# Module III

## Correlation And Regression

### Chapter 1:
#### Curve fitting

let $x$ be an independent variable & $y$ be a variable depending on $x$. Here, we say that $y$ is a function of $x$ & write it as $\boxed{y = f(x)}$.

If $f(x)$ is a known fun, then for any allowble values $x_1, x_2 \dots x_n$ of $x$, we can find the corresponding values $y_1, y_2 \dots y_n$ of $y$ & there by determine the pairs $(x_1, y_1) (x_2, y_2) \dots (x_n, y_n)$ which constitute a bivariate data. These pairs of values of $x$ & $y$ give us $n$ points on the curve $y = f(x)$.

Suppose considel the converse prbm, Suppose we are given $n$ values $x_1, x_2 \dots x_n$ of an independent variable $x$ & the corresponding values $y_1, y_2 \dots y_n$ of a variable $y$ depending on $x$. Then the pairs $(x_1, y_1) (x_2, y_2) \dots (x_n, y_n)$ gives us $n$ points in $xy$ plane.

Generally it is not possible to find the

actual curve $y=f(x)$ that passes through these points.

**Q** Hence we try to find a curve that serves as best approximation to the curve $y=f(x)$. Such a curve is referred as **curve of best fit**. The process of determining a curve of best fit → **curve fitting**. The method for curve fitting is called method of least squares.

→ **Method of least squares:-**

This method is for finding the unknown coefficients in a curve that serves as best approximation to the curve $y=f(x)$.

By A.M Legendre & C.F Gauss, the Principle of least squares says that the sum of squares of the error for the observed values i.e the corresponding estimated values should be the least.

b/s the observed values i.e the corresponding estimated values should be the least.

Suppose to fit a $k^{th}$ deg curve, given by $y=a_0+a_1x+a_2x^2\cdots a_kx^k$ to $\cdots$ (1) the given $n$ observations $(x_1,y_1)(x_2,y_2)\sim(x_n,y_n)$. The curve has $k+1$ unknown const. & hence if $n\geq k+1$.

we get $k+1$ eq on subset the values of $(x_i,y_i)$ in (1). This gives unique soln to the values. However if $n>k+1$: no unique soln of least fit. (D)

now let, $y_e = a_0+a_1x+a_2x^2+\cdots a_kx^k$ be the estimated values of $y$ when $x$ takes the value $x_i$. But the corresponding observed value of $y$ is $y_i$.

hence if $e_i$ is the residual as error for this point.

$$e_i = y_i - y_e = y_i - a_0 - a_1x_i - a_2x_i^2 \sim a_kx_i^k$$

To make the sum of D min, we have to minimise,

$$S = \sum_{i=1}^{n} e_i^2 = \sum_{i=1}^{n} (y_i-a_0-a_1x_i-a_2x_i^2\cdots a_kx_i^k)^2 \quad (2)$$

By differential calculus, S will have its min-value when,

$$\frac{\partial S}{\partial a_0} = 0, \frac{\partial S}{\partial a_1} = 0 \cdots \frac{\partial S}{\partial a_k} = 0 \cdots$$

which gives $k+1$ eq → Normal eqi
Solving these eq~ we get the best values of $a_0,a_1,\ldots a_k$

substit..an (1). we get the curve of best fit.

① ⇒ **Fitting of a straight line :-**

$$y = a + bx$$

$$\boxed{\begin{aligned} \Sigma y_i &= na + b\Sigma x_i \\ \Sigma x_i y_i &= a\Sigma x_i + b\Sigma x_i^2 \end{aligned}} \Bigg\} \text{ Normal eq.}$$

Q) Fit a straight line

$$x = 1 \quad 2 \quad 3 \quad 4 \quad 5$$
$$y = 14 \quad 13 \quad 9 \quad 5 \quad 2$$

Estimate the value of $y$ when $x = 3.5$

A) $N = 5$

| $x_i$ | $y_i$ | $x_i^2$ | $x_i y_i$ |
|---|---|---|---|
| 1 | 14 | 1 | 14 |
| 2 | 13 | 4 | 26 |
| 3 | 9 | 9 | 27 |
| 4 | 5 | 16 | 20 |
| 5 | 2 | 25 | 10 |
| $\Sigma x_i = 15$ | $\Sigma y_i = 43$ | $\Sigma x_i^2 = 55$ | $\Sigma x_i y_i = 97$ |

$\Sigma y_i = na + b\Sigma x_i \longrightarrow 43 = 5a + 15b$ —①

$\Sigma x_i y_i = a\Sigma x_i + b\Sigma x_i^2 \longrightarrow 97 = 15a + 55b$ —②

①×3     $15a + 45b = 129$
②         $15a + 55b = 97$

$$-10b = 32$$
$$b = -3.2$$

$15a + 55 \times -3.2 = 97$
$15a - 176 = 97$
$15a = 19 + 176$
$15a = 273$
$$a = 18.2$$

$y = 18.2 - 3.2 \times 3.5$
$= 18.2 - 11.2$
$$y = 7$$

$y = a + bx$
$x = 3.5$

* Suppose to have a straight line that serves as best approximation to the actual curve $y = f(x)$ passing through given points $(x_1,y_1), (x_2,y_2) \dots (x_n,y_n)$. This line will be referred as a **line of best fit.**

$$y = a + bx \quad \text{—①}$$

$a, b \rightarrow$ parameters to be determined.
Let '$y_i$' be the value of $y$ corresponding to the value $x_i$ of $x$ as determined by eq —①. The value $y_e \rightarrow$ **Estimated value of** $y$ when $x = x_i$, the observed value of $y$.

Then the diff "$y_i - y_e \rightarrow$ **Residual/error.**
By the principle of least squares, we have

$$S = \Sigma (y_i - y_e)^2 \quad \text{—②}$$

Determine $a$ & $b$, so that $S$ is min.
Two necessary conditions for this are,
$$\frac{\partial s}{\partial a} = 0 \,, \quad \frac{\partial s}{\partial b} = 0.$$

using —② these conditions yield the following eq.

$\Sigma(y_i - a - bx_i) = 0$ or $\Sigma y_i = na + b\Sigma x_i$ —③

$\Sigma(y_i - a - bx_i)x_i = 0$ or $\Sigma x_i y_i = a\Sigma x_i + b\Sigma x_i^2$ —④

③ & ④ $\rightarrow$ **Normal eq.** has determining $a$ & $b$.

Putting the values of $a$ & $b$ & $c$
determined → ①, we get the eq of line
of best fit for the given data.

② **Fitting of parabola :-**

$$\boxed{y = a + bx + cx^2}$$

Suppose curve with to true a parabola as
the curve of best fit for a data
consisting of 'n' given pairs $(x_i, y_i), i = 1, 2, \ldots n$

Q) of best fit in the form $y = a + bx + cx^2$ —①
where $a, b, c$ = constants.

Let $\bar{y}$ de the value of 'y' corresponding
to the value '$x_i$' of $x$ determined by
eq—①, then sum of error squares of error
b/w observed value of 'y' & estimated
value of 'y' is given by :
$y_e$ → estimated value.

using eq—①
$$S = \sum_{i=1}^{n} \left(y_i - (\sum_{i=1}^{n} y_i - y_e)\right)^2$$
$$S = \sum_{i=1}^{n} \left(y_i - (a+bx+cx^2)\right)^2 \quad —②$$
so that 'S' is least

3 necessary conditions for this are,
$$\frac{\partial S}{\partial a} = 0 \quad \frac{\partial S}{\partial b} = 0 \quad \frac{\partial S}{\partial c} = 0$$

using eq—② This conditions yield the following
normal eq—

$$\boxed{\sum y_i = na + b\sum x_i + c\sum x_i^2} \quad —③$$
$$\boxed{\sum x_i y_i = a\sum x_i + b\sum x_i^2 + c\sum x_i^3} \quad —④$$
$$\boxed{\sum x_i^2 y_i = a\sum x_i^2 + b\sum x_i^3 + c\sum x_i^4} \quad —⑤$$

Solve eq —③, —④, —⑤ has determining $a, b, c$.
Putting the values $a, b, c$, so determining
in eq—① we get the eq of parabola of
best fit for the given data.

Q) fit a parabola.

$$y = a + bx + cx^2 \quad —①$$ to the following data.

$x$ : 1   2   3   4   5   6   7
$y$ : 2.3   5.2   9.7   16.5   29.4   35.5   54.4.

A) $n = 7$

| $x_i$ | $y_i$ | $x_i^2$ | $x_i^3$ | $x_i^4$ | $x_iy_i$ | $x_i^2 y_i$ |
|---|---|---|---|---|---|---|
| 1 | 2.3 | 1 | 1 | 1 | 2.3 | 2.3 |
| 2 | 5.2 | 4 | 8 | 16 | 10.4 | 20.8 |
| 3 | 9.7 | 9 | 27 | 81 | 29.1 | 87.3 |
| 4 | 16.5 | 16 | 64 | 256 | 66.0 | 264.0 |
| 5 | 29.4 | 25 | 125 | 625 | 147 | 735 |
| 6 | 35.5 | 36 | 216 | 1296 | 213 | 1278 |
| 7 | 54.4 | 49 | 343 | 2401 | 380.8 | 2665.6 |
| $\sum x_i = 28$ | $\sum y_i = 153$ | $\sum x_i^2 = 140$ | $\sum x_i^3 = 784$ | $\sum x_i^4 = 4676$ | $\sum x_iy_i = 848.6$ | $\sum x_i^2 y_i = 5053$ |

$$153 = 7a + 28b + 140c \quad —①$$
$$848.6 = 28a + 140b + 784c \quad —②$$
$$5053 = 140a + 784b + 4676c \quad —③$$

$\frac{153}{7}$   $\frac{848.6}{28}$   $\frac{5053}{140}$

$$28a + 140b + 784c = 848.6$$
$$28a + 112b + 560c = 612$$
$$28b + 224c = 236.6 \quad —④$$

## Left page

$$eq-② \times 5 \rightarrow 140a + 700b + 3920c = 4243$$

$$140a + 784b + 4676c = 5053$$
$$\underline{140a \pm 700b \pm 3920c = 4243}$$
$$84b + 756c = 810 \quad —⑤$$

$$eq-④ \times 3 \rightarrow 84b + 672c = 709.8$$

$$84b + 756c = 810$$
$$\underline{84b \pm 672c = 709.8}$$
$$84c = 100.2$$
$$c = 1.1928 \rightarrow 1.193$$

$$84b + 672 \times 1.193 = 709.8$$
$$84b + 801.696 = 709.8$$
$$84b = 709.8 - 801.696$$
$$84b = -91.896$$
$$b = -1.094$$

$$7a + 28b + 140c = 153$$
$$7a + 28 \times -1.094 + 140 \times 1.193 = 153$$
$$7a - 30.632 + 167.02 = 153$$
$$7a + 136.388 = 153$$
$$7a = 153 - 136.388$$
$$7a = 16.612$$
$$a = 2.37$$

$$a = 2.37$$
$$b = -1.094$$
$$c = 1.193$$

## Right page

③ **Fitting a parabola** :-    $y = ab^x$.

$$\boxed{\begin{array}{l} \log(ab) = \log a + \log b \\ \log(a^b) = b \log a. \end{array}}$$

Fitting of a curve of the form $y = ab^x$
Suppose we wish to fine a curve
whose eq $\rightarrow y = ab^x$ —Ⓐ

As the curve is of best fit for a data a
consisting of pairs $(x_i, y_i)i = 1,2,....n$ taking
log on both sides of eq —Ⓐ we get,

$$y = ab^x$$
$$\log y = \log a + x \log b$$

$$\log y = y$$
$$\log a = A$$
$$\log b = B$$

$$y = A + Bx \quad —②$$

Is a linear eq, is the normal eq thet out
yield A & B —

$$\boxed{\begin{array}{l} \Sigma y_i = nA + B\Sigma x_i \quad —③ \\ \Sigma x_i y_i = A\Sigma x_i + B\Sigma x_i^2 \quad —④ \end{array}}$$

where $y_i = \log y_i$
Solving this eq we are obtain A & B—

$$a = antilog\ A$$
$$b = antilog\ B$$

Substituting these values of a & b in eq—Ⓐ
we obtain the curve best fit for
the given data.

1) Fit a curve: $y = ab^x$

x : 1   2   3   4   5   6   7   8
y: 1.0   1.2   1.8   2.5   3.6   4.7   6.6   9.1

A) $n = 8$

($y = ab^x$)

| $x_i$ | $y_i$ | $y_i' = \log_e y_i$ | $x_i^2$ | $x_i y_i'$ |
|---|---|---|---|---|
| 1.0 | 1.0 | 0.000 | 1 | 0 |
| 2 | 1.2 | 0.1823 | 4 | 0.3646 |
| 3 | 1.8 | 0.5878 | 9 | 1.7634 |
| 4 | 2.5 | 0.9163 | 16 | 3.6652 |
| 5 | 3.6 | 1.2809 | 25 | 6.40045 |
| 6 | 4.7 | 1.5475 | 36 | 9.285 |
| 7 | 6.6 | 1.8870 | 49 | 13.209 |
| 8 | 9.1 | 2.2082 | 64 | 17.6656 |

$\Sigma x_i = 36$   $\Sigma y_i = 30.5$   $\Sigma y_i' = 8.61$   $\Sigma x_i^2 = 204$   $\Sigma x_i y_i' = 52.3573$

$\Sigma y_i' = nA + B\Sigma x_i \rightarrow 8.61 = 8A + 36B$   ×36
$\Sigma x_i y_i' = A\Sigma x_i + B\Sigma x_i^2 \rightarrow 52.3573 = 36A + 204B$   ×8

$309.96 = 288A + 1296B$
$418.8584 = 288A + 1632B$

$288A + 1632B = 418.8584$
$288A + 1296B = 309.96$

$336B = 108.8984$

$B = 0.3241$

$8A + 36 \times 0.3241 = 8.61$
$8A + 11.6676 = 8.61$
$8A = 8.61 - 11.6676$
$8A = -3.0576$
$A = -0.3822$

$a = \text{antilog } A \rightarrow$ step + in → no
$= \text{antilog}(-0.3822) = 0.6684$

$b = \text{antilog } B$
$= \text{antilog}(0.3241) = 1.3828$

$y = ab^x$
$= (0.6684)(1.3828)^x$

---

H/W
1) Fit a parabola —

x : 1   2   3   4   5   6   7   8   9   10   11
y: 2   6   7   8   10   11

$x = 4.5$   $y = 6$.

A) $n = 9$.

| $x_i$ | $y_i$ | $x_i^2$ | $x_i^3$ | $x_i^4$ | $x_i y_i$ | $x_i^2 y_i$ |
|---|---|---|---|---|---|---|
| 1 | 2 | 1 | 1 | 1 | 2 | 2 |
| 2 | 6 | 4 | 8 | 16 | 12 | 24 |
| 3 | 7 | 9 | 27 | 81 | 21 | 63 |
| 4 | 8 | 16 | 64 | 256 | 32 | 128 |
| 5 | 10 | 25 | 125 | 625 | 50 | 250 |
| 6 | 11 | 36 | 216 | 1296 | 66 | 396 |
| 7 | 11 | 49 | 343 | 2401 | 77 | 539 |
| 8 | 10 | 64 | 512 | 4096 | 80 | 640 |
| 9 | 9 | 81 | 729 | 6561 | 81 | 729 |

| $\Sigma x_i =$ | $\Sigma y_i =$ | $\Sigma x_i^2 =$ | $\Sigma x_i^3 =$ | $\Sigma x_i^4 =$ | $\Sigma x_i y_i =$ | $\Sigma x_i^2 y_i =$ |
|---|---|---|---|---|---|---|
| 45 | 74 | 285 | 2025 | 15333 | 421 | 2771 |

$74 = 9a + 45b + 285c$ —— ①

$421 = 45a + 285b + 2025c$ —— ②

$2771 = 285a + 2025b + 15333c$ —— ③

① × 5 → $370 = 45a + 225b + 1425c$

$$45a + 285b + 2025c = 421$$
$$\underline{45a + 225b + 1425c = 370}$$
$$60b + 600c = 51 \quad \text{——④}$$

② × 45
③ × 45

$128259 + 91125b + 689985c = 1246695$

$\underline{12825a + 81225b + 577125c = 1119785}$

$$9900b + 112860c = 4710 \quad \text{——⑤}$$

$60b + 600c = 51 \qquad \times 9900$

$9900b + 112860c = 4710 \qquad \times 60$

$594000b + 6771600c = 282600$

$\underline{594000b + 5940000c = 504900}$

$$831600c = -222300$$
$$c = -0.2673$$

$60b + 600 \times -0.2673 = 51$

$60b - 160.38 = 51$

$60b = 51 + 160.38$

$60b = 211.38$

$b = 3.523$

---

$454a + 285 \cdot 13.5 \; 23 + 2025 \times 0.2673 = 421$

$45a + 1004.055 - 541.2825 = 421$

$45a + 462.7725 = 421$

$45a = 421 - 462.7725$

$45a = -41.7725$

$a = 0.9282$

$$a = 0.9282 \qquad b = 3.523 \qquad c = -0.2673$$

---

④ **Fitting of a curve $y = ax^b$ :-**

Suppose we wish to true a curve whose eq is of the form $y = ax^b$ —①

As the curve is best fit for a data consisting of the pairs $(x_i, y_i)^{\circ} = 1, 2, \dots n$

Taking log on both sides of eq ①
we get.

$\log y = \log a + b \log x$

$\log y = A$
$\log a = A$
$\log x = X$

$u = A + bx$ —② & a linear eq

$\Sigma u_i = nA + b\Sigma x_i$ —③

$\Sigma x_i u_i = A\Sigma x_i + b\Sigma x_i^2$ —④

⑤ **Fitting of a curve $y = ae^{bx}$ :-**

$y = ae^{bx}$ —①

If $y = ae^{bx}$ is the curve of best fit for a data.

taking log

$\log y = \log a + bx \log e$ —②

$u = \log y$
$A = \log a$
$B = b \log e$

$u = A + Bx$ —③

$$\Sigma u_i^0 = nA + B\Sigma x_i^0 \quad —(4)$$

$$\Sigma x_i u_i = A\Sigma x_i^0 + B\Sigma x_i^2 \quad —(5)$$

Solving eq —(4) & —(5) we get A & B

$$\boxed{\begin{array}{l} a = antilog\ A \\ b = antilog\ \dfrac{B}{\log e} \end{array}} \quad \dfrac{y_0}{\log\left(\dfrac{B}{e}\right)\log e} \dfrac{y_0}{x_i^0}$$

Substituting a, b in eq —① we get the curve of best fit for the given pairs of observations.

**⑥ Fitting of a curve** $y = ax^2 + \dfrac{b}{x}$ :-

To fit this curve for a set of observations $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$. The normal eq —

$$\boxed{\begin{array}{l} \Sigma y = a\Sigma x^2 + b\Sigma \dfrac{1}{x} \\ \Sigma xy = a\Sigma x^3 + nb \end{array}}$$

Solve 2 eq for a & b; Substituting a, b in eq, we get the required curve of best fit.

**1)** Fit a straight line (by least □'s) to following data —

| x : | 1 | 2 | 3 | 4 | 5 |
|-----|---|---|---|---|---|
| y : | 35 | 68 | 100 | 138 | 170 |

$$y = a + bx.$$

$$\Sigma y = na + b\Sigma x_i^0 \quad —①$$
$$\Sigma x_i y_i = a\Sigma x_i^0 + b\Sigma x_i^2 \quad —②$$

| $x_i^0$ | $y_i^0$ | $x_i^2$ | $x_i y_i^0$ |
|-----|-----|-----|-----|
| 1 | 35 | 1 | 35 |
| 2 | 68 | 4 | 136 |
| 3 | 100 | 9 | 300 |
| 4 | 138 | 16 | 552 |
| 5 | 170 | 25 | 850 |
| $\Sigma x_i^0=15$ | $\Sigma y_i=511$ | $\Sigma x_i^2=55$ | $\Sigma x_i y_i = 1873$ |

$$511 = 5a + 15b \quad ×3$$
$$1873 = 15a + 55b$$

$$15a + 45b = 1533$$
$$\underline{15a + 55b = 1873}$$
$$\dfrac{15a + 45 \times 34 = 1533}{}$$
$$15a + 1530 = 1533$$
$$15a = 1533 - 1530$$
$$15a = 3$$
$$\underline{\underline{a = 0.2}}$$

$$-10b = 340$$
$$\underline{\underline{b = 34}}$$

Line of best fit —
$$y = a + bx$$
$$\underline{\underline{y = 0.2 + 34\,x}}$$

---

**3)**

| x : | 1 | 2 | 3 | 4 | 5 | 6 |
|-----|---|---|---|---|---|---|
| y : | 2.4 | 3 | 3.6 | 4 | 6 | 8 |

**48** 

| x : | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|-----|---|---|---|---|---|---|---|
| y : | 80 | 90 | 92 | 83 | 94 | 99 | 92 |

| x : | 0 | 1 | 2 | 3 | 4 |
|-----|---|---|---|---|---|
| y : | 1 | 1.8 | 3.3 | 4.5 | 6.3 |

$$n = 5$$

**2.A)**

| xi | yi | xi² | xiyi |
|---|---|---|---|
| 0 | 1 | 0 | 0 |
| 1 | 1.8 | 1 | 1.8 |
| 2 | 3.3 | 4 | 6.6 |
| 3 | 4.5 | 9 | 13.5 |
| 4 | 6.3 | 16 | 25.2 |
| $\Sigma xi=10$ | $\Sigma yi=16.9$ | $\Sigma xi^2=30$ | $\Sigma xiyi=47.1$ |

n = 5

$16.9 = 5a + 10b$

$47.1 = 10a + 30b$ ×3

$15a + 30b = 50.7$

$10a + 30b = 47.1$

$5a = 3.6$    $a = \dfrac{3.6}{5} = 0.72$

$5 \times 0.72 + 10b = 16.9$

$3.6 + 10b = 16.9$

$10b = 16.9 - 3.6$

$10b = 13.3$

$b = 1.33$

$$y = 0.72 + 1.33x$$

---

**3.A)**

| xi | yi | xi² | xi yi |
|---|---|---|---|
| 1 | 2.4 | 1 | 2.4 |
| 2 | 3 | 4 | 6 |
| 3 | 3.6 | 9 | 10.8 |
| 4 | 4 | 16 | 16 |
| 5 | 6 | 25 | 30 |
| 6 | 8 | 36 | 48 |
| $\Sigma yi=27$ | | $\Sigma xi^2=91$ | $\Sigma xiyi=113.2$ |

$27 = 6a + 21b$ ×21

$113.2 = 21a + 91b$ ×6    n=6.

$126a + 441b = 567$

$126a + 546b = 679.2$

$105b = 112.2$

$b = 1.068$

$6a + 21 \times 1.068 = 27$

$6a + 22.428 = 27$

$6a = 27 - 22.428$

$6a = 4.572$

$a = 0.762$

$y = a + bx$

$$y = 0.762 + 1.068x$$

## 4·A)

| $x_i$ | $y_i$ | $x_i^2$ | $x_i y_i$ |
|---|---|---|---|
| 1 | 80 | 1 | 80 |
| 2 | 90 | 4 | 180 |
| 3 | 92 | 9 | 276 |
| 4 | 83 | 16 | 332 |
| 5 | 94 | 25 | 470 |
| 6 | 99 | 36 | 594 |
| 7 | 92 | 49 | 644 |
| | $\sum y_i = 630$ | $\sum x_i^2 = 140$ | $\sum x_i y_i = 2576$ |

$\sum x_i = 28$

$630 = 7a + 28b \qquad \times 28$

$2576 = 28a + 140b \qquad \times 7 \qquad n=7$

$196a + 980b = 18032$

$196a + 784b = 17640$

$\overline{196b = 392}$

$b = 2$

$7a + 28 \times 2 = 630$

$7a + 56 = 630$

$7a = 630 - 56$

$7a = 574$

$\underline{\underline{a = 82}}$

$y = a + bx$

$\underline{\underline{y = 82 + 2x}}$

---

1) Derive the least □ eq's for fitting a curve of the type $y = ax + \dfrac{b}{x}$ to a set of $n$ points $(x_i,y_i)_i = 1,2,3\ldots n$.

A) The error of estimate $(e_i)$ for $i^{th}$ point $(x_i,y_i)$ is given by,

$e_i = (y_i - ax_i - \dfrac{b}{x_i})$, $\qquad \hat{y} = ax + \dfrac{b}{x}$

According to Principle of least □, we have to determine the values of a & b, so that the sum of the □'s of errors

$S = \sum e_i^2 = \sum_{i=1}^{n}\left(y_i - ax_i - \dfrac{b}{x_i}\right)^2$ is min.

consequently the normal eq's —

$\dfrac{\partial S}{\partial b} = 0 => -2\sum_{i=1}^{n}\dfrac{1}{x_i}\left(y_i - ax_i - \dfrac{b}{x_i}\right)$ —①

$\dfrac{\partial S}{\partial a} = 0 => -2\sum x_i\left(y_i - ax_i - \dfrac{b}{x_i}\right)$

$\sum \dfrac{y_i}{x_i} = na + \dfrac{b}{\sum x_i^2} => na + \dfrac{b}{\sum x_i^2}$

$\sum x_i y_i = a\sum x_i^2 + nb$

H/w)

2) $y = a e^{bx}$

| x : | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| y : | 1.6 | 4.5 | 13.8 | 40.2 | 125.0 | 300 |

3) $y = ax^2 + \dfrac{b}{x}$

| x : | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| y : | 1.51 | 0.99 | 3.88 | 7.66 |

4) $y = \dfrac{a}{x} + bx$

| x : | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|
| y : | 5.40 | 6.30 | 8.20 | 10.30 | 12.60 | 14.40 | 17.30 | 19.58 |

**2.A)**

| $x_i$ | $y_i$ | $x_i^2$ | $u_i = \log y_i$ | $x_i u_i$ |
|---|---|---|---|---|
| 1 | 1·6 | 1 | 0·4700 | 0·47 |
| 2 | 4·5 | 4 | 1·5040 | 3·008 |
| 3 | 13·8 | 9 | 2·6246 | 7·8798 |
| 4 | 40·2 | 16 | 3·6938 | 14·7752 |
| 5 | 125·0 | 25 | 4·8283 | 24·1415 |
| 6 | 300 | 36 | 5·7037 | 34·2222 |
| 21 | 448·62 | 91 | 18·8244 | 84·4907 |

$18.8244 = 6A + 21B$    × 21

$84.4907 = 21A + 91B$    × 6

$126A + 441B = 375.3124$
$126A + 441B = 506.9424$

$105B = 11.6318$

$B = 1.06316$

$a = \text{antilog}(-0.58366)$
$= 0.5577$

$b = \text{Blog } e$
$= 1.05$

$105B = 131.0631$
$B = 1.06316$

$6A + 21 \times 1.06316 = 18.8244$
$6A + 22.32636 = 18.8244$
$6A = -3.50196$
$A = -0.5 \, 8366$

$(\text{value} \to \text{log } 2.71) \to 35$

$y = 0.5578^{1.05x}$

---

**3.A)**

| $x$ | $y$ | $x^2$ | $x^3$ | $x^2y$ | $1/x$ |
|---|---|---|---|---|---|
| 1 | 1·51 | 1 | 1 | 1·51 | 1 |
| 2 | 0·99 | 4 | 8 | 1·98 | 0·5 |
| 3 | 3·88 | 9 | 27 | 11·64 | 0·33 |
| 4 | 7·66 | 16 | 64 | 30·64 | 0·25 |
| 10 | 14·04 | 30 | 100 | 45·77 | 2·08 |

$14·04 = 30a + 2·08b$   × 100
$45·77 = 100a + 4b$   × 30

---

**4.A)**

$30a + 2·08b = 1404$
$30600a + 120b = 1373·1$

$88b = 30·9$

$b = 0·35$

$100a + 4 \times 0·35 = 45·77$
$100a + 1·4 = 45·77$
$100a = 44·37$
$a = 0·4437$

$y = 0.4437 x^2 + \dfrac{0.35}{x}$

$y = \dfrac{a}{x} + bx$

$y_i = \dfrac{a}{x_i^2} + bx_i$

$S = \varepsilon r_i^2 = \left(y_i - \dfrac{a}{x_i^2} - bx_i\right)^2$

$\dfrac{\partial S}{\partial a} = 0 \Rightarrow -2\dfrac{1}{x_i^2}\left(y_i - \dfrac{a}{x_i^2} - bx_i\right) = 0$ —①

$\dfrac{\partial S}{\partial b} = 0 \Rightarrow -2 x_i\left(y_i - \dfrac{a}{x_i^2} - bx_i\right) = 0$ —②

$-① \Rightarrow \dfrac{y_i}{x_i^2} = a\dfrac{1}{x_i^4} + nb$

$-② \Rightarrow x_i y_i = na + b x_i^2$

| $x_i$ | $y_i$ | $x_i^2$ | $\dfrac{1}{x_i^2}$ | $x_i y_i$ | $\dfrac{y_i}{x_i^2}$ |
|---|---|---|---|---|---|
| 1 | 5·40 | 1 | 1 | 5·4 | 5·4 |
| 2 | 6·30 | 4 | 0·25 | 12·6 | 3·15 |
| 3 | 8·20 | 9 | 0·111 | 24·6 | 2·73 |
| 4 | 10·30 | 16 | 0·0625 | 41·2 | 2·575 |
| 5 | 12·60 | 25 | 0·04 | 63 | 2·52 |
| 6 | 14·90 | 36 | 0·0277 | 89·4 | 2·48 |

| 7 | 17·30 | 4.9 | 0·0204 | 19·41 | 2·4·7 |
|---|---|---|---|---|---|
| 8 | 19·50 | 64 | 0·015665 | 156 | 2·43 |
|   |   |   | 1·527325 | 5133 | 23·755 |

$\sum x_i y_i = na + b \sum c_i^2$

$\sum \dfrac{y_i}{x_i} = a \sum \dfrac{1}{x_i^2} + nb$

$513·3 = 8a + 20·6b$

$23·755 = 1·527a + 8b \qquad \times 1·527$

$12·2168 + 311·508b = 783·8091$

$\dfrac{12·2168 \,+\, 6a\,b \;=\; 190·04}{247·508b = 593·769}$ ( )

$b = 2·398$

$12·216a + 64 \times 2·398 = 190·04$

$12·216a + 153·472 = 190·04$

$12·216a = 36·568$

$a = \dfrac{2·99}{x} \dashrightarrow (ii)$

$\underline{y = 2·99 + 2·398x}$

---

## Q: Correlation & Regression

**Correlation :-** It is a statistical measure for binding out degree of association b/w 2/more variable by association mean the tendency of the 2 variables to move together.

*x* tend to accombined by the corresponding movements in the other wi *y*. Then *x* & *y* are said to be correlated. The movements may be in the same (dir)/ opposite (dir). $(x\uparrow, y\uparrow) \to$ same (dir), $x\downarrow, y\uparrow\to$ opp.)

(c) Said to be +ve/-ve according as the movements are in the same (in this movement are in the opposite (dir). If 'y' is unaffected by any change in 'x'. then x & y are said to be un correlated.

**LR conner's Definition :-**

If 2 more quantities vary in sympathy, so that movements in the 1 tend to be accombined by the corresponding movements in the other, then they are said to be correlated.

→ **Linear (c) :-**

(c) may be linear /non-linear. The variation in 'x' bears a constant ratio to the corresponding amount of variation in 'y'. then (c) b/s

If measure[s] the relation b/w 2 variables $x$ & $y$ is given the formula,

$$r_{xy} = \frac{\sum_{i=1}^{n}(x_i - \bar{x})(y_i - \bar{y})}{n \; \sigma_x \; \sigma_y} \quad \text{---} \; (1).$$

$x_i, y_i$ $(i = 1, 2, \dots, n) = 2$ sets of values of variables.

$\bar{x}, \bar{y}$ = means

$\sigma_x, \sigma_y$ = standard deviation of $x$ & $y$.

$$\bar{x} = \frac{1}{n}\sum_{i=1}^{n}x_i \qquad \bar{y} = \frac{1}{n}\sum_{i=1}^{n}y_i$$

$$\sigma_x^2 = \frac{1}{n}\sum_{i=1}^{n}(x_i - \bar{x})^2 = \frac{1}{n}\sum_{i=1}^{n}(x_i^2 - \bar{x}^2)$$

$$\sigma_y^2 = \frac{1}{n}\sum_{i=1}^{n}(y_i - \bar{y})^2 = \frac{1}{n}\sum_{i=1}^{n}(y_i^2 - \bar{y}^2).$$

If $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$ be 'n' pairs of observations on 2 variables $x$ & $y$. then covariance of $x$ & $y$ is,

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

Cov indicates the joint variations b/w 2 variables. ∴ (c) coefficient b/w x & y is-

$$r = \frac{\text{Cov}(x,y)}{\sigma_x \; \sigma_y}$$

'r' can be written in diff. forms -

$$x_i = x - \bar{x}$$
$$y_i = y - \bar{y}$$

$$r = \frac{\sum x_i' y_i'}{n \; \sigma_x \; \sigma_y}$$

from ①

$$r = \frac{1}{n} \cdot \frac{\sum_{i=1}^{n}x_i' y_i'}{\sigma_x \cdot \sigma_y}$$

$$= \frac{\sum x_i' y_i'}{\sqrt{\sum x_i'^2} \times \sqrt{\sum y_i'^2}}$$

cancel 1/n from numerator & denomenator.

$$r = \frac{\sum_{i=1}^{n}x_i' y_i'}{\sqrt{\sum x_i'^2} \times \sqrt{\sum y_i'^2}}.$$

$$\text{Cov}(x, y) = \frac{1}{n}\sum(x_i - \bar{x})(y_i - \bar{y})$$

$$= \frac{1}{n}\left[\sum(xy_i - x\bar{y} - y_i\bar{x} + \bar{x}\bar{y})\right]$$

$$\left[(a-b)(c-d) = ac - ad - bc + bd\right]$$

$$= \frac{\sum x_i y_i}{n} - \bar{y}\frac{\sum x_i}{n} - \bar{x}\frac{\sum y_i}{n} + \bar{x}\bar{y}$$

$$= \frac{\sum x_i y_i}{n} - \bar{x}\bar{y} - \bar{x}\bar{y} + \bar{x}\bar{y}$$

$$= \frac{\sum x_i y_i}{n} - \bar{x}\bar{y}$$

$$\bar{x} = \frac{\sum x_i}{n}$$

$$\text{Cov}(x, y) = \frac{\sum x_i y_i}{n} - \frac{\sum x_i}{n} \times \frac{\sum y_i}{n}$$

$$r = \frac{\text{Cov}(x, y)}{\sigma_x \; \sigma_y} = \frac{\frac{\sum x_i y_i}{n} - \frac{\sum x_i}{n} \times \frac{\sum y_i}{n}}{\sqrt{\frac{\sum x_i^2}{n} - \left(\frac{\sum x_i}{n}\right)^2} \times \sqrt{\frac{\sum y_i^2}{n} - \left(\frac{\sum y_i}{n}\right)^2}}$$

$$\sigma_x^2 = \frac{\sum x_i^2}{n} \qquad \sigma_y^2 = \frac{\sum y_i^2}{n}$$

multiply by $n^2$,

$$\boxed{r = \frac{n\sum x_i y_i - \sum x_i \sum y_i}{\sqrt{n\sum x_i^2 - (\sum x_i)^2}\sqrt{n\sum y_i^2 - (\sum y_i)^2}}}$$

★ **Theorem :-** The (c) coefficient is independent of change of origin & scale of measurement.

**Proof**

Let $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$ be a set of 'n' pairs of observations.

$$r_{uv} = \frac{\frac{1}{n}\sum(x_i-\bar{x})(y_i-\bar{y})}{\sqrt{\frac{1}{n}\sum(x_i-\bar{x})^2}\sqrt{\frac{1}{n}\sum(y_i-\bar{y})^2}} \quad \text{——①}$$

Let us transform $x_i^o = u_i$, $y_i^o = v_i$

$$u_i^o = \frac{x_i^o - x_0}{c_1}, \qquad v_i^o = \frac{y_i^o - y_0}{c_2} \quad \text{——②}$$

$x_0, y_0, c_1, c_2$ are arbitrary constants.

from ②, we have

$$x_i^o = c_1 u_i^o + x_0 \qquad y_i^o = c_2 v_i^o + y_0$$

$$\bar{x} = x_0 + c_1\bar{u} \qquad \bar{y} = y_0 + c_2\bar{v}.$$

$$x_i^o - \bar{x} = c_1(u_i^o - \bar{u})$$
$$y_i^o - \bar{y} = c_2(v_i^o - \bar{v})$$

Substitute these into ①

$$r_{uv} = \frac{\frac{1}{n}\sum c_1(u_i^o - \bar{u}) \cdot c_2(v_i^o - \bar{v})}{\sqrt{\frac{1}{n}\sum c_1^2(u_i^o - \bar{u})^2}\sqrt{\frac{1}{n}\sum c_2^2(v_i^o - \bar{v})^2}}$$

$$= \frac{\frac{1}{n}\sum c_1 c_2(u_i^o - \bar{u})(v_i^o - \bar{v})}{\sqrt{\frac{1}{n}\sum(u_i^o - \bar{u})^2}\sqrt{\frac{1}{n}\sum(v_i^o - \bar{v})^2}}$$

$$r_{uv} = \frac{\frac{1}{n}\sum(u_i^o - \bar{u})(v_i^o - \bar{v})}{\sqrt{\frac{1}{n}\sum(u_i^o - \bar{u})^2}\sqrt{\frac{1}{n}\sum(v_i^o - \bar{v})^2}}$$

$$= \frac{\frac{1}{n}\sum u_i^o v_i^o - \bar{u}\bar{v}}{\sqrt{\frac{1}{n}\sum u_i^2 - \bar{u}^2}\sqrt{\frac{1}{n}\sum v_i^2 - \bar{v}^2}}$$

$$\boxed{r_{uv} = \frac{n\sum u_i v_i - \sum u_i \sum v_i}{\sqrt{n\sum u_i^2 - (\sum u_i)^2}\sqrt{n\sum v_i^2 - (\sum v_i)^2}}}$$

**Remark** – In general $\sigma_x = |c| \sigma_u$
$$\sigma_y = |d| \sigma_v$$

$$\therefore r_{xy} = \frac{cd}{|c|\cdot|d|} \cdot r_{uv}$$

Now $\frac{cd}{|c|\cdot|d|} = +1$ or $-1$ according as $c$, $d$ have same (opposite) sign.

(i) $r_{xy} = \pm r_{uv}$ according as $c$, $d$ have same (opposite) sign.

**Note** – In actual computations, we can take $c_1 = c_2 = 1$ so we assume

$u_i^o = x_i^o - x_0$ & $v_i^o = y_i^o - y_0$,

$x_0, y_0 \rightarrow$ should be chosen that most numerically $> x_i^o$ & $y_i^o$.

$\therefore$ If $u_i$ & $v_i$ respectively,

(c) coefficient remains unchanged

$$r_{uv} = \frac{\sum_{i=1}^{4}(u_i - \bar{u})(v_i - \bar{v})}{n\,\sigma_u\,\sigma_v}$$

(c) coefficient remains unchanged

Hence, $r_{uv}$ simplified as,

$$r_{uv} = \frac{\text{Cov}(u,v)}{\sigma_u \sigma_v} \qquad \left[\because \text{ 2nd } \text{Cov}(u,v) \atop \text{like } \text{Cov}(x,y)\right]$$

$\Rightarrow$ Derivation of Spearman's formula for Rank (c) coefficient :–

$$R = 1 - \frac{6\sum d^2}{n(n^2-1)}$$

Let $(x_1, y_1), \dots, (x_n, y_n)$ be the ranks of $n$ individuals in 2 characters of these 2 characters of

Eolcualrel Spearman's Rank Procoefficient ($c$)
$R$ is the product moment $(c)$
coefficient b/o these ranks,

$$R' = \frac{cov(x,y)}{\sigma_x \sigma_y}$$

$$cov(x,y) = \frac{\sum \{(x_i - \bar{x})(y_i - \bar{y})\}}{M}$$

$n = n \text{ road no.}$

$x_1, x_2, \dots x_n$ are $1, 2, 3, \dots n$ is same order.

$$\sum x = 1 + 2 + \dots n = \frac{n(n+1)}{2}$$

$$\sum x^2 = 1^2 + 2^2 \dots n^2 = \frac{n(n+1)(2n+1)}{6}$$

$$\bar{x} = \frac{\sum x}{n} = \frac{n+1}{2}$$

$$\frac{\sum x^2}{n} = \frac{(n+1)(2n+1)}{36}$$

$$\frac{n(n+1)}{2} = \frac{n+1}{2}$$

$$\therefore \sigma_x^2 = \frac{\sum x^2}{n} - \left(\frac{\sum x}{n}\right)^2 = \frac{(n+1)(2n+1)}{6} - \frac{(n+1)^2}{4}$$

$$x = 4(n+1)(2n+1) - 6(n+1)^2 \, (a+b)^2$$

$$x = \frac{4(n+1)(2n+1) - 6(n^2+2n+1)}{24}$$

$$= \frac{4(2n^2 + n + 2n+1) - 6(n^2+2n+1)}{24}$$

$$= \frac{8n^2 + 4n + 8n + 4 - 6n^2 - 12n - 6}{24}$$

$$= \frac{8n^2 + 4n + 4 - 6n^2 - 12n - 6}{24}$$

$$= \frac{8n^2 - 6n^2}{= 2n^2} = \frac{2n^2 - 2}{24} = \frac{2(n^2-1)}{24}_{12} = \frac{n^2-1}{12}$$

Similarly,
$$\bar{y} = \frac{n+1}{2}$$

$$\sigma_y^2 = \frac{n^2-1}{12}$$

Let $d_i = x_i - y_i$

$$d_i^2 = (x_i - \bar{x}) - (y_i - \bar{y})$$

$$\therefore \frac{\sum d_i^2}{n} = \frac{\sum [(x_i - \bar{x}) - (y_i - \bar{y})]^2}{n}$$

$$= \frac{\sum (x_i - \bar{x})^2 + \sum (y_i - \bar{y})^2 - 2\sum(x_i-\bar{x})(y_i-\bar{y})}{n} \quad (a-b)^2$$

$$= \sigma_x^2 + \sigma_y^2 - 2 \, cov(x,y)$$

$$2 \, cov(x,y) = \frac{n^2-1}{12} + \frac{n^2-1}{12} - \frac{\sum d_i^2}{n}$$

$$2 \, cov(x,y) = \frac{2(n^2-1)}{12} - \frac{\sum d_i^2}{n}$$

$$\therefore cov(x,y) = \frac{n^2-1}{12} - \frac{\sum d_i^2}{2n}$$

$$\therefore cov(x,y) = \frac{n^2-1}{12} - \frac{\sum d_i^2}{2n}$$

from eq $-(1)$,

$$R = \frac{\frac{n^2-1}{12} - \frac{\sum d_i^2}{2n}}{\sqrt{\frac{n^2-1}{12}} \times \sqrt{\frac{n^2-1}{12}}}$$

$$\sqrt{x} \times \sqrt{} = x$$

$$= \frac{n^2-1}{12} - \frac{\sum d_i^2}{2n} = \frac{\frac{n^2-1}{12} \cdot \frac{n^2-1}{12}}{\frac{n^2-1}{12}} - \frac{\frac{\sum d_i^2}{2n}}{\frac{n^2-1}{12}}$$

$$\boxed{R = 1 - \frac{6\sum d_i^2}{n(n^2-1)}}$$

$$= 1 - \frac{\frac{6}{d}}{\frac{a}{d}} = \frac{a\,d}{b\,c} = \frac{\sum d_i^2 \times 6}{2n(n^2-1)}$$

$$= 1 - \frac{6\sum d_i^2}{n(n^2-1)}$$

→ Limits of (c) coefficient :-

Now find the limits of (c) coefficient b/w 2 variables if we show that it lies b/w $-1$, $\sum + 1$ :

(i-c) $-1 \le \gamma_{xy} < +1$

**Proof**

Let $(x_i, y_i) \ldots (x_n, y_n)$ be given pairs of observations.

$$\gamma_{xy} = \frac{\frac{1}{n}\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\frac{1}{n}\sum(x_i - \bar{x})^2}\sqrt{\frac{1}{n}\sum(y_i - \bar{y})^2}} \quad —①$$

$$x_i' = x - \bar{x}$$
$$y_i' = y - \bar{y}$$

$$\sigma_x^2 = \frac{1}{n}\sum(x_i - \bar{x})^2 = \frac{\sum x_i'^2}{n}$$
$$\sigma_y^2 = \frac{1}{n}\sum(y_i - \bar{y})^2 = \frac{\sum y_i'^2}{n}$$

then let $x_i'$, $y_i'$.

$$\gamma_{xy} = \frac{\sum x_i' y_i'}{n\,\sigma_x\,\sigma_y} \quad —②$$

Now split eq.

$$\sum_{i=1}^{n}\left(\frac{x_i'}{\sigma_x} \pm \frac{y_i'}{\sigma_y}\right)^2 = \frac{\sum x_i'^2}{\sigma_x^2} + \frac{\sum y_i'^2}{\sigma_y^2} \pm 2\sum\frac{x_i'y_i'}{\sigma_x \sigma_y} \quad —③$$

$$(a \pm b)^2$$

$$—③ \div n$$

$$= \frac{n\sigma_x^2}{\sigma_x^2} + \frac{n\sigma_y^2}{\sigma_y^2} \pm 2\frac{n\sigma_x\sigma_y}{\sigma_x\sigma_y}\gamma_{xy}$$

$$= n + n \pm 2n\gamma_{xy}$$

$$= 2n \pm 2n\gamma_{xy} = 2n(1 \pm \gamma_{xy})$$

Left hand side of above identity is the sum of □ of n no. it tends to +ve/0.

$$1 \pm \gamma_{xy} > 0 \quad (\text{or}) \quad \gamma_{xy} \le 1 \ \& \ \gamma_{xy} \ge -1$$

Hence, $-1 \le \gamma_{xy} \le +1$

(or) $-1 \le \gamma_{xy} \le +1$

∴ (c) coefficient lies b/w $-1$ & $+1$.

1) find coefficient of (c),

| x : | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| y : | 6 | 8 | 11 | 9 | 12 | 10 | 14 |

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

$$\boxed{y = \frac{\sum(x-\bar{x})(y-\bar{y})}{\sqrt{\sum(x-\bar{x})^2}\sqrt{\sum(y-\bar{y})^2}}}$$

$$\Rightarrow \boxed{\frac{\sum xy}{\sqrt{\sum x^2}\sqrt{\sum y^2}}}$$

$$\bar{x} = \frac{\sum x}{n} \qquad \bar{y} = \frac{\sum y}{n}$$

# 1)

| X | Y | $x_i = x - \bar{x}$ | $y_i = y - \bar{y}$ | $x_i y$ | $x^2$ | $y^2$ |
|---|---|---|---|---|---|---|
| 1 | 6 | -3 | -4 | 12 | 9 | 16 |
| 2 | 8 | -2 | -2 | 4 | 4 | 4 |
| 3 | 11 | -1 | 1 | -1 | 1 | 1 |
| 4 | 9 | 0 | -1 | 0 | 0 | 1 |
| 5 | 12 | 1 | 2 | 2 | 1 | 4 |
| 6 | 10 | 2 | 0 | 0 | 4 | 0 |
| 7 | 14 | 3 | 4 | 0 | 9 | 16 |
| 28 | 70 | 0 | 0 | 29 | 28 | 42 |

$$\bar{x} = \frac{\sum x_i}{n} = \frac{28}{7} = 4 \qquad \bar{x} = \frac{28}{7} = 4$$

$$\bar{y} = \frac{\sum y}{n} = \frac{70}{7} = 10$$

$$\delta = \frac{1}{n} \sum x_i^\circ - \bar{x} \; . \; \sum y_i^\circ - \bar{y}$$

$$\sqrt{\frac{1}{n}(\sum x_i^\circ - \bar{x})^2} \cdot \sqrt{\frac{1}{n}(\sum y_i^\circ - \bar{y})^2}$$

$$\delta = \frac{\frac{1}{n} \sum x_i^\circ - \bar{x} \; . \; \sum y_i^\circ - \bar{y}}{\sqrt{\frac{1}{n}(\sum x_i^\circ - \bar{x})^2} \cdot \sqrt{\frac{1}{n}(\sum y_i^\circ - \bar{y})^2}}$$

$$= \frac{\sum x_i^\circ - \bar{x} \; . \; \sum y_i^\circ - \bar{y}}{\sqrt{\frac{1}{n}(\sum x_i^\circ - \bar{x})^2 \cdot \sum (y_i^\circ - \bar{y})^2}}$$

$$= \frac{\sum x_i^\circ - \bar{x} \; . \; \sum y_i^\circ - \bar{y}}{\sqrt{\sum(\sum x_i^\circ - \bar{x})^2 \cdot \sum(\sum y_i^\circ - \bar{y})^2}}$$

$$x_i^\circ - \bar{x} = x \qquad y_i^\circ - \bar{y} = y$$

$$\boxed{Y = \frac{\sum x \; y}{\sqrt{x^2} \cdot \sqrt{y^2}}}$$

$$\delta = \frac{29}{\sqrt{28} \cdot \sqrt{42}} \qquad x^2 = 28 \qquad y^2 = 42$$

$$= \frac{29}{\sqrt{28 \cdot 42}}$$

$$= \frac{29}{5.291 \times 6.928} = \frac{29}{36.650}$$

$$= \frac{29}{34.289} = 0.845$$

# 2)

Real Karl Pearson's coefficient of correlation b/w 2 variables $x$ & $y$ 0.28. Their covariance of $x$ & $y$ is 7.6. Of the variance of $x$ is 9, find standard deviation of $y$ series.

$$Y_{xy} = \frac{\text{Cov}(x,y)}{\sigma_x \, \sigma_y} \qquad \boxed{\sqrt{\text{variance}} = \text{Standard dev.}}$$

$$0.28 = \frac{7.6}{\sigma_x \, \sigma_y} \qquad \boxed{\text{std. dev} \cdot \text{and} \; \text{variance}}$$

$$0.28 = \frac{7.6}{3 \cdot \sigma_y} \qquad \Rightarrow \sigma_x^2 = 9 \quad \sigma_x = 3 \quad \sigma_y = ?$$

$$\sigma_y = \frac{7.6}{3 \times 0.28} \qquad Y_{xy} = 0.28 \qquad \text{cov}(x,y) = 7.6$$

$$\sigma_y = \frac{7.6}{0.84} = 9.047$$

# 3)

Calculate pearson's coefficient of correlation of from the following taking 100 & 50 as the assumed avg of $x$ & $y$.

| $x$: | 104 | 111 | 104 | 114 | 118 | 117 | 105 | 108 | 106 | 100 | 104 | 105 |
|------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| $y$: | 57 | 55 | 47 | 45 | 45 | 50 | 64 | 63 | 66 | 62 | 69 | 61 |

**1)**

| X | Y | $u = x-\bar{x}'$ | $v = y-\bar{y}$ | $u^2$ | $v^2$ | $uv$ |
|---|---|---|---|---|---|---|
| 104 | 57 | 4 | 7 | 16 | 49 | 28 |
| 111 | 55 | 11 | 5 | 121 | 25 | 55 |
| 104 | 47 | 4 | -3 | 16 | 9 | -12 |
| 114 | 45 | 14 | -5 | 196 | 25 | -70 |
| 118 | 45 | 18 | -5 | 324 | 25 | -90 |
| 117 | 50 | 17 | 0 | 289 | 0 | 0 |
| 105 | 64 | 5 | 14 | 25 | 196 | 70 |
| 108 | 63 | 8 | 13 | 64 | 169 | 104 |
| 106 | 66 | 6 | 16 | 36 | 256 | 96 |
| 100 | 62 | 0 | 12 | 0 | 144 | 0 |
| 104 | 69 | 4 | 19 | 16 | 361 | 76 |
| 105 | 61 | 5 | 11 | 25 | 121 | 55 |
|   |   | 96 | 84 | 1128 | 1380 | 312 |

$n = 12$

$$r = \frac{n\sum u_i v_i - \sum u_i \sum v_i}{\sqrt{n\sum u_i^2 - \sum(u_i)^2}\sqrt{n\sum v_i^2 - \sum(v_i)^2}}$$

$$= \frac{12 \times 312 - 96 \times 84}{\sqrt{12 \times 1128 - (96)^2}\sqrt{12 \times 1380 - (84)^2}}$$

$$= \frac{3744 - 8064}{\sqrt{13536 - 9216}\,\sqrt{16560 - 7056}}$$

$$= \frac{-4320}{\sqrt{13440}\,\sqrt{16476}} = \frac{-4320}{\sqrt{13536 - 9216}\,\sqrt{16560 - 7056}}$$

$$= \frac{-4320}{\sqrt{4320}\,\sqrt{9504}}$$

$$= \frac{-4320}{\text{...}} = -0.674$$

$$= \frac{-4320}{64.7595493}$$

$$= -0.674$$

**4)** Cal the (C) of correlation for the following ages of husbands & wifes.

Age of (H) (x) = 23, 27, 28, 29, 30, 31, 33, 35, 36, 39
Age of (w) (y) = 18, 22, 23, 24, 25, 26, 28, 29, 30, 32

| X | Y | $\bar{x} = x-\bar{x}$ | $\bar{y} = y-\bar{y}$ | $x_i^2$ | $y_i^2$ | $x_i y_i$ |
|---|---|---|---|---|---|---|
| 23 | 18 | -8.7 | -7.7 | 65.61 | 59.29 | 62.37 |
| 27 | 22 | -4.7 | -3.7 | 16.81 | 13.69 | 15.18 |
| 28 | 23 | -3.1 | -2.7 | 9.61 | 7.29 | 8.37 |
| 29 | 24 | -2.1 | -1.7 | 4.41 | 2.89 | 3.57 |
| 30 | 25 | -1.1 | -0.7 | 1.21 | 0.49 | 0.77 |
| 31 | 26 | -0.1 | 0.3 | 0.01 | 0.09 | -0.03 |
| 33 | 28 | 1.9 | 2.3 | 3.61 | 5.29 | 4.37 |
| 35 | 29 | 3.9 | 3.3 | 15.21 | 10.89 | 12.67 |
| 36 | 30 | 4.9 | 4.3 | 24.01 | 18.49 | 21.07 |
| 39 | 32 | 7.9 | 6.3 | 62.41 | 39.69 | 49.77 |
|   |   |   |   | 202.9 | 158.1 | 178.36 |

$$\bar{x} = \frac{311}{10} = 31.1$$

$$\bar{y} = \frac{257}{10} = 25.7$$

$$r = \frac{\sum x_i y_i}{\sqrt{\sum x_i^2}\,\sqrt{\sum y_i^2}}$$

$$= \frac{178.36}{\sqrt{202.9}\,\sqrt{158.1}}$$

$$= \frac{178.36}{14.2442971 \times 12.5737822\text{...}}$$

$$= \frac{178.36}{\text{...}}$$

$$= \frac{178.36}{179.10469} = 0.995$$

5) Cal (c) of correlation.

X : 6  2  10  4  8
Y : 9  11  5  8  7

$$\gamma = \frac{n\Sigma xy - \Sigma x \Sigma y}{\sqrt{n\Sigma x^2 - (\Sigma x)^2}\sqrt{n\Sigma y^2 - (\Sigma y)^2}}$$

| X | Y | $x^2$ | $y^2$ | x y |
|---|---|---|---|---|
| 6 | 9 | 36 | 81 | 54 |
| 2 | 11 | 4 | 121 | 22 |
| 10 | 5 | 100 | 25 | 50 |
| 4 | 8 | 16 | 64 | 32 |
| 8 | 7 | 64 | 49 | 56 |
| 30 | 40 | 220 | 340 | 214 |

$$\bar{x} = \frac{30}{} \qquad \bar{y} = \frac{40}{}$$

$$\gamma = \frac{n\Sigma xy - \Sigma x \Sigma y}{\sqrt{n\Sigma x^2 - (\Sigma x)^2}\sqrt{n\Sigma y^2 - (\Sigma y)^2}}$$

$$= \frac{5 \times 214 - 30 \times 40}{\sqrt{5 \times 220 - (30)^2}\sqrt{5 \times 340 - (40)^2}} = \frac{1070 - 1200}{\sqrt{200}\sqrt{100}}$$

$$= \frac{-130}{141.4213562 \times 10}$$

$$= \frac{-130}{141.4213562}$$

$$= -0.919$$

6) 

| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|
Marks in maths | 78 | 36 | 98 | 25 | 75 | 82 | 90 | 62 | 65 | 39
Marks in stats | 84 | 51 | 91 | 60 | 68 | 62 | 86 | 58 | 53 | 47

Cal. Rank correlation coefficient.

| Rollno | Mathematics Mark | Ranks R1 | Statistics Mark | Ranks R2 | $d$ | $d^2$ |
|---|---|---|---|---|---|---|
| 1 | 78 | 4 | 84 | 3 | 1 | 1 |
| 2 | 36 | 9 | 51 | 9 | 0 | 0 |
| 3 | 98 | 1 | 91 | 1 | 0 | 0 |
| 4 | 25 | 10 | 60 | 6 | 4 | 16 |
| 5 | 75 | 5 | 68 | 4 | 1 | 1 |
| 6 | 82 | 3 | 62 | 5 | -2 | 4 |
| 7 | 90 | 2 | 86 | 2 | 0 | 0 |
| 8 | 62 | 7 | 58 | 7 | 0 | 0 |
| 9 | 65 | 6 | 53 | 8 | -2 | 4 |
| 10 | 39 | 8 | 47 | 10 | -2 | 4 |
|  |  |  |  |  |  | 30 |

$$R = 1 - \frac{6\Sigma d^2}{n(n^2-1)}$$

$$= 1 - \frac{6 \times 30}{10(10^2-1)}$$

$$= 1 - \frac{180}{10 \times 99}$$

$$= 1 - \frac{180}{990}$$

$$= 1 - 0.1818$$

$$= 0.8181$$

→ Spearman's formula for repeated ranks :-

If in a series 2|more individuals have the same score when we find the avg of the ranks of these individuals ^allote these avg rank to each of them.

$$eg = 98 \rightarrow 7, 8, 9.$$

⟹ The score 98 occurs 3 times, there is a time of 7, 8, 9 place.

$$= \frac{7+8+9}{3} = \frac{24}{3} = 8$$

$$\boxed{S \cdot \text{formula} = 1 - \frac{6 \sum d^2 + \frac{t^3-t}{12}}{n(n^2-1)}}$$

1) find rank (co) coefficient

| Series A : | 115 | 109 | 112 | 87 | 98 | 120 | 98 | 100 | 98 | 118 |
|---|---|---|---|---|---|---|---|---|---|---|
| Series B : | 75 | 73 | 85 | 70 | 76 | 82 | 65 | 73 | 68 | 80 |

ⓐ

| Series A | | Series B | | $d = x-y$ | $di^2$ |
|---|---|---|---|---|---|
| Score | Rank(x) | Score | Rank(y) | | 4 |
| 115 | 3 | 75 | 5 | -2 | 4 |
| 109 | 5 | 73 | 6·5 | -1·5 | 2·25 |
| 112 | 4 | 85 | 1 | 3 | 9 |
| 87 | 10 | 70 | 8 | 2 | 4 |
| 98 | 8 | 76 | 4 | 4 | 16 |
| 120 | 1 | 82 | 2 | -1 | 1 |
| 98 | 8 | 65 | 10 | 2 | 4 |
| 100 | 6 | 73 | 6·5 | 0·5 | 0·25 |
| 98 | 8 | 68 | 9 | 1 | 1 |
| 118 | 2 | 80 | 3 | -1 | 1 |
| | | | | | 42·5 |

based rank 73 → 2 times = $\frac{6·5}{2}$ = 6·5

98 occurs 3 times = $\frac{}{3}$

$$R = 1 - \frac{6\left[\Sigma d^2 + \frac{t^3-t}{12}\right]}{n(n^2-1)}$$

$$\Sigma \frac{t^3-t}{12} = 2 + 0.5 = 2.5$$

(a) 
$$t = 3$$
$$\frac{3^3-3}{12} = \frac{24}{12} = 2$$

(b) 
$$t = 2$$
$$\frac{2^3-2}{12} = \frac{6}{12} = 0.5$$

$$R = 1 - 6 \times \frac{\left[42.5 + 2.5\right]}{10(10^2-1)}$$

$$= 1 - 6 \times \frac{\left[45\right]}{10 \times 99}$$

$$= 1 - \frac{270}{10 \times 99}$$

$$= 1 - \frac{8}{35}$$

$$= 1 - \frac{3}{11} = \frac{11-3}{11} = \frac{8}{11} = 0.72$$

2) Cal Rank (C) coefficient from following data specifying the ranks of 7 students in 2 subjects.

(x) Rank in 1st Sub : 1   2   3   4   5   6   7
(y) Rank in 2nd Sub : 4   3   2   6   5   1   7

3) The coefficient of rank (C) of marks obtained by 10 students in 2 subjects was computed as 0.5. It was later discovered that the difference in 2 subjects was wrongly taken as 3 instead of 7. Find the correct coefficient of rank (C)?

---

A) 
$$R = 0.5$$
$$n = 10$$

$$R = 1 - \frac{6 \Sigma d^2}{n(n^2-1)}$$

$$0.5 = 1 - \frac{6 \Sigma (x-y)^2}{10 \times 99}$$

$$0.5 = 1 - \frac{6 \Sigma (x-y)^2}{990}$$

$$(0.5-1) 990 = 6 \Sigma (x-y)^2$$

$$-495 = -6 \Sigma (x-y)^2$$

$$\Sigma (x-y)^2 = \frac{495}{6} = 82.5$$

$$0.5 \times 990 = 6 \Sigma (x-y)^2$$

$$0.5 = \frac{990 - 6 \Sigma (x-y)^2}{990}$$

$$\Rightarrow 82.5 - \frac{3^2 + 7^2}{}$$

wrong → correct →

$$= 82.5 - 9 + 49 = 122.5$$

$$R = 1 - \frac{6 \times 122.5}{990}$$

$$= \frac{990 - 6 \times 122.5}{990} = \frac{990 - 735}{990}$$

$$R = 0.2575$$

$$d = x_i - y_i$$
$$\frac{3}{7} \checkmark$$
$$(x_i - y_i)^2 = 82.5 \checkmark$$
$$82.5 - 3^2 + 7^2$$
$$122.5$$
$$\Sigma (x_i - y_i)^2$$

→ Stop here

2.A)

| $x_i$ | $y_i$ | $d_i = x_i - y_i$ | $d_i^2$ |
|---|---|---|---|
| 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 |
| | | | $20$ |

$R = 1 - \dfrac{6 \Sigma d_i^2}{n(n^2-1)}$

$= 1 - \dfrac{6 \times 20}{7(48)}$

$= 1 - \dfrac{120}{336} = 1 - 0.357 = 0.643$

Q @ (R) (variable $x^2 \cdot r^2 = x/6$

---

$+ 3x^2 = 12 x^2 = \dfrac{9x^2 - 1}{2\sqrt{x}(1+x^2)}$

$\dfrac{2\sqrt{x}(1+3x^2)}{} = \dfrac{9x^2 - 1}{2\sqrt{x}(1+x^2)}$

$\Rightarrow$ Regression :— (R)

Suppose are given $n$ pairs of values of 2 variables $x$ & $y$. If we fit a straight line to this data by taking $x$ as independent variable & $y$ as dependent variable., then the extraight line obtained → Regression line of $y$ on $x$. ≈ regression line of $x$ on $y$.

The reciprocal of its slope $c$ → regression coefficient of $x \equiv x$ on $y$.

→ Eq. boy (R) lines :—

Let $y = a + bx$ — ① be the eq of (R) line of $y$ on $x$.

If $a$ & $b$ are determined by the normal eq obtained by principle of least $\square$'s.

$\Sigma y_i = na + b \Sigma x_i$ — ②
$\Sigma x_i y_i = a \Sigma x_i + b \Sigma x_i^2$ — ③

② ÷ n .

$\dfrac{\Sigma y_i}{n} = a + \dfrac{b \Sigma x_i}{n}$

$\bar{y} = a + b\bar{x}$ — ④

---

It is not required that ①

regression (xx) by its while the other is true. The value of (xx)

otherwise ① (cannot exceed unity)

∴ If I of the (R) (xx) exceed unity, the other must be K than unity.

∴ If the (R) (xx) exceed unity, the other must be K than unity.

## 2. A)

$\bar{x}, \bar{y}$ - are the means of $x$ & $y$ series

(i)-(ii)

$y - \bar{y} = bx - b\bar{x}$

$y - \bar{y} = b(x - \bar{x})$ —⑤

② × $\sum x_i$ - ① × n

$\sum x_i y_i = na \sum x_i + b \sum x_i^2$ —⑥

$n \sum x_i y_i = na \sum x_i + nb \sum x_i^2$ —⑦

⑦-⑥

$n \sum x_i y_i - \sum x_i y_i = nb \sum x_i^2 - b \sum x_i^2$

$n \sum x_i y_i - \sum x_i \sum y_i = b(n \sum x_i^2 - b \sum x_i^2)$

$n \sum x_i^2 - \sum x_i^2 = b(n \sum x_i^2 - \sum x_i^2)$ → $\left(\frac{1}{n}, n^2\right)$

$b = \dfrac{n \sum x_i y_i - \sum x_i \sum y_i}{n \sum x_i^2 - \sum x_i^2}$

$= \dfrac{\dfrac{\sum x_i y_i}{n} - \dfrac{\sum x_i \sum y_i}{n^2}}{\dfrac{\sum x_i^2}{n} - \dfrac{\sum x_i^2}{n^2}}$

$b = \dfrac{\dfrac{\sum x_i y_i}{n} - \bar{x}\bar{y}}{\dfrac{\sum x_i^2}{n} - \bar{x}^2} = \dfrac{cov(x,y)}{\sigma_x^2} = \dfrac{P_{xy}}{\sigma_x^2}$

Sub ② & ⑤

The eqⁿ of y on x. eqⁿ-

$y - \bar{y} = \dfrac{P_{xy}}{\sigma_x^2}(x - \bar{x})$ —⑧

---

Simultaneously when $x$ is depending on $y$, the ® eqⁿ of $x$ on $y$ is obtained as -

$(x - \bar{x}) = \dfrac{P_{xy}}{\sigma_y^2}(y - \bar{y})$ —⑨

let us denote $\dfrac{P_{xy}}{\sigma_x^2}$ as $b_{yx}$ &

$\dfrac{P_{yx}}{\sigma_y^2}$ as $b_{xy}$

Thus $b_{yx} = \dfrac{P_{xy}}{\sigma_x^2}$ as , $b_{xy} = \dfrac{P_{xy}}{\sigma_y^2}$

Hence $b_{yx}$ → ® coefficient of $y$ on $x$

and $b_{xy}$ → ® coefficient of $x$ on $y$.

The (R) eqⁿ of $y$ on $x$ -

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

The (R) eqⁿ of $x$ on $y$ -

$(x - \bar{x}) = b_{xy}(y - \bar{y})$     $x → y$

**Remarks :-**

* slope of ® line of y on x, $b_{yx} = \dfrac{Y_{xy}}{\sigma_x}$
* slope of ® line of x on y

    Reciprocal of $b_{xy}$ → $\dfrac{\sigma_y}{Y_{xy}}$

* since ® $b_{yx} = r\left(\sigma_y / \sigma_x\right)$ eqⁿ of x are $\dfrac{\sigma_x}{\sigma_y}$ are
  the reciprocal of $b_{xy} = r\left(\sigma_y / \sigma_x\right)$ eqⁿ of y are same sign as that of $b_{yx}$
* since cor $b_{yx} = r(\sigma_y / \sigma_x)$ eqⁿ of x are $\dfrac{\sigma_y}{\sigma_x}$ are same sign

    the, it follows that r has same sign as that of $b_{yx}$

2. A)

* Since $b_{xy} = x\left(\dfrac{\sigma_x}{\sigma_y}\right)$ are equally
  find $\sqrt{ant} \ (b_{xy}) \ (b_{xy}) = x^2$.
* Since AM is always $> GM$ for any
  2 no... are true
  $\therefore$ AM of $b_{xy}, b_{yx}$ is always $>$
  $\dfrac{1}{2}(b_{yx} + b_{xy}) = \sqrt{b_{yx} \times b_{xy}} = |x|$.

* coefficient of ©.
* 2 lines of ® always pass through
  the point $(\bar{x}, \bar{y})$.

* (R) eq of $y$ on $x$ is used for estimation,
  the value of $y$ for a given value of
  $x$. Ee the ® eq of $x$ on $y$ is used
  for estimating $x$ for a specified
  value $B, y$.

* Correlation

| | Regression |
|---|---|
| * (c) means the ship b/s 2/more variables | ® means act of returning to the avg value. |
| * (c) measures the degree of relation ship b/s the variables | ® measures the nture of relation ship b/s variables |
| * There may be nonsense correlation b/s 2 variables | No such correspondance to culcutta is |

* very useful has further mathematical treatment.

| res col has further mathematical treatment. |

ⓓ Find the most likely price in Bombay
corresponding to the price $\bar{x} = 70$ out
calcutta. —

1) Find the most likely price in Bombay
corresponding to the price $\bar{x} = 70$ out
calcutta —

avg price in calcutta 65, avg price at
Bombay 67, SD out calcutta 2.5, SD at
Bombay 3.5, (c) of correlation b/s the
prices in 2 cities is 0.8.

[ $x \to$ price of calcutta
  $y \to$ " of Bombay.

$\bar{x} = 65$       $\bar{y} = 67$,
$\sigma_x = 2.5$     $\sigma_y = 3.5$.     $x = 0.8$.

$\therefore$ Eq of the ® of $y$ on $x$,
$y - \bar{y} = \dfrac{x \cdot \sigma_y}{\sigma_x}(x - \bar{x})$

$y - 67 = \dfrac{0.8 \times 3.5}{2.5}(x - 65)$

where $x = 70$

$y - 67 = 1.12(x - 65)$

$\therefore y = 72.6$.

that most likely price in Bombay
corresponding to the price of Rs. 70
at culcutta is Rs. 72.6

$(y-67)=1.12(70-65)$
$y-67 = 1.12(5)$
$y-67 = 5.6$
$y = 5.6 + 67$
$y = 72.6$

2) If $x_i$ & $y_i$ are devi- of $n$ pairs of values of $x$ & $y$ from means (i.e if $x_i^0 = x_i - \bar{x}$ & $y_i^0 = y_i - \bar{y}$, respectively means (i.e if $x_i^0 = x_i - \bar{x}$ & $y_i^0 = y_i - \bar{y}$.

P. that :-

a) $\gamma = 1 - \dfrac{1}{2n} \sum \left(\dfrac{x_i^0}{\sigma_x} - \dfrac{y_i^0}{\sigma_y}\right)^2$

b) $\gamma = -1 + \dfrac{1}{2n} \sum \left(\dfrac{x_i^0}{\sigma_x} - \dfrac{y_i^0}{\sigma_y}\right)^2$

hence deduce that $-1 \le \gamma \le +1$.

A) a)

$\sum \left(\dfrac{x_i^0}{\sigma_x} - \dfrac{y_i^0}{\sigma_y}\right)^2 = \sum \left[\dfrac{x_i^{0^2}}{\sigma_x^2} - \dfrac{2x_i^0 y_i^0}{\sigma_x \sigma_y} + \dfrac{y_i^{0^2}}{\sigma_y^2}\right]$

$= \dfrac{n\sigma_x^2}{\sigma_x^2} - \dfrac{2n\gamma}{\sigma_x\sigma_y} + \dfrac{n\sigma_y^2}{\sigma_y^2}$

$\boxed{\begin{aligned} x_i^{0^2} &= n\sigma_x^2 \\ y_i^{0^2} &= n\sigma_y^2 \\ \dfrac{x_i y_i^0}{\sigma_x \sigma_y} &= n\gamma \end{aligned}}$

$= 2n - 2n\gamma = 2n(1-\gamma)$

$\therefore 2n\gamma = 2n - \sum \left(\dfrac{x_i^0}{\sigma_x} - \dfrac{y_i^0}{\sigma_y}\right)^2$

$\gamma = 1 - \dfrac{1}{2n} \sum \left(\dfrac{x_i^0}{\sigma_x} - \dfrac{y_i^0}{\sigma_y}\right)^2$

b)

$\gamma = -1 + \dfrac{1}{2n} \sum \left(\dfrac{x_i^0}{\sigma_x} - \dfrac{y_i^0}{\sigma_y}\right)^2$

$= -1 + \dfrac{1}{2n} (2n + 2n\gamma) = \gamma$

$\dfrac{1}{2n} \sum \left(\dfrac{x_i^0}{\sigma_x} + \dfrac{y_i^0}{\sigma_y}\right)^2$ are

$\dfrac{1}{2n} \sum \left(\dfrac{x_i^0}{\sigma_x} - \dfrac{y_i^0}{\sigma_y}\right)^2$ are

Since $\sum$ are non -ve.

$\therefore \gamma \le 1$ & $\gamma \ge -1$

(i.e) $-1 \le \gamma \le +1$

3) If $r$ is (c) of correlation b/w n pairs of values of $x$ & $y$ & $\sigma_x, \sigma_y$, are SD of x & y & $(x-y)$. P. that.

$\gamma = \dfrac{\sigma_x^2 + \sigma_y^2 - \sigma_{x-y}^2}{2\sigma_x \sigma_y}$

$\sum \left(\dfrac{x_i^0}{\sigma_x} - \dfrac{y_i^0}{\sigma_y}\right)^2 =$

$= \dfrac{n\sigma_x^2 - 2n\gamma + n\sigma_y^2}{\sigma_x^2}$

$= 2n - 2n\gamma$

$= 2n(1-\gamma) \to (2n-2n\gamma)$

$\to 2n(1-\gamma)$

**ⅰ)** Let $u_i = x_i - y_i$

$$\bar{u} = \bar{x} - \bar{y}$$

where $\bar{x}, \bar{y}, \bar{u}$ are mean of $x, y, u$.

$$u_i - \bar{u} = (x_i - \bar{x}) - (y_i - \bar{y}) \qquad (a)$$

$(a-b)^2$

$$(u_i - \bar{u})^2 = (x_i - \bar{x})^2 - 2(x_i - \bar{x})(y_i - \bar{y}) + (y_i - \bar{y})^2 \qquad (b)$$

$$\sum (u_i - \bar{u})^2 = \sum (x_i - \bar{x})^2 - 2\sum (x_i - \bar{x})(y_i - \bar{y}) + (y_i - \bar{y})^2$$

$$\sum (u_i - \bar{u})^2 = \sum (x_i - \bar{x})^2 - 2\sum (x_i - \bar{x})(y_i - \bar{y}) + \sum (y_i - \bar{y})^2$$

(ⅲ) →

$$\sum (u_i - \bar{u})^2 = \sum (x_i - \bar{x})^2 + \sum (y_i - \bar{y})^2 - 2\sum (x_i - \bar{x})(y_i - \bar{y})$$

$$n\sigma_u^2 = n\sigma_x^2 + n\sigma_y^2 - 2\sum (x_i - \bar{x})(y_i - \bar{y})$$

$$\sigma_{x-y}^2 = \sigma_x^2 + \sigma_y^2 - \frac{2\sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

**(ⅱ)**

$$\sigma_{x-y}^2 = \sigma_x^2 + \sigma_y^2 - 2 r \sigma_x \sigma_y$$

$$\sigma_{x-y}^2 = \sigma_x^2 + \sigma_y^2 - 2 r \sigma_x \sigma_y$$

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

---

**ⅲ)** If $z = ax + by$, $\text{co } x$ is (c) of correlation b/s $x$ & $y$, $s.\text{that} -$

Deduce that,

$$\sigma_z^2 = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab r \sigma_x \sigma_y$$

$$\sigma_{x-y}^2 = \sigma_x^2 + \sigma_y^2 - 2 r \sigma_x \sigma_y \qquad (ⅰ)$$

$$\sigma_{x+y}^2 = \sigma_x^2 + \sigma_y^2 + 2 r \sigma_x \sigma_y \qquad (ⅱ)$$

$$\sigma_{x+y}^2 = \sigma_x^2 + \sigma_y^2 + 2 r \sigma_x \sigma_y \qquad (ⅲ)$$

∴ that $\sigma_{x+y}$ is $> \sigma_{x-y}$ then $\sigma_{x-y}$ according as $r$ is +ve / -ve.

**ⅲ)** $z = ax + by$, $\bar{z} = a\bar{x} + b\bar{y}$

Let $z_i = ax_i + by_i$

$$z_i - \bar{z} = a(x_i - \bar{x}) + b(y_i - \bar{y})$$

$$(z_i - \bar{z})^2 = a^2 (x_i - \bar{x})^2 + b^2 (y_i - \bar{y})^2 + 2ab(x_i - \bar{x})(y_i - \bar{y})$$

$\sum$ to all,

$$\sum (z_i - \bar{z})^2 = a^2 \sum (x_i - \bar{x})^2 + b^2 \sum (y_i - \bar{y})^2 + 2ab \sum (x_i - \bar{x})(y_i - \bar{y})$$

**(ⅱ)**

$$n\sigma_z^2 = \frac{na^2 \sigma_x^2 + nb^2 \sigma_y^2 + 2ab \sum (x_i - \bar{x})(y_i - \bar{y})}{n}$$

$$\sigma_z^2 = \frac{a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{n}}{}$$

$$\sigma_z^2 = a^2 \sigma_x^2 + b^2 \sigma_y^2 + 2ab r \sigma_x \sigma_y$$

→ (ⅰ) Proved

$a = 1$    $b = 1$ → ①

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

$a = 1, \ b = 1$ → ①

$\sigma_x^2 + \sigma_y^2 + 2\sigma_x\sigma_y$

$\sigma_x^2 + y \ > < \ \sigma_{x-y}$

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

$\sigma \ > \ \text{or} \ < 0$

Thus $\sigma_{x+y} \geq \sigma_{x-y}$ according

as $r$ is +ve / -ve.

5) If $\theta$ is acute angle b/w $\theta$ (R) lines relating the variables $x$ & $y$. Attach—

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

Indicate the significance of the cases of
(i) $r = 0$ & $r = \pm 1$.

d) $b_{xy} = r\dfrac{\sigma_y}{\sigma_x}$

$b_{xy} = r \cdot \dfrac{1}{r} \cdot \dfrac{\sigma_y}{\sigma_x}$

---

$$\tan\theta = \frac{b_{xy} - b_{yx}}{1 + b_{yx}\,b_{xy}}$$

$$= \frac{\dfrac{1}{r}\dfrac{\sigma_y}{\sigma_x} - \dfrac{r\sigma_y}{\sigma_x}}{1 + \dfrac{r\sigma_y}{\sigma_x} \cdot \dfrac{1}{r}\dfrac{\sigma_y}{\sigma_x}} \qquad \frac{\sigma_y}{\sigma_x}$$

$$= \frac{\dfrac{1}{r}\dfrac{\sigma_y}{\sigma_x} - \dfrac{r\sigma_y}{\sigma_x}}{1 + \dfrac{\sigma_y^2}{\sigma_x^2}}$$

$$= \frac{\dfrac{1}{r}\sigma_y\sigma_x - r\sigma_y\sigma_x}{\sigma_x^2 + \sigma_y^2} \times \frac{\sigma_x^2}{\sigma_x^2}$$

$$= \frac{\dfrac{1}{r}\sigma_y\sigma_x - r\sigma_y\sigma_x}{\sigma_x^2 + \sigma_y^2}$$

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

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

when $r = \pm 1$

$\theta = 0$

$\therefore$ when $r = 0$, $\theta = \dfrac{\pi}{2}$

when $r = \pm 1$

$\theta = 0$

$\dfrac{1 - \frac{1}{1}}{1} = 0$

$\therefore$ the lines (R) coincide