

## Correlation:

Correlation is the degree of the relationship between two or more variables. It does not explain the cause behind the relationship.

e.g.: price and supply, demand and supply, income and expenditure are correlated.

### Types of correlation

Positive correlation: Correlation is said to be positive when two variables move in same direction.

Ex-1 If income and expenditure of a household may be increasing or decreasing simultaneously. If so, there is positive correlation.

Ex-2  $y = a + bn$  (e.g.  $y = 2x + 3$ )

Negative correlation: Correlation is said to be negative when two variables move in opposite direction.

e.g. Determinants of Quantity demanded.

$$Q_d = f(P)$$

Ex-1 Prices and demand for a commodity move in the opposite direction.

Ex-2  $y = a - bn$ .

Type II Based on number of variables.

i) Simple correlation.

ii) Multiple correlation.

iii) Partial correlation.

**Simple correlation:**

When two variables are studied it is called simple correlation.

**Multiple correlation:**

When at least three variables are studied and their relationships are simultaneously worked out it is a case of multiple correlation.

e.g.: Determinants of Quantity demanded.

$$Q_d = f(P, P_c, P_s, t, y)$$

where  $Q_d$  stands for quantity demanded,  $f$  stands for function,  $P$  is the price of the goods,  $P_c$  is the price of competitive goods,  $P_s$  is the price of substituting goods,  $t$  is the taste and preference,  $y$  is the income.

**Partial correlation:**

When more than two variables are studied keeping the other variables constant, it is called Partial Correlation.

Type III Based on change in proportion.

**Linear correlation:** Correlation is said to be linear when one variable moves with the other variable in fixed proportion.

**Non linear correlation:** Correlation is said to be non linear when one variable moves with the other variable in changing proportion.

$$\text{e.g. } y = a + bx^2$$

## Scale of correlation coefficient

Value of r.

$0 \leq r \leq 0.19$

Very low

Correlation.

$0.2 \leq r \leq 0.39$

Low correlation.

$0.4 \leq r \leq 0.59$

Moderate.

$0.6 \leq r \leq 0.79$

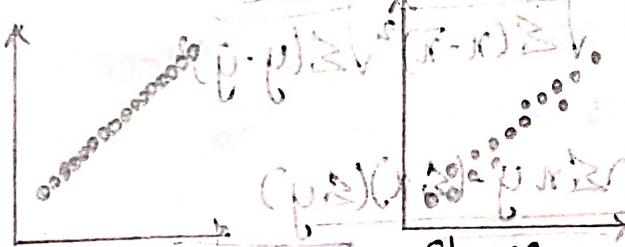
High correlation.

$0.8 \leq r \leq 1.0$

Very high

Correlation.

$$(p-p_1)(x-x_1) = r : \text{positive trend}$$



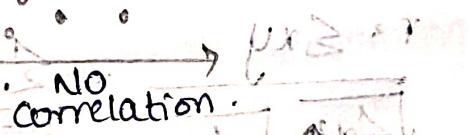
Perfect  
Positive

Strong  
Positive  
correlation.

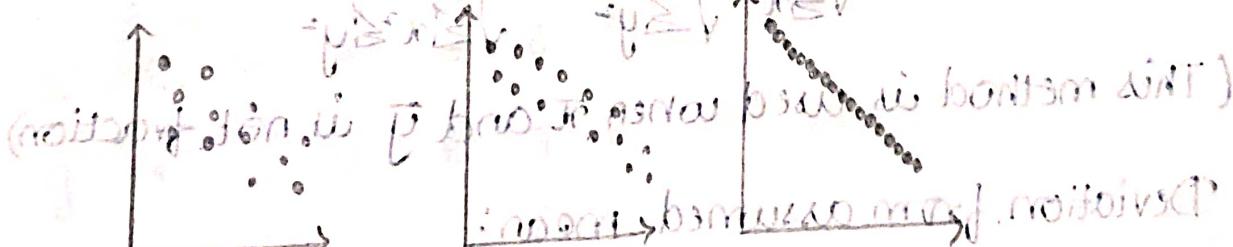
Weak.  
Positive.  
correlation.

Correlation: (absolute value is constant)

(absolute value



No  
correlation.



Weak  
negative.

Strong  
Negative

Perfect  
negative

Correlation

Correlation.

$$(p-p_1)(x-x_1) + (p-p_2)(x-x_2) = f(x) = 0$$

Karl Pearson's coefficient correlation is a popular method of calculating correlation. Arithmetic mean and standard deviation are the basis for its calculation. The correlation coefficient (r) also called as the linear correlation coefficient measures the strength and direction of a linear relationship between two variables. The value of r lies between -1 and 1.

Deviation from Actual mean

$$\text{Indirect method: } r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}}$$

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

(This method is used when given variable are small magnitude)

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

(This method is used when  $\bar{x}$  and  $\bar{y}$  is not fraction)

Deviation from assumed mean:

Direct method

$$r = \frac{n \sum d_x d_y - (\sum d_x)(\sum d_y)}{\sqrt{n \sum d_x^2 - (\sum d_x)^2} \sqrt{n \sum d_y^2 - (\sum d_y)^2}}$$

(This method is used when  $\bar{x}$  and  $\bar{y}$  is fraction or the series is large)

Karl Pearson's coefficient of correlation:

Karl Pearson's coefficient of correlation is a mathematical and most popular method of calculating correlation.

Arithmetic mean and standard deviation are the basis for its calculation. The correlation coefficient (r) also called as the linear correlation coefficient measures the strength and direction of a linear relationship between two variables. The value of r lies between -1 and 1.

Deviation from Actual mean

$$\text{Indirect method: } r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2} \sqrt{\sum (y - \bar{y})^2}}$$

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

(This method is used when given variable are small magnitude)

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}}$$

(This method is used when  $\bar{x}$  and  $\bar{y}$  is not fraction)

Deviation from assumed mean:

Direct method

$$r = \frac{n \sum d_x d_y - (\sum d_x)(\sum d_y)}{\sqrt{n \sum d_x^2 - (\sum d_x)^2} \sqrt{n \sum d_y^2 - (\sum d_y)^2}}$$

(This method is used when  $\bar{x}$  and  $\bar{y}$  is fraction or the

series is large)

$d_x = x - A$ ,  $d_y = y - B$  where  $A, B$  are constants.

$n$ : number of pair of observations

$\sum d_x$  = sum of deviation of  $x$ .

$\sum d_y$  = sum of deviation of  $y$ .

$\sum d_x^2$  = sum of deviation square of  $x$ .

$\sum d_y^2$  = sum of deviation square of  $y$ .

$\sum d_x d_y$  = sum of product of deviation  $x$  and  $y$ .

The Pearson correlation is the most used measurement for a linear relationship between two variables.

The stronger the correlation between these two data sets, the closer it will be to  $+1$  or  $-1$ .

$$r = \frac{\sum d_x d_y}{\sqrt{\sum d_x^2} \sqrt{\sum d_y^2}}$$

Calculate the KPCC. from the following data.

$x$	$y$	$x^2$	$y^2$	$xy$
2	18	4	324	36
4	12	16	144	48
5	10	25	100	50
6	8	36	64	48
8	7	64	49	56
11	5	121	25	55
$\sum x = 36$		$\sum y = 60$	$\sum x^2 = 266$	$\sum y^2 = 706$
$\sum xy = 293$				

$$r = \frac{6(266) - (36)(60)}{\sqrt{6(266) - (36)^2} \sqrt{6(706) - (60)^2}}$$

$$r = \frac{6(266) - (36)(60)}{\sqrt{6(266) - (36)^2} \sqrt{6(706) - (60)^2}}$$

$$\frac{1758 - 2160}{\sqrt{1596 - 1296} \cdot \sqrt{4236 - 3600}}$$

$$\frac{1758 - 2160}{\sqrt{300} \cdot \sqrt{636}}$$

$$\frac{-402}{(26.38)(25.2)} = -0.92$$

$$\frac{-402}{680} = -0.59$$

Note:

With precision limit of 0.0042, the value

		Positive Correlation			
negative correlation.		Strong -ve	Weak +ve	weak +ve	Strong +ve
-1.00	-0.50	0.00	0.01	0.02	0.03
		0.04	0.05	0.06	0.07
		0.08	0.09	0.10	0.11
		0.12	0.13	0.14	0.15
		0.16	0.17	0.18	0.19
		0.20	0.21	0.22	0.23
		0.24	0.25	0.26	0.27
		0.28	0.29	0.30	0.31
		0.32	0.33	0.34	0.35
		0.36	0.37	0.38	0.39
		0.40	0.41	0.42	0.43
		0.44	0.45	0.46	0.47
		0.48	0.49	0.50	0.51
		0.52	0.53	0.54	0.55
		0.56	0.57	0.58	0.59
		0.60	0.61	0.62	0.63
		0.64	0.65	0.66	0.67
		0.68	0.69	0.70	0.71
		0.72	0.73	0.74	0.75
		0.76	0.77	0.78	0.79
		0.80	0.81	0.82	0.83
		0.84	0.85	0.86	0.87
		0.88	0.89	0.90	0.91
		0.92	0.93	0.94	0.95
		0.96	0.97	0.98	0.99

Perfect negative correlation.

Karl Pearson correlation coefficient b/w Age and weight.

Age	1	2	3	4	5	6	7	8	9	10
weight	3	4	9	6	7	12	10	11	13	14

Age (x)	weight (y)	$x^2$	$y^2$	$xy$	$\sum x$	$\sum y$	$n = 10$
1	3	1	9	3	15	32	
2	4	4	16	8	26	36	
3	9	9	81	27	41	45	
4	6	16	36	24	50	42	
5	12	25	144	60	75	76	
$\sum x = 15$	$\sum y = 32$	$\sum x^2 = 55$	$\sum y^2 = 254$	$\sum xy = 115$			

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

$$\sqrt{5(55) - (15)^2} \sqrt{5(254) - (32)^2}$$

$$\frac{585 - 480}{\sqrt{275 - 225} \sqrt{508 - 1024}} = \frac{105}{\sqrt{50} \sqrt{246}}$$

$$\frac{105}{\sqrt{50} \sqrt{246}} = \frac{105}{\sqrt{1000} \sqrt{246}} = \frac{105}{\sqrt{24600}} = \frac{105}{156.8} = 0.68$$

$0.68 \Rightarrow$  shows +ve correlation

$$\frac{105}{110.85} = 0.95 \Rightarrow$$

calculate correlation coefficient b/w death and birth rate for the data.

Birth rate Death  $n = n_i - \bar{x}$   $y_i - \bar{y}$   $x^2 - \bar{x}^2$   $y^2 - \bar{y}^2$

	15	-6	-7	36	49	42
24	20	-4	-2	16	4	8
26	22	2	0	4	0	0
32	24	9	2	9	4	6
33	27	5	5	25	25	25
35	24	0	2	0	4	0
30	24	0	0	90	86	81

$$r = \frac{\sum xy}{\sqrt{n^2 - \sum x^2} \sqrt{n^2 - \sum y^2}} = \frac{81}{\sqrt{90} \sqrt{86}}$$

$$= \frac{81}{9.4 \times 9.2} = 0.93.$$

$$\text{Ans} = \frac{81}{91} = 0.89$$

$$r = 0.89$$

Ans = 0.89

Problem: Psychological tests of intelligence and business analytic ability were applied to 10 students. There is a record of ungrouped data showing Intelligence Ratio (IR) and business analytic Ability Ratio (BAAR). Calculate the coefficient of correlation.

Students	IR (xi)	BAAR (yi)	$\bar{x} = \frac{\sum x_i}{n} = \frac{990}{10} = 99$	$\bar{y} = \frac{\sum y_i}{n} = \frac{980}{10} = 98$	$r = \frac{\sum xy - \bar{x}\bar{y}}{\sqrt{\sum x^2 - \bar{x}^2} \sqrt{\sum y^2 - \bar{y}^2}}$
A	105	101	6	3	0.59
B	104	103	5	5	0.59
C	102	100	3	2	0.59
D	101	98	2	0	0.59
E	100	95	1	-3	0.59
F	99	96	0	-2	0.59
G	98	104	-1	+6	0.59
H	96	92	-3	-6	0.59
I	93	97	-6	-1	0.59
J	92	94	-7	-4	0.59
					$\sum x^2 = 170$
					$\sum y^2 = 140$

$$\bar{x} = \frac{\sum x_i}{n} = \frac{990}{10} \Rightarrow \bar{x} = 99.$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{980}{10} \Rightarrow \bar{y} = 98.$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{92}{\sqrt{170} \sqrt{140}} = \frac{92}{\sqrt{154}} = 0.59.$$

Problem: Psychological tests of intelligence and business analytic ability were applied to 10 students. There is a record of ungrouped data showing Intelligence Ratio (IR) and business analytic Ability Ratio (BAAR). Calculate the coefficient of correlation.

Students	IR (xi)	BAAR (yi)	$\bar{x} = \frac{\sum xi}{n} = \frac{990}{10} = 99$	$\bar{y} = \frac{\sum yi}{n} = \frac{980}{10} = 98$	$x_i - \bar{x}$	$y_i - \bar{y}$	$x_i^2$	$y_i^2$	$xy$
A	105	101	6	3	36	9	1296	81	315
B	104	103	5	5	25	25	625	625	275
C	102	100	3	2	9	4	81	16	24
D	101	98	2	0	4	0	16	0	0
E	100	95	1	-1	-3	1	9	1	-3
F	99	96	0	-2	0	4	0	4	0
G	98	104	-1	+6	1	36	-1	1296	36
H	96	92	-3	-6	-9	36	81	144	16
I	93	97	-6	-1	-36	1	36	1	-6
J	92	94	-7	-4	-49	16	49	256	21
					170	140	990		

$$\bar{x} = \frac{\sum xi}{n} = \frac{990}{10} = 99.$$

$$\bar{y} = \frac{\sum yi}{n} = \frac{980}{10} = 98.$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \cdot \sum y^2}} = \frac{92}{\sqrt{170 \cdot 140}} = \frac{92}{\sqrt{1825 \cdot 13}} = 0.59.$$

Problem Find if there is any significant correlation between heights and weights for following data.

height $x_i$	weight $y_i$	$\bar{x} = \frac{\sum x_i}{n} = \frac{540}{90} = 60$	$\bar{y} = \frac{\sum y_i}{n} = \frac{1080}{90} = 120$	$\sum x_i - \bar{x}$	$\sum y_i - \bar{y}$	$x^2 - \bar{x}^2$	$xy$
57	113			-3	-7	9	49
57	117			-1	-3	1	3
59	126						
62	126						
63	126						
64	130			5	9	25	81
65	129						
55	111					4	16
58	116			-2	-4	1	8
57	112			-3	-8	9	72
						<u>109</u>	<u>472</u>
							<u>216</u>

$\therefore r = 1.00$

$$\bar{x} = \frac{\sum x_i}{n} = \frac{540}{90} = 60$$

$$\bar{y} = \frac{\sum y_i}{n} = \frac{1080}{90} = 120$$

$$r = \frac{\sum xy}{\sqrt{\sum x^2} \sqrt{\sum y^2}} = \frac{216}{\sqrt{109} \sqrt{472}} = \frac{216}{219.41} = 0.98$$

Spearman's rank correlation coefficient

Ranks are not equal:

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

Where  $d = R_1 - R_2$

$d$  is difference in a pair of ranks

$R$  is rank correlation coefficient

$R_1$  is rank of observations with respect to first variable.

$R_2$  is rank of observations with respect to second variable.

$n$  is the number of paired observations

Ranks are equal:

$$R = 1 -$$

$$\frac{6 \left[ \sum d^2 + \frac{1}{12} (m_1^3 - m_1) + \frac{1}{12} (m_2^3 - m_2) + \frac{1}{12} (m_{12}^3 - m_{12}) \right]}{n(n^2-1)}$$

$$= \frac{6P_1^2 + 6P_2^2 + 6P_{12}^2 - 3(P_1 + P_2 + P_{12})^2}{12 \cdot P_1 \cdot P_2 \cdot P_{12}}$$

Calculate SRCC for the following data.

Ranking economics	Ranking statistics	$d = R_1 - R_2$	$d^2$
1	3	-3	9
2	8	-6	36
3	2	1	1
4	3	0	0
5	7	-2	4
6	4	-1	1
7	9	-2	4
8	10	-1	1
9	1	6	36
10	6	0	0

$$n = 10$$

$$\sum d = 132$$

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

$$1 - \frac{6 \times 132}{10(100-1)}$$

$$1 - \frac{792}{990}$$

$$1 - 0.8 = 0.2$$

$$R_s = 0.2$$

Calculate.

x	17	13	15	16	6	11	14	9	7	12	
y	36	46	35	24	12	18	27	22	12	8	

x	y	$R_1$	$R_2$	$d = R_1 - R_2$	$d^2$
17	36	1	2	-1	1
13	46	5	1	4	16
15	35	3	3	0	0
16	24	2	5	-3	9
6	12	10	8	2	4
11	18	7	7	0	0
14	27	4	4	0	0
9	22	8	6	2	4
7	2	9	10	-1	1
12	8	6	9	-3	9

$$n = 10$$

$$\leq d^2 = 44.$$

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

$$1 - \frac{6 \times 44}{10(100-1)}$$

$$1 - \frac{264}{990}$$

$$1 - 0.26$$

$$R = 0.74$$

9)	x	40	42	35	45	30	50	
	y	110	120	125	130	128	160.	

x	y	$R_1$	$R_2$	$d = R_1 - R_2$	$d^2$
40	110	4	6	-2	4
42	120	3	5	-2	4
35	125	5	4	1	1
45	130	2	2	0	0
30	128	6	3	3	9
50	160.	1	1	0	0

$$n = 6$$

$$\leq d^2 = 18$$

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

$$= 1 - \frac{6 \times 18}{6(36-1)}$$

$$= 1 - \frac{108}{210}$$

$$= 1 - 0.51$$

$$= 0.49$$

calculate between 2 variables. x and y.

$$x \quad y \quad R_1 \quad R_2 \quad d = R_1 - R_2 \quad d^2$$

x	y	R <sub>1</sub>	R <sub>2</sub>	d = R <sub>1</sub> - R <sub>2</sub>	d <sup>2</sup>
68	62	4.5	5.5	-1	1
58	6	7	7	0	0
64	68	2.5	3.5	-1	1
75	45	9	10	-1	1
50	45	9	10	-1	1
64	81	6	5	1	1
80	60	1	6	-5	25
75	68	2.5	3.5	-1	1
40	48	10	9	1	1
55	50	8	8	0	0
64	70	1.6	2	-.4	.16

$$\frac{2+2+2+2}{72} = \frac{8}{72}$$

$$m_1 = 2$$

$$m_2 = 3$$

$$m_3 = 2$$

$$1 - \frac{6 \left[ 72 + \frac{1}{12}(8-2) + \frac{1}{12}(27-3) + \frac{1}{12}(8-2) \right]}{100(100-1)}$$

$$1 - \frac{6 \left[ 72 + \frac{6}{12} + \frac{24}{12} + \frac{6}{12} \right]}{990}$$

$$1 - \frac{6 \left[ 75 \right]}{990}$$

$$1 - \frac{6(75)}{990} = 1 - \frac{6(75)}{990}$$

$$1 - 0.43$$

$$= 0.57$$

Ans = 0.57

Problem : The ranks of 15 students in two subjects A and B are given below. The two numbers within brackets denote the ranks of a student in A and B.

(1, 10) (2, 7) (3, 2) (4, 6) (5, 4) (6, 8) (7, 3) (8, 1)  
 (9, 11) (10, 5) (11, 9) (12, 5) (13, 14) (14, 12) (15, 13)

Find Spearman's rank correlation coefficient.

Rank A ( $R_1$ )	Rank (B) $R_2$	$d = R_1 - R_2$	$d^2$
1	10	-9	81
2	7	-5	25
3	2	6	1
4	6	5	25
5	4	1	1
6	8	-2	4
7	3	4	16
8	1	3	9
9	11	-2	4
10	5	-5	25
11	9	-2	4
12	12	0	0
13	14	-1	1
14	1	2	4
15	13	2	4

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

$$\sum d^2 = 272$$

$$1 - \frac{6 \times 272}{15(225-1)} = 1 - \frac{6 \times 272}{15(224)} = 1 - \frac{1632}{3360} = 1 - 0.48 \\ R = 0.52$$

Judge 1 6 5 10 324 9 7 8.

1

Judge 3 5 8 4 7102 1691.

2

Judge 6 4 9 8 123. 1057.

3

Use the rank correlation coefficient to determine which pair of judges has the nearest approach for judgement of beauties.

	Judge (R <sub>1</sub> )	Judge (R <sub>2</sub> )	Judge (R <sub>3</sub> )
d <sub>1</sub>	25 5 6 5 5 0 - - 1 - 1 - 1	100 - 100 - 100 - 100 - 100 - 100	
d <sub>2</sub>	5 2 4 2 2 0 1 1 2	100 - 100 - 100 - 100 - 100 - 100	
d <sub>1</sub> <sup>2</sup> + d <sub>2</sub> <sup>2</sup>	0 - - 16 36 64 - 5 (1-1) / 12	100 - 100 - 100 - 100 - 100 - 100	
d <sub>1</sub> <sup>2</sup> + R <sub>2</sub> - R <sub>1</sub>	3 - 4 - 5 - 6 - 7 - 8 - 9 - 10	100 - 100 - 100 - 100 - 100 - 100	
d <sub>1</sub> <sup>2</sup> - R <sub>2</sub> + R <sub>1</sub>	5 - 4 - 3 - 2 - 1 - 0 - 1 - 2	100 - 100 - 100 - 100 - 100 - 100	
Judge 2. d <sub>1</sub> , R <sub>1</sub> - R <sub>2</sub>	-2 - 3 - 6 - 7 - 8 - 9 - 10 - 11	100 - 100 - 100 - 100 - 100 - 100	
Judge 3. d <sub>1</sub> , R <sub>1</sub> - R <sub>2</sub>	6 4 9 8 - 2 3 10 5 7	100 - 100 - 100 - 100 - 100 - 100	
Judge 2	3 5 8 4 7 10 2 - 6 9	100 - 100 - 100 - 100 - 100 - 100	
Judge 1	- 6 5 10 9 2 5 0, 1 8	100 - 100 - 100 - 100 - 100 - 100	

$$R_{12} = 1 - \frac{6 \sum d_1^2}{n(n^2-1)}$$

$$R_{12} = 1 - \frac{6(100)}{10(99)}$$

$$R_{12} = 0.21.$$

$$R_{23} = 1 - \frac{6 \sum d_2^2}{n(n^2-1)}$$

$$1 - \frac{6 \times 214}{990}$$

$$R_{23} = -0.29.$$

$$R_{31} = 1 - \frac{6 \sum d_3^2}{n(n^2-1)}$$

$$1 - \frac{6(60)}{990} = R_{31} = 0.64.$$

∴ The correlation coefficient  $R_{31} = 0.64$  is highest, the judges 3 and 1 have nearest approach for judgement of beauties.

Find the RCC from the following data

marks in  
calculus  
marks in  
statistics

15 20 28 12 40 60 80 80.

40 30 50 30 20 10 30 60,

marks in  
calculus.

	marks in calculus	$R_1$	$R_2$	$d = R_1 - R_2$	$d^2$
15	40	7	36	4	16
20	30	5.5	15	0.5	0.25
28	50	4	2	2	4
12	30	8	5	3	9
40	20	3	7	-4	16
60	10	2	8	-6	36
20	30	5.5	15	0.5	0.25
80.	60.	1	1	0.	0.

$$\sum d^2 = 81.5$$

$$m_1 = 2$$

$$m_2 = 3$$

$$R_2 = 1 - \frac{6}{12} \left[ \sum d^2 + \frac{1}{12} (m_1^3 - m_1) + \frac{1}{12} (m_2^3 - m_2) \right]$$

$$= 1 - \frac{6}{12} \left[ 81.5 + \frac{1}{12} (8 - 2) + \frac{1}{12} (27 - 3) \right]$$

$$= 1 - \frac{6}{12} \left[ 81.5 + 0.5 + 2 \right] = 1 - \frac{6}{12} \left[ 84 \right] = 1 - \frac{504}{8(63)}$$

$$= 1 - 1 = 0$$

16) calculate RCC after making adjustment at tied rank.

x	y	R <sub>1</sub>	R <sub>2</sub>	d	d <sup>2</sup>
18	13	3	5.5	-2.5	6.25
33	13	5	5.5	-0.5	0.25
40	24	4	1	3	9
9	6	10	8.5	1.5	2.25
16	15	8	10	4	16
16	4	8	10	-2	4
65	20	1	2	-1	1
24	9	6	7	-1	1
16	6	8	8.5	-0.5	0.25
57	14	2	3	-1	1
					<u>41</u>

$$m_1 = 3$$

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

$$m_2 = 2$$

$$\frac{5+6}{2} = \frac{11}{2} = 5.5$$

$$m_3 = 2$$

$$\frac{8+9}{2} = \frac{17}{2} = 8.5$$

$$R = 1 - \frac{6}{n(n^2-1)} \left[ 3d^2 + \frac{1}{12}(m_3^3 - m_1) + \frac{1}{12}(m_2^3 - m_2) + \frac{1}{12}(m_1^3 - m_3) \right]$$

$$= 1 - \frac{6}{990} \left[ 41 + \frac{24}{12} + \frac{6}{12} + \frac{6}{12} \right]$$

$$= 1 - \frac{6}{990} [41 + 2 + 0.5 + 0.5] = 1 - \frac{6}{990} [44]$$

$$1 - \frac{264}{990} = 1 - 0.26 \\ = 0.74$$

## Regression Analysis

Regression is the measure of the average relationship between two or more variable in terms of the original units of the data. It is a statistical tools with help of which the unknown values of one variable can be estimated from known values of another variable.

Two kinds of regression may be studied on the basis of :

Changes in proportions:-

Linear regression: Regression is said to be linear when one variable move with other variable in fixed proportion.

Non linear regression: Regression is said to be non linear when one variable move with other variable in changing proportions.

Number of variation:

Simple regression: When only two variables are studied it is a single regression.

$$\text{eg: } Y = a + bx + \epsilon.$$

Where  $Y$  = Dependent variable.

$x$  = Independent (explanatory) variable.

$a$  = Intercept.

$b$  = Slope.

$\epsilon$  = Error term.

Multiple regression: When at least three variables are studied and their relationship are simultaneously worked out, it is a case of multiple regression.

$$\text{equation: } y = a + b_1x_1 + c_2x_2 + d_3x_3 + \epsilon$$

where,  $y$  = Dependent variable.

$x_1, x_2, x_3$  = Independent (explanatory) variable.

$a$  = Intercept.

$b, c, d$  = Slope.

$\epsilon$  = Error term.

Partial regression:

When more than two variables are studied, keeping the other variables constant, it is called Partial regression.

The methods of least squares:

$$y = a + bx$$

where  $y$  = Dependent variable

$x$  = Independent (explanatory) variable.

$a$  = Intercept.

$b$  = Slope.

Regression equation of  $y$  on  $x$

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

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

Regression equation  $x$  on  $y$ :  $n - \pi = b_{xy}x + a_{xy}$   
 $a_{xy} = \bar{y} - b_{xy}\bar{x}$

Short cut method:

$$S_{xy} = S_{x'y'} - \frac{\sum xy}{n}$$

$$S_{xx} = \sum x^2 - \frac{(\sum x)^2}{n}$$

$$b = \frac{S_{xy}}{S_{xx}}$$

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

The regression equation is  $y = a + bx$ .

Deviations taken from assumed mean values  
of  $x$  and  $y$ :

Regression equation  $x$  and  $\bar{y}$ :  $x - \bar{x} = b_{xy}(y - \bar{y})$

$$\text{where } b_{xy} = \frac{n \sum d_x d_y - (\sum d_x)(\sum d_y)}{n \sum d_y^2 - (\sum d_y)^2}$$

$n$  = number of observations.

$d_x = x - A$ ,  $A$  is assumed mean of  $x$ .

$d_y = y - B$ ,  $B$  is assumed mean of  $y$ .

Regression equation of  $x$  on  $y$ :  $y - \bar{y} = b_{yx}(x - \bar{x})$ .

$$\text{where } b_{yx} = \frac{n \sum d_x d_y - (\sum d_x)(\sum d_y)}{n \sum d_x^2 - (\sum d_x)^2}$$

$n$  = number of observations of  $x$ .

$d_x = x - A$ ,  $A$  is assumed mean of  $x$ .

$d_y = y - B$ ,  $B$  is assumed mean of  $y$ .

Deviation taken from actual mean values of  $x$  and  $y$ .

Regression equation  $x$  and  $y$ :  $n - \pi = b_{xy}$

where  $b_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (y - \bar{y})^2}$ .

Regression equation  $y$  on  $x$ :  $y - \bar{y} = b_{xy}(x - \bar{x})$

where  $b_{xy} = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2}$ .

Regression coefficient in term of correlation coefficient:

Regression equation  $x$  on  $y$ :  $x - \bar{x} = r \frac{\sigma_x}{\sigma_y} (y - \bar{y})$

where  $\bar{x}$  = mean of  $x$ ;  $\bar{y}$  = mean of  $y$ .  
r = correlation between the  $x$  and  $y$  variables.

Regression equation  $y$  on  $x$ :  $y - \bar{y} = r \frac{\sigma_y}{\sigma_x} (x - \bar{x})$ .

The correlation coefficient is the geometric mean of two regression coefficients.  $r = \sqrt{b_{xy} \times b_{yx}}$ .

from the following data find two regression equations.

x	1	2	3	4	5	6	7	8	9
y	2	4	7	6	5	6	5	4	3

Sol

x	y	$x^2$	$y^2$	$xy$
1	2	1	4	2
2	4	4	16	8
3	7	9	49	21
4	6	16	36	24
5	5	25	25	25
6	6	36	36	36
7	5	49	25	35
$\sum x = 28$	$\sum y = 35$	$\sum x^2 = 140$	$\sum y^2 = 191$	$\sum xy = 151$

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

$$n \sum y^2 - (\sum y)^2$$

$$\frac{\sum (151) - (28)(35)}{\sum (191) - (35)^2}$$

$$\frac{1057 - 980}{1337 - 1225} = \frac{77}{112} = 0.69 = b_{xy}$$

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

$$\bar{y} = \frac{\sum y}{n} = \frac{35}{7} = 5$$

Regression equation  $x$  on  $y$

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

$$x - 4 = 0.69(y - 5)$$

$$x - 4 = 0.69y - 3.45$$

$$x = 0.69y - 3.45 + 4$$

$$x = 0.69y + 0.55$$

$$\therefore b_{xy} = 0.39.$$

Regression equation  $y$  on  $x$ .

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

$$y - 5 = 0.39(x - 4)$$

$$y - 5 = 0.39x - 1.56$$

$$y = 0.39x - 1.56 + 5$$

$$y = 0.39x + 3.44.$$

- 2) From the following data find two regression equations

$x \quad 3 \quad 4 \quad 8 \quad 7 \quad 2$

$y \quad 11 \quad 12 \quad 9 \quad 3 \quad 5$

$x$	$y$	$x^2$	$y^2$	$xy$
3	11	9	121	33
4	12	16	144	48
8	9	64	81	72
7	3	49	9	21
2	5	4	25	10
$\sum x = 24$	$\sum y = 40$	$\sum x^2 = 142$	$\sum y^2 = 380$	$\sum xy = 184$

$$b_{xy} = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum y^2 - (\sum y)^2}$$

$$\frac{s(184) - (24)(40)}{s(142) - (24)^2} = \frac{920 - 960}{710 - 576}$$

$$\frac{-40}{300} = -0.13$$

$$b_{xy} = -0.13$$

$$\bar{x} = \frac{\sum x}{n} = \frac{24}{5} = 4.8$$

$$\bar{y} = \frac{\sum y}{n} = \frac{40}{5} = 8$$

Regression eq  $x$  and  $y$

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

$$x - 4.8 = -0.13(y - 8)$$

$$x - 4.8 = -0.13y + 1.04$$

$$x = -0.13y + 1.04 + 4.8$$

$$x = -0.13y + 5.84$$

$$b_{yx} = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum x^2 - (\sum x)^2}$$

$$\frac{s(184) - (24)(40)}{s(142) - (24)^2} = \frac{920 - 960}{710 - 576}$$

$$\frac{-40}{134} = -0.29$$

$$b_{yx} = -0.29$$

$$\Rightarrow \frac{-40}{134}$$

$$\Rightarrow -0.29$$

Regression equation  $y$  on  $x$

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

$$y - 8 = -0.29(x - 4.8)$$

$$y - 8 = -0.29x + 1.39$$

$$y = -0.29x + 9.39$$

Problem: Price indices of cotton and wool are given below for twelve months of a year. Obtain the equation of the line of regression between the indices.

Price index of cotton 78 77 85 88 87 82 81 77 76 83

Price index of wool 84 82 82 85 89 90 88 92 83 89 91

$x$	$y$	$d_x = x_i - A$	$d_x^2$	$d_y = y_i - B$	$d_y^2$	$dx dy$
78	84	-5	25	-4	16	-20
77	82	-6	36	-6	36	36
85	82	2	4	-6	36	-12
88	85	5	25	-3	9	-15
87	89	4	16	1	1	4
82	90	-1	1	2	4	-2
81	88-B	-2	4	0	0	0
77	92	-6	36	4	16	-24
76	83	-7	49	-5	25	35
83-A	89	0	0	1	1	0
97-A	98	14	196	10	100	140
93	99-B	10	100	11	121	110
<hr/>	<hr/>	<hr/>	<hr/>	<hr/>	<hr/>	<hr/>
1004	1061	8	492	5	365	292

$$b_{yx} = \frac{n \sum d_x d_y - (\sum d_x)(\sum d_y)}{n \sum d_x^2 - (\sum d_x)^2}$$

$$\frac{12(292) - (8)(5)}{12(492) - (8)^2}$$

$$\frac{3504 - 40}{5904 - 631}$$

$$\frac{3464}{5840} = 0.59.$$

$$b_{yx} = 0.59.$$

Regression equation  $y$  on  $x$ .

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

$$y - 88.4 = 0.59(x - 83.6)$$

$$y = 0.59x - 49.3 + 88.4$$

$$y = 0.59x + 39.1$$

$$b_{xy} = \frac{n \sum d_x d_y - (\sum d_x)(\sum d_y)}{n \sum d_y^2 - (\sum d_y)^2}$$

$$= \frac{12(292) - (8)(5)}{12(365) - (5)^2}$$

$$\frac{3504 - 40}{4380 - 250}$$

$$\frac{3464}{4355} = 0.79.$$

$$b_{xy} = 0.79.$$

Regression equation  $x$  on  $y$ .

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

$$x - 83.6 = 0.79(y - 88.4) \Rightarrow x = 0.79y - 69.8 + 83.6$$

$$x = 0.79y + 13.8$$

Problem: The following data relate to the scores obtain by a salesman of a company in an intelligence test and their weekly sales (Rs. in 1000).

Sales man	A	B	C	D	E	F	G	H	I	J
-----------	---	---	---	---	---	---	---	---	---	---

Test scores	50	60	50	60	80	50	80	40	70	
-------------	----	----	----	----	----	----	----	----	----	--

Weekly sales	30	60	40	50	60	30	70	50	60	
--------------	----	----	----	----	----	----	----	----	----	--

x	y	$\alpha_x = x_i - A$	$\alpha_x^2$	$\alpha_y = y_i - A$	$\alpha_y^2$	$\alpha_x \alpha_y$
50	30	-10	100	-20	400	200
60	60	0	0	10	100	0
50	40	-10	100	-10	100	100
60	50	0	0	0	0	0
80	60	20	400	10	100	200
50	30	-10	100	-20	200	200
80	70	20	400	20	400	400
40	50	-20	400	0	0	0
70	60	10	100	10	100	100
<hr/>	<hr/>	<hr/>	<hr/>	<hr/>	<hr/>	<hr/>
540.	450.	0	1600	0	1600	1200

$$b_{yx} = \frac{n \sum dx dy - \sum dx \cdot \sum dy}{n \sum dx^2 - (\sum dx)^2}$$

$$\frac{9(1200) - 0}{9(1600)} = \frac{9(1200)}{9(1600)} \cdot \frac{3}{4} = 0.75$$

① Regression equation  $y = a + b_1 x$

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

$$y - 50 = 0.75(x - 60)$$

$$y - 50 = 0.75x - 45$$

$$y = 0.75x - 45 + 50$$

$$y = 0.75x + 5$$

when  $x = 65$

$$y = 0.75(65) + 5$$

$$48.75 + 5$$

$$y = 53.75$$

Prob: The following data based on 450 students are given for marks in statistics and Economics.

at a certain exam

Mean marks in statistics: 40

Mean marks in Economics: 48

S.D of marks in Economics: 2.56

Sum of the product of deviation of marks from their respective mean: 48075

Obtain equation of the two lines of regression, and estimate the average marks in Economics of candidates who obtained 50 marks in statistics.

Sol) Given  $n = 450$   
let the marks in statistics and economics be  $x, y$   
 $\bar{x} = 40, \bar{y} = 48$   
 $\sigma_x = 12, \sigma_y = 256$   
 $\sigma_y = \sqrt{256} = 16$

$$\text{Edxdy} = 42075$$

$$r = \frac{42075}{450 \times 16 \times 12} = 0.45$$

Regression equation  $y$  on  $x$ :

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

$$y - 48 = 0.48 \cdot \left( \frac{16}{12} \right) (x - 40)$$

$$y - 48 = 0.48 \left( \frac{4}{3} \right) (x - 40)$$

$$y - 48 = 0.48 (1.3) (x - 40)$$

$$y - 48 = 0.624 (x - 40)$$

$$y = 48 + 0.624x - 24.96$$

$$y = 0.624(50) - 24.96 + 48$$

$$y = 31.2 + 23.04$$

$$y = 581.06$$