ at $x = 38$, ⟹ $y = 36.6 \cong 37$

Method II:
$y = a + bx$ → equ of required line.
$\Sigma y = aN + b\Sigma x$ and $\Sigma xy = a\Sigma x + b\Sigma x^2$

Q 2) Find the coeff. of regression lines and hence equ of regression
lines for the following data,

| X | 78 | 36 | 98 | 25 | 75 | 82 | 90 | 62 | 65 | 39 |
|---|---|---|---|---|---|---|---|---|---|---|
| Y | 84 | 51 | 91 | 60 | 68 | 62 | 86 | 58 | 53 | 47 |

Estimate value of Y when $x = 50$ and value of X
when $y = 90$.

Ex 3) obtain the equations of line of regression for the following data.

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

→ $n=8$ $\Sigma x = 544$ $\Sigma y = 552$ $\Sigma xy = 37560$ $\Sigma x^2 = 37028$, $\Sigma y^2 = 38132$

① Regression line y on x is $y = ax + b$ ⟹ $\Sigma y = a\Sigma x + bn$, $\Sigma 1 = n$
$\Sigma xy = a\Sigma x^2 + b\Sigma x$ → By Solving $a = 0.66$ & $b = 23.67$
$Y = (0.66)x + 23.67$

② Regression line x on y ⟹ $x = ay + b$ ⟹ $\Sigma x = a\Sigma y + bn$
$\Sigma xy = a\Sigma y^2 + b\Sigma y$ → By Solving $a = 0.55$ & $b = 30.36$.
$x = (0.55)y + 30.36$

Ex 4) Fit the straight line of the form $y = ax + b$ to the following data,

| X | 1 | 3 | 5 | 7 | 8 | 10 |
|---|---|---|---|---|---|---|
| Y | 8 | 12 | 15 | 17 | 18 | 20 |

$\delta n$ $n = 6$ $\Sigma y = a\Sigma x + b$ & $\Sigma xy = a\Sigma x + b\Sigma x$
$\Sigma x = 34$ $\Sigma y = 90$ $\Sigma xy = 582$, $\Sigma x^2 = 248$
$a = 1.300$, $b = 7.63$

## Regression lines

**Ex ① Find ① the lines of regression ① coeff of correlation for**

(A)

| X | 5 | 10 | 12 | 13 | 16 | 17 | 20 | 25 |
|---|---|----|----|----|----|----|----|----|
| Y | 6 | 19 | 22 | 24 | 27 | 29 | 33 | 37 |

Ans:-
$y = 0.8x + 13.23$
$x = \cdot$     $y =$
$r = \cdot$

**Ex ② The heights in cms of fathers (x) and of the eldest sons (y) are given below. Then find the lines of regression. Also estimate the height of the eldest son If the height of the father is 172 cms. and the height of father if the height of son is 173 cm. Also find coeff of correlation (r).**

| X | 165 | 160 | 170 | 163 | 173 | 158 | 178 | 168 | 173 | 170 | 175 | 180 |
|---|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Y | 173 | 168 | 173 | 165 | 175 | 168 | 173 | 165 | 180 | 170 | 173 | 178 |

$\{$ Ans:- ① $y = (1.016)x - 5.123$  ② $x = (0.476)y + 98.98$  ③ $169.97$
④ $173.45$  ⑤ $r = 0.696$ $\}$

---

**Fitting of parabola [second degree curve]**

Note ⇒
We know second degree curve i.e. (fitting of parabola) equation
is  $y = ax^2 + bx + c$  —①

$\Sigma y = a\Sigma x^2 + b\Sigma x + c\Sigma 1$     But $\Sigma 1 = N$.

$\Sigma y = a\Sigma x^2 + b\Sigma x + cN$  —②     $\Big|$ $\Sigma x^2 y = a\Sigma x^4 + b\Sigma x^3 + c\Sigma x^2$ —④

$\Sigma xy = a\Sigma x^3 + b\Sigma x^2 + c\Sigma x$ —③ $\Big|$

By solving eq ②, ③ & ④ the obtain values of a, b, & c.

**Ex ① By the method of least square method find the best values of a, b, c in the second degree curve is. $y = ax^2 + bx + c$ to fit the following data**

| X | -2 | -1 | 0 | 1 | 2 |
|---|----|----|----|----|----|
| Y | -3.150 | -1.390 | 0.620 | 2.880 | 5.378 |

Ans:—
$Y = (0.1233)x^2 + (2.1326)x$
$+ (0.621)$

**Ex ② Fit a second degree curve (parabola) to the following data and estimate y when $x = 80$.**

| X | 10 | 20 | 30 | 40 | 50 | 60 | 70 |
|---|----|----|----|----|----|----|----|
| Y | 20 | 60 | 70 | 80 | 90 | 100 | 100 |

at $x = 80°, \Rightarrow y = 94.286$

[Ans:- Put $u = (x - 40)/10$
$v = y/10$
$v = \cdot$   $Y = (0.1233)x$
$Y = (0.2381)x^2 + (1.2143)x + (8.381)$]

**Ex ③ Fit the second degree curve (parabola) to the following data**

| year (X) | 1965 | 1966 | 1967 | 1968 | 1969 | 1970 | 1971 | 1972 |
|----------|------|------|------|------|------|------|------|------|
| profit in crores Rs (Y) | 125 | 140 | 165 | 195 | 200 | 215 | 220 | 230 |

Ans:- put $x = (x - 1968.5) \times 2$, $Y = Y$, $Y = (-0.40)x^2 + (7.68)x + 194.68$