

Curve fitting

Let x be an independent variable and y be a variable depending on x ; Here we say that y is a function of x and write it as $y = f(x)$. If $f(x)$ is a known function, then for any allowable 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 and thereby 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 and y give us n points on the curve $y = f(x)$.

Suppose we consider the converse problem. That is, suppose we are given n values x_1, x_2, \dots, x_n of an independent variable x and corresponding values y_1, y_2, \dots, y_n of a variable y depending on x . Then the pairs $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$ give us n points in the xy -plane. Generally, it is not possible to find the actual curve $y = f(x)$ that passes through these points. Hence we try to find a curve that serves as best approximation to the curve $y = f(x)$. Such a curve is referred to as the curve of best fit. The

process of determining a curve of best fit is called curve fitting. The method generally employed for curve fitting is known as the method of least squares which is explained below.

Method of least squares

This is a method for finding the unknown coefficients in a curve that serves as best approximation to the curve $y = f(x)$. The basic ideas of this method were created by A.M. Legendire and C.F. Gauss.

"The principle of least squares says that the sum of the squares of the error between the observed values and the corresponding estimated values should be the least."

Suppose it is desired to fit a k -th degree curve given by

$$y = a_0 + a_1 x + a_2 x^2 + \dots + a_k x^k \quad \dots (1)$$

to the given pairs of observations $(x_1, y_1), (x_2, y_2) \dots (x_n, y_n)$. The curve has $k + 1$ unknown constants and hence if $n = k + 1$ we get $k + 1$ equations on substituting the values of (x_i, y_i) in equation (1). This gives unique solution to the values $a_0 a_1 a_2 \dots a_n$. However, if $n > k + 1$, no unique solution is possible and we use the method of least squares.

Now let

$y_e = a_0 + a_1 x + a_2 x^2 + \dots + a_k x^k$ be the estimated value 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 or error for this point,

$$e_i = y_i - y_e = y_i - a_0 - a_1 x_i - a_2 x_i^2 - \dots - a_k x_i^k$$

To make the sum of squares minimum, we have to minimise.

$$S = \sum_{i=1}^n e_i^2 = \sum_{i=1}^n (y_i - a_0 - a_1 x_i - a_2 x_i^2 - \dots - a_k x_i^k)^2 \quad \dots (2)$$

By differential calculus, S will have its minimum value when

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

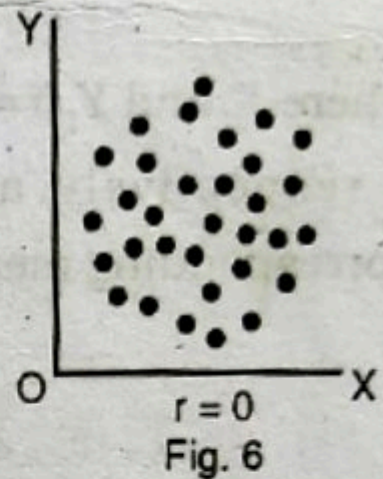
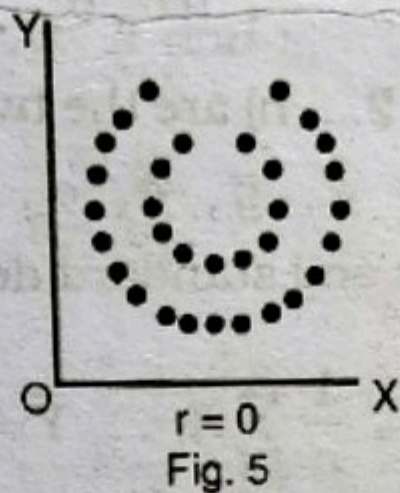
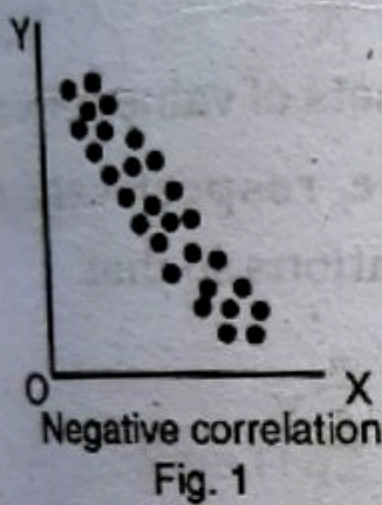
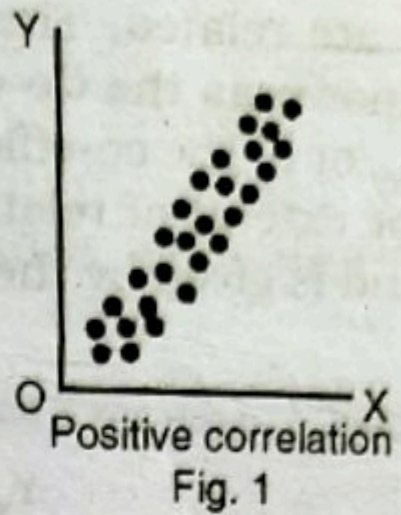
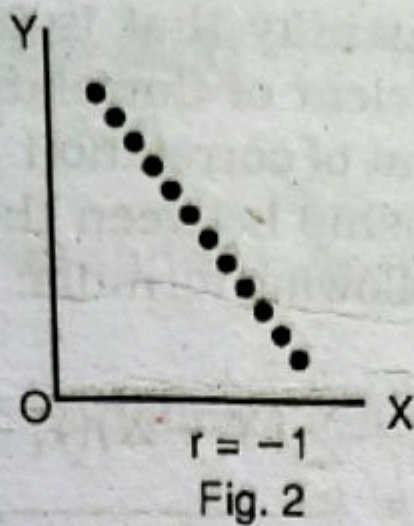
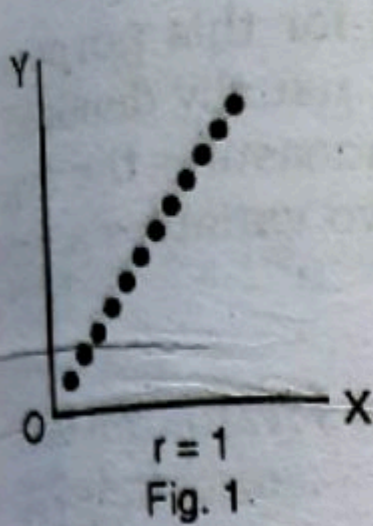
Scatter Diagram

The existence of correlation can be shown graphically by means of a *scatter diagram*. Statistical data relating to simultaneous movements (or variations) of two variables can be graphically represented by points. One of the two variables, say X, is shown along the horizontal axis OX and the other variable Y along the vertical axis OY. All the pairs of values of X and Y are now shown by points (or dots) on the graph paper. This diagrammatic representation of bivariate data is known as scatter diagram.

The scatter diagram of these points and also the direction of the scatter reveals the nature and strength of correlation between the two variables. The following are some scatter diagrams showing different types of correlation between two variables.

In Fig. 1 and 3, the movements (or variations) of the two variables are in the same direction and the scatter diagram shows a linear path. In this case, correlation is positive or direct.

In Fig. 2 and 4, the movements of the two variables are in opposite directions and the scatter shows a linear path. In this case correlation is negative or indirect.



In Fig. 5 and 6 points (or dots) instead of showing any linear path lie around a curve or form a swarm. In this case correlation is very small and we can take $r = 0$.

In Fig. 1 and 2, all the points lie on a straight line. In these cases correlation is perfect and $r = +1$ or -1 according as the correlation is positive or negative.

Problem 5.

Scatter diagram.

correlation (assoc coeff)
regression (at coeff)

1. Fit a straight line of the form $y = a + bx$ to the following data by the method of least squares.

x	0	1	3	6	8
y	1	3	2	5	4

Here, there are 2 variables a and b so we have 2 normal equations.

$$y = a + bx$$

$$\sum y = na + b \sum x \quad \text{--- 1st normal eqn}$$

$$\sum xy = a \sum x + b \sum x^2 \quad \text{--- 2nd normal eqn.}$$

x	y	x^2	xy	y^2
0	1	0	0	1
1	3	1	3	9
3	2	9	6	4
6	5	36	30	25
8	4	64	32	16
<u>18</u>	<u>15</u>	<u>110</u>	<u>71</u>	<u>55</u>

$$\sum y = na + b \sum x \Rightarrow 15 = 5a + 18b \quad \text{--- (1)}$$

$$\sum xy = a \sum x + b \sum x^2 \Rightarrow 71 = 18a + 110b \quad \text{--- (2)}$$

From (1) we get,

$$a = \frac{15 - 18b}{5} \quad \text{--- (3)}$$

Sub (3) in (2)

$$\begin{aligned} 71 &= 18 \left(\frac{15 - 18b}{5} \right) + 110b \\ &= 54 - \frac{324b}{5} + 110b \end{aligned}$$

$$\begin{aligned} 17 &= 110b - \frac{324b}{5} \\ &= b \left(110 - \frac{324}{5} \right) \end{aligned}$$

$$17 = b \times \frac{226}{5}$$

$$b = \frac{85}{226} = 0.376$$

$$a = \frac{15 - 18 \times 0.376}{5}$$

$$a = 1.6464$$

$$\therefore a = 1.6464 \quad \& \quad b = 0.376$$

where $x = at + by$

normal equations given by,

$$\sum x = na + b \sum y \quad \text{--- (4)}$$

$$\sum xy = a \sum y + b \sum y^2 \quad \text{--- (5)}$$

From (4), & (5)

$$18 = 5 \times a + 15b \quad \text{--- (6)}$$

$$71 = 15a + 55b \quad \text{--- (7)}$$

$$18 - 15b = 5a$$

$$a = \frac{18 - 15b}{5}$$

Sub in (7).

$$71 = 15 \left(\frac{18 - 15b}{5} \right) + 55b$$

$$= 54 - 45b + 55b$$

$$17 = 10b$$

$$b = \underline{1.7}$$

$$a = \frac{18 - 15 \times 1.7}{5} = -1.5$$

$$\underline{a = -1.5 \text{ and } b = 1.7}$$

2. Given a table of values for the function. Fit a parabola (second degree polynomial) to this data.

x	1.0	1.5	2.0	2.5	3.1	4.0
y	5.1	5.3	5.6	5.7	5.9	6.1
x	1.0	1.5	2.0	2.5	3.1	4.0
y	1.1	1.3	1.6	2.0	3.4	4.2

A). Here there are 3 variables, \therefore 3 normal equations.

$$y = a + bx + cx^2 \quad \text{--- (1)}$$

$$\sum y = na + b\sum x + c\sum x^2 \quad \text{--- (2)}$$

$$\sum xy = a\sum x + b\sum x^2 + c\sum x^3 \quad \text{--- (3)}$$

$$\sum x^2 y = a\sum x^2 + b\sum x^3 + c\sum x^4 \quad \text{--- (4)}$$

x	y	x^2	x^3	x^4	xy	$x^2 y$
1.0	1.1	1	1	1	1.1	1.1
1.5	1.3	2.25	3.375	5.0625	1.95	2.925
2.0	1.6	4	8	16	3.2	6.4
2.5	2.0	6.25	15.625	39.0625	5	12.5
3.1	3.4	9.61	29.791	92.3521	10.54	32.674
4.0	4.2	16	64	256	16.8	67.2
<u>14.1</u>	<u>13.6</u>	<u>39.11</u>	<u>121.791</u>	<u>409.471</u>	<u>38.59</u>	<u>122.799</u>

$$(2) \Rightarrow$$

$$13.6 = 6a + 14.1b + 39.11c \quad \text{--- (5)}$$

$$(3) \Rightarrow$$

$$38.59 = 14.1a + 39.11b + 121.79c \quad \text{--- (6)}$$

$$(4) \Rightarrow$$

$$122.799 = 39.11a + 121.79b + 409.47c \quad \text{--- (7)}$$

Solve the three equations using gaussian elimination.

$$\begin{bmatrix} 6 & 14.1 & 39.11 & 13.6 \\ 14.1 & 39.11 & 121.79 & 38.59 \\ 39.11 & 121.79 & 409.47 & 122.799 \end{bmatrix}$$

$$R_1 \rightarrow R_1/6$$

$$\begin{bmatrix} 1 & 2.35 & 6.52 & 2.27 \\ 14.1 & 39.11 & 121.79 & 38.59 \\ 39.11 & 121.79 & 409.47 & 122.79 \end{bmatrix}$$

$$R_2 \rightarrow R_2 - 14.1 R_1$$

$$R_3 \rightarrow R_3 - 39.11 R_1$$

$$\begin{bmatrix} 1 & 2.35 & 6.52 & 2.27 \\ 0 & 5.97 & 29.86 & 6.58 \\ 0 & 29.88 & 154.47 & 34.01 \end{bmatrix}$$

$$R_2 \rightarrow \frac{1}{5.97} R_2$$

$$\begin{bmatrix} 1 & 2.35 & 6.52 & 2.27 \\ 0 & 1 & 5 & 1.10 \\ 0 & 29.88 & 154.47 & 34.01 \end{bmatrix}$$

$$R_3 \rightarrow R_3 - 29.88 R_2$$

$$\begin{bmatrix} 1 & 2.35 & 6.52 & 2.27 \\ 0 & 1 & 5 & 1.10 \\ 0 & 0 & 5.07 & 1.142 \end{bmatrix}$$

$$5.07c = 1.142$$

$$c = 0.225$$

$$b + 5c = 1.10$$

$$b + 5 \times 0.225 = 1.10$$

$$b = 1.10 - 1.125$$

$$= -0.025$$

$$a + 2.35b + 6.52c = 2.27$$

$$a + 2.35 \times (-0.025) + 6.52 \times 0.225 = 2.27 \quad a =$$

$$a + -0.059 + 1.467 = 2.27$$

$$a = 2.27 - 1.408$$

$$= 0.862$$

$$a = 0.862, b = -0.029, c = 0.225$$

regression line
 intersection of
 regression lines
 (\bar{x}, \bar{y})
 sign of regression
 coeff same.
 sign of correlation coeff
 same as the regression
 coefficients.

? $20x + 3y = 5 \rightarrow$ regression line of y on x

$$3y = 5 - 20x$$

$$y = \frac{5}{3} - \frac{20}{3}x$$

$-\frac{20}{3} \rightarrow$ regression coefficient.

? Fit a straight line to the following set of data, find the regression lines, regression coefficients, \bar{x}, \bar{y} also calculate the correlation coeff for the data and comment on the result.

x	1	2	3	4	5	6	7
y	0.5	2.5	2.0	4.0	3.5	6.0	5.5

A.

x	y	x^2	xy	y^2
1	0.5	1	0.5	0.25
2	2.5	4	5	6.25
3	2.0	9	6	4
4	4.0	16	16	16
5	3.5	25	17.5	12.25
6	6.0	36	36	36
7	5.5	49	38.5	30.25
$\frac{28}{28}$	$\frac{24}{24}$	$\frac{140}{140}$	$\frac{119.5}{119.5}$	$\frac{105}{105}$

$y = a + bx$

Normalized eqns,

$$\sum y = na + b\sum x$$

$$\sum xy = a\sum x + b\sum x^2$$

$$24 = 7a + 28b \Rightarrow 28b = 24 - 7a$$

$$19.5 = 28a + 140b$$

$$19.5 = 28a + 140 \left(\frac{24 - 7a}{28} \right)$$

$$= 28a + 5(24 - 7a)$$

$$= 28a + 120 - 35a$$

$$= -7a + 120$$

$$7a = 120 - 19.5$$

$$a = 0.07$$

$$b = 0.8396$$

regression line of y on x, $y = a + bx$

$$\Rightarrow y = 0.07 + 0.8396x$$

regression line of x on y, $x = c + dy$

normal equations,

$$\sum x = nc + \sum yd$$

$$\sum xy = c \sum y + d \sum y^2$$

$$28 = 7c + 24d \quad \text{--- (1)}$$

$$19.5 = 24c + 105d \quad \text{--- (2)}$$

From (1), we get,

$$7c = 28 - 24d$$

$$c = \frac{(28 - 24d)}{7}$$

Sub c in (2)

$$19.5 = 24 \left(\frac{28 - 24d}{7} \right) + 105d$$

$$= \frac{672 - 576d}{7} + 105d$$

$$= 96 - \frac{576d}{7} + 105d$$

$$23.5 = 22.714d$$

$$d = 1.0346$$

$$c = \frac{28 - 24(1.0346)}{7}$$
$$= \underline{\underline{0.4528}}$$

Regression line of x on y , $\Rightarrow x = 0.45 + (1.03)d$

Regression coeff. of x on $y = 1.0346$

" " " y on $x = 0.8396$

correlation coeff = ^{geometric} mean of regression coeffs

$$= r = 0.93$$

Since correlation coeff = 0.93 ≈ 1

\therefore positive correlation.

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

$$\bar{y} = \frac{\sum y}{n} = \frac{24}{7} = \underline{\underline{3.4286}}$$

correlation.

when the changes in one variable are associated or followed by changes in other is called correlation.

If an increase (or decrease) in the values of one variable corresponds to an increase (or decrease) in the other, the correlation is said to be positive. If the increase (or decrease) in one corresponds to the decrease (or increase) in the other, the correlation is said to be negative.

If there is no relationship indicated b/w the variables, they are said to be independent or uncorrelated.

Regression Lines

Suppose we are given n pairs of values $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ of two variables x and y . If we fit a straight line to this data by taking x as independent variable and y as dependent variable. then the straight line obtained is called the regression line of y on x . Its slope is called

the coeff. of y on x . Similarly if we fit a straight line to the data by taking y as independent variable and x as dependent variable, the line obtained is the regression line of x on y , the reciprocal of its slope is called the regression coeff. of x on y .

Equation for regression lines

Let $y = a + bx$ be the equation of the regression line of y on x , where a and b are determined by solving the normal equations obtained by the principle of least squares.

$$\begin{aligned} \sum y &= na + b \sum x \\ \sum xy &= a \sum x + b \sum x^2 \end{aligned}$$

here b is called regression coeff. of y on x

similarly,

$x = a + by$ is the eqn of regression line of x on y ,

normal equations,

$$\begin{aligned} \sum x &= na + b \sum y \\ \sum xy &= a \sum y + b \sum y^2 \end{aligned}$$

here b is called regression coeff. of x on y .

correlation coeff: G.M of regression coefficients

$$r = \frac{\sum XY}{\sqrt{(\sum X^2 \sum Y^2)}} \Rightarrow r = \frac{\sum XY}{\sqrt{(\sum X^2 \sum Y^2)}}$$

where $X = x - \bar{x}$, $Y = y - \bar{y}$ (center)

$\sigma_x = \text{S.D of } x \text{ series}$, $\sigma_y = \text{S.D of } y \text{ series}$.

In a partially destroyed lab record, only the lines of regression of y on x and x on y are available as $4x - 5y + 33 = 0$ and $20x - 9y = 107$ respectively. Calculate \bar{x} , \bar{y} and the coeff. of correlation b/w them.

Given,

regression line of y on x , $4x - 5y + 33 = 0$

$$\rightarrow -5y = -4x - 33$$

$$5y = 4x + 33$$

regression line of x on y , $20x - 9y = 107$

$$20x = 107 + 9y$$

Intersection of regression lines gives the point, (\bar{x}, \bar{y})

$$4x - 5y = -33 \quad \text{--- (1)}$$

$$20x - 9y = 107 \quad \text{--- (2)}$$

Multiply (1) by 5

$$20x - 25y = -165 \quad \text{--- (3)}$$

$$(3) - (2) \Rightarrow -16y = -272$$

$$y = 17 //$$

$$y = 17 \text{ in (1)} \Rightarrow 4x - 85 = -33$$

$$x = 13 //$$

$$(\bar{x}, \bar{y}) = (13, 17)$$

correlation coeff: G.M of regression coefficients

regression line of y on x ,

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

regression line of x on y ,

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

$$\text{correlation coeff: } \sqrt{\frac{4}{5} \cdot \frac{9}{20}} = \sqrt{\frac{9}{25}}$$

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

Note:

- * correlation coefficient always lies between -1 and 1
- * both the sign of regression coeff. of y on x and regression coeff. of x on y are same.
- * sign of correlation coeff is same as the sign of regression coeff.

(regression coeff. or sign) -ve \Rightarrow correlation coeff. positive

(~~interchange~~) change the sign of correlation coeff.

Rank correlation coefficient.

Rank correlation is based on the rank or the order and not on the magnitude of the variable.

If the ranks assigned to individuals range from 1 to n , then the Karl Pearson's correlation coeff. b/w a series of ranks is called Rank correlation coefficient.

Edward Spearman's formula for Rank correlation coefficient (R) is given by,

$$R = 1 - \frac{6 \sum d_i^2}{n(n^2-1)} \quad \text{or} \quad 1 - \frac{6 \sum d_i^2}{(n^3-n)}$$

where d is the difference b/w the ranks of the 2 series, and n is the no. of individuals in each series.

10 participants in a contest are ranked by 2 judges as follows,

calculate rank correlation coefficient.

X = 1 6 5 10 3 2 4 9 7 8

Y = 6 4 9 8 1 2 3 10 5 7

X	Y	d	d ²
1	6	5	25
6	4	2	4
5	9	4	16
10	8	2	4
3	1	2	4
2	2	0	0
4	3	1	1
9	10	1	1
7	5	2	4
8	7	1	1
		<u>1</u>	<u>60</u>

$$R = 1 - \frac{62d^2}{n^3 - n}$$

$$= 1 - \frac{6 \times 60}{10^3 - 3}$$

$$= \underline{\underline{0.6363}}$$

2. The judges A, B, C gives the following ranking. Find which pairs of judges has common approach.

A	1	6	5	10	3	2	4	9	7	8
B	3	6	8	4	7	10	2	1	6	9
C	6	4	9	8	1	2	3	10	5	7

A)

A	B	C	d_{AB}	d_{AB}^2	d_{BC}	d_{BC}^2	d_{AC}	d_{AC}^2
1	3	6	2	4	3	9	5	25
6	6	4	0	0	2	4	2	4
5	8	9	3	9	1	1	4	16
10	4	8	6	36	4	16	2	4
3	7	1	4	16	6	36	2	4
2	10	2	8	64	8	64	0	0
4	2	3	2	4	1	1	1	1
9	1	10	8	64	9	81	1	1
7	6	5	1	1	1	1	2	4
8	9	7	1	1	2	4	1	1
				200		217		60

$$R_{AB} = 1 - \frac{62d_{AB}^2}{n^3 - n}$$

$$= 1 - \frac{6 \times 200}{10^3 - 10} = 0.21$$

$$RBC = 1 - \frac{6 \sum d^2 c^2}{n^3 - n}$$

$$= 1 - \frac{6 \times 217}{10^3 - 10}$$

$$= -0.315$$

$$RAC = 1 - \frac{6 \times 60}{10^3 - 10}$$

$$= 0.63$$

Since $R(A, C)$ is maximum, the pair of judges A and C have the nearest common approach.

problems using the formula, Coeff. of correlation

$$= \frac{\sum XY}{\sqrt{(\sum X^2 \sum Y^2)}}$$

psychological test of intelligence and engineering ability were applied to 10 students. There is a record of ungrouped data showing intelligence ratio (I.R) and engineering ratio (E.R). calculate the coeff. of correlation.

Student	A	B	C	D	E	F	G	H	I	J
I.R	105	104	102	101	100	99	98	96	93	92
E.R	101	103	100	98	95	96	104	92	97	94

Student	x	$x = x - \bar{x}$	y	$y = y - \bar{y}$	x^2	y^2	xy
A	105	6	101	3	36	9	18
B	104	5	103	5	25	25	25
C	102	3	100	2	9	4	6
D	101	2	99	0	4	0	0
E	100	1	95	-3	1	9	-3
F	99	0	96	-2	0	4	0
G	99	-1	104	6	1	36	-6
H	96	-3	92	-6	9	36	18
I	93	-6	97	-1	36	1	6
J	92	-7	94	-4	49	16	28
	990		980		170	140	92

$$r = \frac{\sum xy}{\sqrt{(\sum x^2)(\sum y^2)}} = \frac{92}{\sqrt{170 \times 140}}$$

$$\bar{x} = \frac{990}{10} = 99$$

$$\bar{y} = \frac{980}{10} = 98$$

$$= \underline{\underline{0.89}}$$