

Department of Mathematics University of Ruhuna

A.W.L. Pubudu Thilan

Chapter 6

Joint distribution of data

- Up to now we have considered distribution of only one variable.
- But in practice we meet problems with two or more variables.
- If we have problem of two variables, dependence between two variables is important.

Concept of correlation

- In statistics, dependence refers to any statistical relationship between two random variables or two sets of data.
- A correlation is a single number that describes the degree of dependence between two variables.
- **Eg:** Yield depends on the amount of fertilizer used.
- **Eg:** Sales of a product depends on price charged.







- Possible correlations range from +1 to -1.
- A zero correlation indicates that there is no **linear** relationship between the variables.
- A correlation of -1 indicates a **perfect negative** correlation.
- A correlation of +1 indicates a **perfect positive** correlation.

Scatter plot

- A scatter plot is a type of mathematical diagram using Cartesian coordinates to display values for two variables for a set of data.
- The data are displayed as a collection of points, each having the value of one variable determining the position on the horizontal axis and the value of the other variable determining the position on the vertical axis.



The correlation is described according to the following ways,

- 1 Positive or negative.
- 2 Simple, partial or multiple.
- 3 Linear or non linear.

- If an increase (or decrease) of values of one variable is associated with an increase (or decrease) in the corresponding values of the other variable, they are said to be correlated and the correlation between these two variables is said to be **positive** or **direct**.
- **Eg:** People who do more revision get higher exam results.



1.2 Negative correlation

- If an increase (or decrease) of values of one variable is associated with an decrease (or increase) in the corresponding values of the other variable, they are said to be correlated and the correlation between these two variables is said to be negative or inverse.
- Eg: There is a negative correlation between TV viewing and class grades. Students who spend more time watching TV tend to have lower grades.



2.1 Simple correlation

- If we are studying only two variables the correlation between them is a simple correlation.
- In simple correlation, we measure the strength of the linear relationship between two variables, without taking into consideration the fact that both these variables may be influenced by a third variable.



2.2 Partial correlation

- Partial correlation analysis involves studying the linear relationship between two variables after excluding (or held constant) the effect of one or more independent factors.
- The correlation between X and Y, with the effects of Z removed (or held constant) is called the partial correlation of X and Y.



2.3 Multiple correlation

In **multiple correlation**, we study the effects of all the independent variables simultaneously on a dependent variable.



3.1 Linear correlation

- Correlation is said to be linear if the ratio of change is constant.
- The amount of output in a factory is doubled by doubling the number of workers is the example of linear correlation.
- If all the points on the scatter diagram tends to lie near a line which are look like a straight line, the correlation is said to be linear.



3.1 Non linear correlation

- Correlation is said to be non linear if the ratio of change is not constant.
- If all the points on the scatter diagram tends to lie near a smooth curve, the correlation is said to be non linear (curvilinear).



- If the points form a straight line it indicates a perfect correlation.
- If the points form a band of some width it indicates imperfect correlation between two variables.
- The direction of the band shows the nature of the correlation (positive or negative).

- The width of the band gives an idea of the degree of correlation.
- The narrower the band the greater is the degree of correlation.
- When the points are scattered in all directions, it indicates that there is no correlation between the variables (i.e variables are uncorrelated).

If two variables are independent \Rightarrow They are uncorrelated. Uncorrelated \Rightarrow Independent.

Mathematical methods of measuring correlation

- The degree or level of correlation is measured with the help of correlation coefficient.
- For population data, the population correlation coefficient is defined and it is denoted by ρ.
- For **sample data**, the sample correlation coefficient is defined and it is denoted by *r*.

- The joint variation of X and Y is measured by the **population covariance** of X and Y.
- The population covariance of X and Y denoted by Cov(X, Y) is defined as:

$$Cov(X, Y) = rac{\sum_{i=1}^{N} (x_i - \mu_x)(y_i - \mu_y)}{N}.$$

- The Cov(X, Y) may be positive, negative or zero.
- The covariance has the same units in which X and Y are measured.
- When Cov(X, Y) is divided by σ_X and σ_Y, we get the population correlation coefficient ρ,

$$\rho = \frac{Cov(X, Y)}{\sigma_X \sigma_Y}.$$

- ρ is free of the units of measurement.
- It is a pure number and lies between -1 and +1.
- If $\rho = \pm 1$, it is called perfect correlation.
- If X and Y are independent, then the correlation coefficient of X and Y is equal to zero.

- The joint variation of X and Y is measured by the **sample covariance** of X and Y.
- The sample covariance of X and Y denoted by Cov(X, Y) is defined as:

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

■ When *Cov*(*X*, *Y*) is divided by *S*_{*X*} and *S*_{*Y*}, we get the sample correlation coefficient *r*,

$$r=\frac{Cov(X,Y)}{S_XS_Y}.$$

- r is free of the units of measurement.
- It is a measure of strength of the linear relation between X and Y variables.

 On the other hand it is called as Karl Pearsons coefficient of correlation or product-moment correlation coefficient and denoted by,

$$r = \frac{\sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \overline{x})^2 \sum_{i=1}^{n} (y_i - \overline{y})^2}} \text{ for } n \text{ pairs } (x_i, y_i).$$

Sample correlation coefficient(r)Formula 2

It can also be written as,

$$r = \frac{\sum_{i=1}^{n} x_i y_i - \frac{\sum_{i=1}^{n} x_i \sum_{i=1}^{n} y_i}{n}}{\sqrt{(\sum_{i=1}^{n} x_i^2 - n\overline{x}^2)(\sum_{i=1}^{n} y_i^2 - n\overline{y}^2)}}.$$

- r is a pure number and lies between -1 and +1.
- The sign of *r* determines the nature of correlation.
- If *r* is positive it indicates a positive correlation.
- If *r* is negative it indicates a negative correlation.
- The magnitude of *r* determines the degree of correlation.

Sample correlation coefficient(r) Properties of $r \Rightarrow Cont...$

- If r = +1, it indicates a positive perfect correlation.
- If r = -1, it indicates a negative perfect correlation.
- If r = 0, the data are uncorrelated.
- If -1 < r < 0, it indicates negative imperfect correlation.
- If r < -0.75, it indicates a strong negative correlation.
- If 0 < *r* < 1, it indicates positive imperfect correlation.
- If r > 0.75, it indicates a strong positive correlation.

Scatter plots for different degree of correlations



- (a) represents a perfect correlation of +1, all the points fall on a perfectly straight line with a positive slope.
- (b) represents a strong correlation where the behavior of one variable is similar, but not identical to the behavior of the other variable.
- (c) a correlation of .50 represents a moderately strong positive relationship.
- (d) relationship is weak, so the coefficient is only .25.

- (e) correlation is zero; the x and y variables are not linearly related.
- (f) coefficient is also zero. This is because the correlation coefficient measures a linear association, while the relationship in Figure (f) is curvilinear.

Note:

Figures (g), (h), (i) are mirror images of Figures (a), (b), (c). All the correlation coefficients are negative.

Example

The daily income and the daily expenditure of ten employees of factory are recorded as below.

Daily income	Daily expenditure		
100	98		
101	99		
102	99		
102	97		
100	95		
99	92		
97	95		
98	91		
96	90		
95	91		

Describe the relation using Scatter plot and Karl Pearsons coefficient of correlation.

Example Solution



Example Solution⇒Cont...

Xi	Уi	$(x_i - \overline{x})$	$(y_i - \overline{y})$	$(x_i - \overline{x})(y_i - \overline{y})$	$(x_i - \overline{x})^2$	$(y_i - \overline{y})^2$
100	98	1	3	3	1	9
101	99	2	4	8	4	16
102	99	3	4	12	9	16
102	97	3	2	6	9	4
100	95	1	0	0	1	0
99	92	0	-3	0	0	9
97	95	-2	0	0	4	0
98	91	-1	-4	4	1	16
96	90	-3	-5	15	9	25
95	91	-4	-4	16	16	16
				64	54	111

Example Solution⇒Cont...

$$\overline{x} = \frac{\sum_{i=1}^{10} x_i}{10} = \frac{100 + 101 + 102 + \dots + 95}{10}$$

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

$$\overline{y} = \frac{\sum_{i=1}^{10} y_i}{10} = \frac{98 + 99 + 99 + \dots + 91}{10}$$

$$= 94.7 \simeq 95$$

$$r = \frac{\sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})}{\sqrt{\sum_{i=1}^{n} (x_i - \overline{x})^2 \sum_{i=1}^{n} (y_i - \overline{y})^2}} \text{ for } n \text{ pairs } (x_i, y_i)$$

$$r = \frac{64}{\sqrt{54 \times 111}}$$

$$= 0.83$$

It indicates a positive imperfect correlation.

Rank correlation coefficient

- We used the product-moment correlation coefficient to measure the strength of a linear association between two variables.
- But the product-moment correlation coefficient is less appropriate when the points on a scatter graph seem to follow a curve or when there are outliers on the graph.
- The rank correlation coefficient is appropriate for those data.

• The rank correlation coefficient is defined as:

$$\rho = 1 - \frac{6\sum_{i=1}^{n} d_i^2}{n(n^2 - 1)};$$
 where

$$d_i$$
 = difference of ranks of i^{th} pair of observations

$$n =$$
 number of observations.

Example 1

Use the raw data in the table below to calculate the correlation between the IQ of a person with the number of hours spent in front of TV per week.

$IQ(x_i)$	Hours of TV per week (y_i)
106	7
86	0
100	27
101	50
99	28
103	29
97	20
113	12
112	6
110	17

Example 1 Solution

Xi	Уi	rank x _i	rank y _i	di	d_i^2
86	0	1	1	0	0
97	20	2	6	-4	16
99	28	3	8	-5	25
100	27	4	7	-3	9
101	50	5	10	-5	25
103	29	6	9	-3	9
106	7	7	3	4	16
110	17	8	5	3	9
112	6	9	2	7	49
113	12	10	4	6	36

Example 1 Solution⇒Cont...

By considering the last column, we can find the value of $\sum_{i=1}^{10} d_i^2$ as 194.

$$\rho = 1 - \frac{6\sum_{i=1}^{n} d_i^2}{n(n^2 - 1)}$$
$$= 1 - \frac{6 \times 194}{10(10^2 - 1)}$$
$$= -0.175757575$$

The correlation between IQ and hours spent watching TV is very low.

Example 2

Paintings completed by a group of ten students are assessed by two independent judges. The marks awarded are shown below. Calculate Spearman's rank correlation coefficient for the marks awarded by the two judges.

Completed paintings	First judge	Second judge
A	30	50
В	35	45
С	35	45
D	40	45
E	50	60
F	55	55
G	60	75
Н	65	70
	70	65
J	80	80

Example 2 Solution

The corresponding table to calculate Spearman's rank correlation coefficient is as follows when higest rank is given for lagest values in both variables.

Paintings	Xi	Уi	rank of x _i	rank of <i>x</i> i	di	d_i^2
А	30	50	10	7	3	9
В	35	45	8.5	9	-0.5	0.25
С	35	45	8.5	9	-0.5	0.25
D	40	45	7	9	-2	4
E	50	60	6	5	1	1
F	55	55	5	6	-1	1
G	60	75	4	2	2	4
Н	65	70	3	3	0	0
	70	65	2	4	-2	4
J	80	80	1	1	0	0

Example 2 Solution⇒Cont...

$$\rho = 1 - \frac{6\sum_{i=1}^{n} d_i^2}{n(n^2 - 1)}$$
$$= 1 - \frac{6 \times 23.5}{10(10^2 - 1)}$$
$$= 0.86$$

- (a) Covariance between X and Y is 12.3 and the variance of X and Y are 18.75 and 12 respectively. Find the correlation coefficient between them.
- (b) The following table shows the GPAs for 12 randomly selected students, and the number of tutorial classes that those students missed.

Pass paper 2011 Cont...

GPA	Missed tutorial classes
2.66	3
2.05	1
2.07	2
2.62	0
1.30	7
3.00	0
3.25	2
2.58	0
2.36	3
2.81	1
3.11	1
2.56	2

- (i) Calculate Spearman's rank correlation coefficient.
- (ii) What can you say about the relationship between GPA and number of missed tutorial classes?

Pass paper 2011 Solution

(a) When Cov(X, Y) is divided by σ_X and σ_Y , we get the coefficient of correlation ρ ,

$$\rho = \frac{Cov(X, Y)}{\sigma_X \sigma_Y}, \\ = \frac{12.3}{\sqrt{18.75}\sqrt{12}} \\ = \frac{12.3}{\sqrt{18.75 \times 12}} \\ = \frac{12.3}{\sqrt{18.75 \times 12}} \\ = \frac{12.3}{\sqrt{225}} \\ = \frac{12.3}{\sqrt{15}} \\ = 0.82$$

Pass paper 2011 Solution \Rightarrow Cont...

(b) (i) The corresponding table for Spearman's rank correlation coefficient is given below.

Xi	Уi	rank of x _i	rank of y _i	di	d_i^2
2.66	3	8	10.5	-2.5	6.25
2.05	1	2	5	-3	9
2.07	2	3	8	-5	25
2.62	0	7	2	5	25
1.30	7	1	12	-11	121
3.00	0	10	2	8	64
3.25	2	12	8	4	16
2.58	0	6	2	4	16
2.36	3	4	10.5	-6.5	42.25
2.81	1	9	5	4	16
3.11	1	11	5	6	36
2.56	2	5	8	-3	9

Pass paper 2011 Solution⇒Cont...

$$\rho = 1 - \frac{6\sum_{i=1}^{n} d_i^2}{n(n^2 - 1)}$$
$$= 1 - \frac{6 \times 385.5}{12(12^2 - 1)}$$
$$= -0.3479$$

(ii) There is weak negative correlation between the two variables.

Thank You