

## Chapter: Correlation -6

md. Abdul Aleem  
2-11-2010

**Bivariate data:** Bivariate data is data for which there are two variables for each observation. As an example, the following bivariate data show the ages of husbands and wives of 10 married couples.

Husband	36	72	37	36	51	50	47	50	37	41
Wife	35	67	33	35	50	46	47	42	36	41

**Correlation:** Correlation is a statistical technique which measure and analyses the degree or extent to which two or more variables fluctuate with reference to one another.

Correlation thus denotes the interdependence amongst variates. The degrees are expressed by a coefficient which ranges between -1 and 1. The direction of change is indicated by + or - signs.

If the **increase (decrease)** in one variable results in the corresponding **increase (decrease)** in the others i.e. if the changes are in the same directions the variables are **positively correlated**. For example, the heights and weights of a group of persons are positively correlated, advertising and sales.

If the **increase (decrease)** in one variable results in the corresponding **decrease (increase)** in the others i.e. if the changes are in the opposite directions the variables are **negatively correlated**. For example, T.V registration and cinema attendance is negatively correlated.

An absence of correlation is indicated by zero.

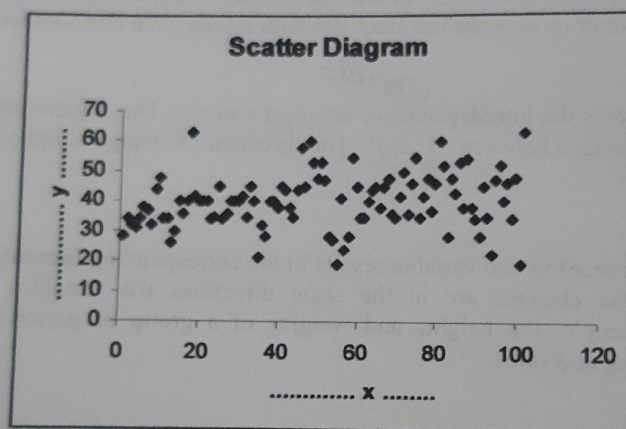
Correlation thus expresses the relationship through a relative measure of change and it has nothing to do with the units in which the variables are expressed.

### Uses

- Economic theory and business studies relationships between variables like price and quantity demanded, advertising, expenditure scales promotion measure etc. The correlation analysis helps in deriving precisely the degree and direction of such relationships.
- The concepts of regression are also based upon the measure of correlation.

✓ **Scatter Diagram:** Scatter diagram (or Dotogram or Scattergram) is a simple and attractive method of diagrammatic represent of bivariate distribution for ascertaining the nature of correlation between the variables. Thus for the bivariate distribution  $(x_i, y_i)$ ,  $i = 1, 2, \dots, n$  if the values of the variables  $X$  and  $Y$  be plotted along the  $X$ -axis and  $Y$ -axis respectively in the  $XY$  plane, the diagram of dots so obtained is known as scatter diagram.

✓ On the other hand, a scatter plot of two variables shows the values of one variable on the  $Y$ -axis and the values of the other variable on the  $X$ -axis. Scatter plots are well suited for revealing the relationship between two variables.



**Types of Correlation:** Correlation is described or classified in several different ways. Three of the most important are:

- ✓ ❖ Positive and negative Correlation ✓
- ✓ ❖ Simple, partial and multiple Correlation ✓
- ✓ ❖ Linear and non-linear Correlation ✓

✓ **Positive and negative correlation:** If two variables changes in the same direction (i.e. if one increases the other also increase or if one decreases the other also decreases) then this is called a positive correlation. For example:



Positive Correlation	
X	Y
10	15
12	20
14	22
18	25
20	37

Positive Correlation	
X	Y
80	50
70	45
60	30
40	20
30	10

If two variables change in the opposite direction (i.e. if one increases, the other decreases and vice versa), then the correlation is called a negative correlation. For example: T.V registrations and cinema attendance.

Negative Correlation	
X	Y
20	40
30	30
40	22
60	15
80	12

Negative Correlation	
X	Y
100	10
90	20
60	30
40	40
30	50

## 2. Simple, Partial and Multiple Correlation

- ◆ When only two variables are studied it is a problem of simple correlation.
- ◆ When three or more variables are studied it is a problem of either multiple or partial correlation.

In multiple correlation three or more variables are studied simultaneously. For example, when we study the relationship between the yield of rice per acre and both the amount of rainfall and the amount of fertilizers used, it is problem of multiple correlation. Similarly the relationship of plastic hardness, temperature and pressure is multivariate.

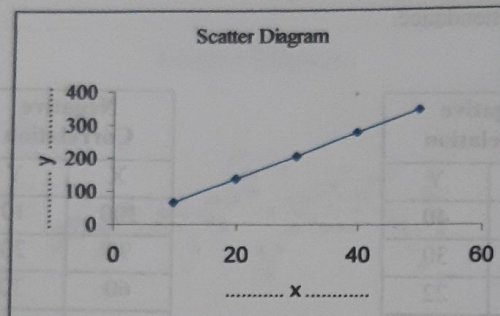
In partial correlation we recognize more than two variables. But consider only two variables to be influencing variable being kept constant. For example, in the rice problem taken above if we limit our correlation analysis of yield and rainfall to periods when a certain average daily temperature existed, it becomes a problem of partial correlation.

**3. Linear and non-linear correlation:** The nature of the graph gives us the idea of the linear type of correlation between two variables. If the graph is in a straight line, the correlation is called a "linear correlation" and if the graph is not in a straight line, the correlation is non-linear and curve-linear.

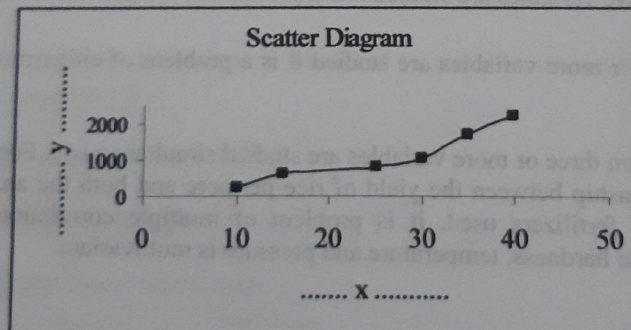
The distinction between linear and non-linear correlation is based upon the constancy of the ratio of change between the variables. If the amount of change in one variable tends to bear a constant ratio to the amount of change in the other variable then the correlation is said to be linear. For example, observe the following two variables X and Y:

X:	10	20	30	40	50
Y:	70	140	210	280	350

It is clear that the ratio of change between the two variables is the same. If such variables are plotted on a graph paper all the plotted points would fall on a straight line.



Correlation would be called non-linear or curvilinear if the amount of change in one variable doesn't bear a constant ratio to the amount of change in the other variable. For example, if we double the amount of rainfall, the production of rice or wheat etc. would not necessarily be doubled.



**Properties of the Coefficient of Correlation:** The following are the important properties of the coefficient of correlation,  $r$ :

- The coefficient of correlation lies between -1 and +1,  $-1 \leq r \leq +1$ .



- ✓ The coefficient of correlation is the geometric mean of the two regression coefficients. Symbolically:  $r = \sqrt{b_{xy} \times b_{yx}}$
- ✓ If X and Y are independent variables then coefficient of correlation is zero. However, the converse is not true.

**Degrees of Correlation:** Through the coefficient of correlation, we can measure the degree or extent of the correlation between two variables. On the basis of the coefficient of correlation we can also determine whether the correlation is positive or negative and also its degree or extent.

- ✓ **Perfect correlation:** If two variables changes in the same direction and in the same proportion, the correlation between the two is **perfect positive**. According to Karl Pearson the coefficient of correlation in this case is +1. On the other hand, if the variables change in the opposite direction and in the same proportion, the correlation is **perfect negative**. Its coefficient of correlation is -1. In practice we rarely come across these types of correlations.
- ✓ **Absence of correlation:** If two series of two variables exhibit no relations between them or change in variable does not lead to a change in the other variable, then we can firmly say that there is **no correlation** or **absurd correlation** between the two variables. In such a case the coefficient of correlation is 0.
- ✓ **Limited degrees of correlation:** If two variables are not perfectly correlated or is there a perfect absence of correlation, then we term the correlation as Limited correlation. It may be positive, negative or zero but lies with the limits  $\pm 1$ .

High degree, moderate degree or low degrees are the three categories of this kind of correlation. The following table reveals the effect (or degree) of coefficient or correlation.

Degrees	Positive	Negative
Absence of correlation →	Zero	0
Perfect correlation →	+ 1	-1
High degree →	+ 0.75 to + 1	- 0.75 to -1
Moderate degree →	+ 0.25 to + 0.75	- 0.25 to - 0.75
Low degree →	0 to 0.25	0 to - 0.25

**Methods of Determining Correlation:** We shall consider the following most commonly used methods.

- ◆ Scatter Plot.
- ◆ Karl Pearson's coefficient of correlation.
- ◆ Spearman's Rank-correlation coefficient.
- ◆ Method of Least Squares.

**Karl Pearson's Coefficient of Correlation:** Of the several mathematical methods of measuring correlation, the Karl Pearson's method, popularly known as Pearsonian coefficient of correlation, is most widely used in practice. The coefficient of correlation is denoted by the symbol  $r$ . If the two variables under study are  $X$  and  $Y$ , the following formula suggested by Karl Pearson can be used for measuring the degree of relationship.

$$r = \frac{\sum XY - \frac{\sum X \sum Y}{N}}{\sqrt{\left\{ \sum X^2 - \frac{(\sum X)^2}{N} \right\} \left\{ \sum Y^2 - \frac{(\sum Y)^2}{N} \right\}}}$$

The value of the coefficient of correlation as obtained by the above formula shall always lie between  $\pm 1$ .

When  $r = +1$ , it means there is perfect positive correlation between the variables.

When  $r = -1$ , it means there is a perfect negative correlation between the variables.

When  $r = 0$ , it means there is no relationship between the variables.

**Example 1:** Calculate the coefficient of correlation between the heights of father and his son for the following data.

Height of father (cm):	165	166	167	168	167	169	170	172
Height of son (cm):	167	168	165	172	168	172	169	171

**Solution:** We know that. Correlation of coefficient



$$r = \frac{\sum XY - \frac{\sum X \sum Y}{N}}{\sqrt{\left[ \sum X^2 - \frac{(\sum X)^2}{N} \right] \left[ \sum Y^2 - \frac{(\sum Y)^2}{N} \right]}}$$

Let us consider the height of father is  $X$  and height of son is  $Y$ .

By using calculator we get,

$$\sum X^2 = 225828 \quad \sum X = 1344 \quad N = 8$$

$$\sum Y^2 = 228532 \quad \sum Y = 1352 \quad \sum XY = 227160$$

$$r = \frac{227160 - \frac{1344 \times 1352}{8}}{\sqrt{\left[ 225828 - \frac{(1344)^2}{8} \right] \left[ 228532 - \frac{(1352)^2}{8} \right]}}$$

$$= 0.603022689 = 0.603$$

So the relation between the height of father and height of son is moderate degree positive.

**Example 2:** The following data consist of observations for the weights of 10 different automobiles (in 1000 pounds) and the corresponding fuel consumptions (gallons per 100 miles).

Weight (x)	Fuel Consumption (y)
3.4	5.5
3.8	5.9
4.1	6.5
2.2	3.3
2.6	3.6
2.9	4.6
2.0	2.9
2.7	3.6
1.9	3.1
3.4	4.9

We would like to find out how  $y$  is correlated to  $x$ .

**Solution:** We know that. Correlation of coefficient

$$r = \frac{\sum XY - \frac{\sum X \sum Y}{N}}{\sqrt{\left[ \sum X^2 - \frac{(\sum X)^2}{N} \right] \left[ \sum Y^2 - \frac{(\sum Y)^2}{N} \right]}}$$

By using calculator we get,

$$\sum X^2 = 89.29 \quad \sum X = 29 \quad N = 10$$

$$\sum Y^2 = 207.31 \quad \sum Y = 43.9 \quad \sum XY = 135.8$$

So the relation between weight and Fuel Consumption is high degree positive.

equ  
res/  
n ec  
n e  
2 +  
1 e

e.  
ac  
d

1

1

5

103

103

5

- |  |  |
|--|--|
|  |  |
|--|--|

- 5

5

- 1011

--	--

103



Employee	Ranked by manager I	Ranked by Manager II
A	10	9
B	2	4
C	1	2
D	4	3
E	3	1
F	6	5
G	5	6
H	8	8
I	7	7
J	9	10

Compute the coefficient of rank correlation and comment on the value.

**Solution: Calculation of Rank Correlation Coefficient**

Employee	Ranked by manager I ( $R_1$ )	Ranked by Manager II ( $R_2$ )	$d^2 = (R_1 - R_2)^2$
A	10	9	
B	2	4	
C	1	2	
D	4	3	
E	3	1	
F	6	5	
G	5	6	
H	8	8	
I	7	7	
J	9	10	
Total			$\sum d_i^2 = 14$

We know that,  $\rho = 1 - \frac{6 \sum d_i^2}{n^3 - n} = 1 - \frac{6 \times 14}{10^3 - 10} = 0.915$

Thus we find that there is a high degree of positive correlation in the ranks assigned by the two managers.

**B. Where Ranks are not given:** When we are given the actual data and not the ranks it will be necessary to assign the ranks. Ranks can be assigned by taking either the highest value as 1 or the lowest value as 1. But whether we start with the lowest value or the highest value we must follow the same method in case of all the variables.

✓

2  
mov  
iable  
the  
ex  
pri

**Solut**

gre  
Y

We know that,  $\rho = 1 - \frac{6\sum d_i^2}{n^3 - n} = 1 - \frac{6 \times 44}{10^3 - 10} = 1 - 0.267 = 0.733$

✓

### Merits

- ✓

### Limitations:

-



- Where the number of observations exceed 30 the calculations becomes quite tedious and require a lot of time. Therefore this method should not applied where  $n$  exceeding 30 unless we are given the ranks and not the actual values of the variable.

#### (4) Method of Least Squares

For finding out correlation by the coefficient method of least squares we have to calculate the values of two regression coefficients that of  $x$  on  $y$  and  $y$  on  $x$ . The correlation coefficient is the square root of the product of two regression coefficients. Symbolically,

$$r = \sqrt{b_{xy} \times b_{yx}}$$

**Coefficient of Determination:** One very convenient and useful way of interpreting the value of coefficient of correlation between two variables is to use the square of coefficient of correlation, which is called coefficient of determination. The coefficient of determination thus equals  $r^2$ .

\*\*\* If the value of  $r = 0.9$ ,  $r^2$  will be 0.81 and this would mean that 81% of the variation in the dependent variable has been explained by the independent variable.