

CHAPTER FOUR

CORRELATION

Mr. Mohammad Manjur Alam (Manju)

Associate Professor

Department of Computer Science and Engineering

International Islamic University Chittagong.

Mobile: 01715966144

FB and Email: manjuralam44@yahoo.com

The primary objective of correlation analysis is to measure the strength or degree of relationship between two or more variables. If the change in one variable affects a change in the other variable, the variables are said to be correlated.

For example, the production of paddy is dependent on the rainfall. Here production of paddy is considered to be a dependent variable.

Types of Correlation

- Positive or negative
- Simple or multiple
- Linear or non-linear

Positive or negative

If the two variables deviate in the same direction, that is if the increase (or decrease) in one results in a corresponding increase (or decrease) in the other, correlation is said to be director positive. But if they constantly deviate in the opposite directions, that is if increase (or decrease) in one results in corresponding decrease (or increase) in the other, correlation is said to be inverse or negative. If the variables are independent, there cannot be any correlation and the variables are said to be zero correlation.

For example, the correlation between (1) the heights and weights of a group of persons, (2) the income and expenditure is positive and the correlation between (1) price and demand of a commodity, (2) the volume and pressure of a perfect gas is negative. And there is no correlation between income and height.

Simple correlation and Multiple Correlation

Correlation only between two variables is called simple correlation. For example, correlation between income and expenditure.

Under Multiple Correlation three or more than three variables are studied. Ex. $Q_d = f(P, PC, PS, t, y)$

Linear correlation and Non Linear correlation

Correlation is said to be linear when the amount of change in one variable tends to bear a constant ratio to the amount of change in the other. The graph of the variables having a linear relationship will form a straight line.

Example: $X = 1, 2, 3, 4, 5, 6, 7, 8,$

$Y = 5, 7, 9, 11, 13, 15, 17, 19,$

$$Y = 3 + 2x$$

The correlation would be non linear if the amount of change in one variable does not bear a constant ratio to the amount of change in the other variable.

Methods of studying simple correlation

1. Scatter Diagram method;
2. Karl Pearson's Coefficient of correlation;
3. Spearman's Rank Correlation

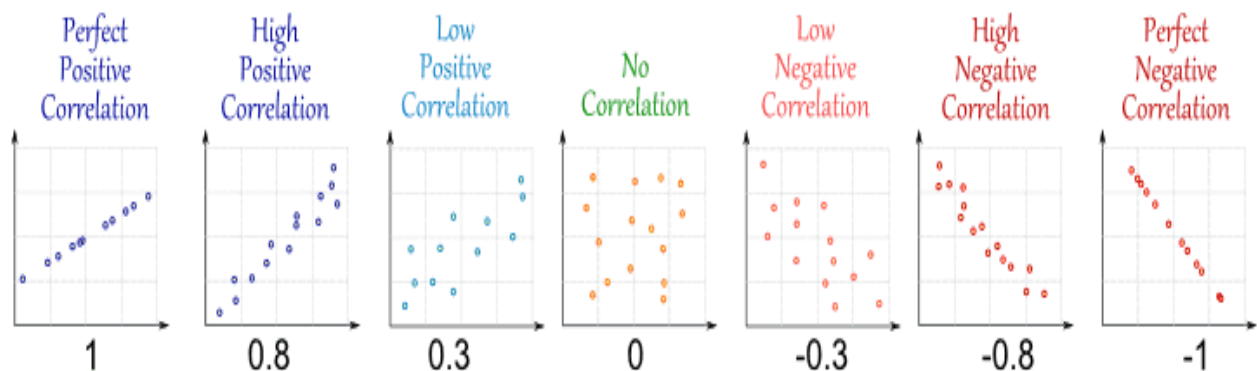
Scatter diagram method

The diagrammatic way of representing bivariate data is called scatter diagram.

Suppose, $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ are n pairs of observations. 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.

Scatter diagrams for different values of r are as follows:

- Scatter diagrams for different values of r are as follows:



Interpret of r

$r = +1$, indicates a perfect positive relationship between x and y. the scatter diagram will be as in fig. 1.1

$r = -1$, indicates a perfect negative relationship between x and y. the scatter diagram will be as in fig. 1.2

$r = 0$, means there is no linear relationship between x and y. In this case the two variables are linearly independent. the scatter diagram will be as in fig. 1.5 and 1.6

$0 < r < 1$, indicates a positive relationship between x and y. In this case the scatter diagram will be as in fig. 1.3

$-1 < r < 0$, indicates a negative relationship between x and y. In this case the scatter diagram will be as in fig. 1.4

Correlation coefficient

The numerical value by which we measure the strength of linear relationship between two or more variables is called correlation coefficient.

Let, $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be the pairs of n observations. Then the correlation coefficient between x and y is denoted by r_{xy} and defined as,

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \dots\dots\dots(1)$$

Equation (1) is also called Karl pearson's coefficient of correlation formula given by 1890.

Algebraically (1) reduces to

$$r = \frac{\sum x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n}}{\sqrt{\left\{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}\right\} \left\{\sum_{i=1}^n y_i^2 - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n}\right\}}}$$

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

Assumptions of Pearson's Correlation Coefficient

- There is linear relationship between two variables, i.e. when the two variables are plotted on a scatter diagram a straight line will be formed by the points.
- Cause and effect relation exists between different forces operating on the item of the two variable series.

Comment on Correlation Coefficient

1 = Perfect positive correlation

$0.7 \leq c < 1$ = Strong positive correlation

$0.4 \leq c < 0.7$ = Fairly positive correlation

$0 < c < 0.4$ = Weak positive correlation

0 = No correlation

$0 > c > -0.4$ = Weak negative correlation

$-0.4 \geq c > -0.7$ = Fairly negative correlation

$-0.7 \geq c < -1$ = Strong negative correlation

-1 = Perfect negative correlation

Properties of correlation coefficient

1. Correlation coefficient is independent of change of origin and scale of measurement.
2. Correlation coefficient lies between -1 to +1. i.e, $-1 < r_{xy} < 1$.
3. Correlation coefficient is symmetric. i.e, $r_{xy} = r_{yx}$
4. Correlation coefficient is the geometric mean of regression coefficients i.e, $r_{xy} = \sqrt{b_{yx} \times b_{xy}}$
5. For two independent variable correlation coefficient is zero
6. It is always unit free.

Advantages of Pearson's Coefficient

- It summarizes in one value, the degree of correlation & direction of correlation also

Limitation of Pearson's Coefficient

- Always assume linear relationship
- Interpreting the value of r is difficult.
- Value of Correlation Coefficient is affected by the extreme values.
- Time consuming methods

Coefficient of Determination

The convenient way of interpreting the value of correlation coefficient is to use of square of coefficient of correlation which is called Coefficient of Determination.

The Coefficient of Determination = r^2 .

Suppose: $r = 0.9$, $r^2 = 0.81$ this would mean that 81% of the variation in the dependent variable has been explained by the independent variable.

The maximum value of r^2 is 1 because it is possible to explain all of the variation in y but it is not possible to explain more than all of it.

Coefficient of Determination = Explained variation / Total variation

An example of Coefficient of Determination

When $r = 0.60$, $r^2 = 0.36$ -----(1)

$r = 0.30$, $r^2 = 0.09$ -----(2)

This implies that in the first case 36% of the total variation is explained whereas in second case 9% of the total variation is explained .

Theorem: Show that Correlation coefficient lies between -1 to +1 i.e, $-1 \leq r_{xy} \leq 1$.

Proof: Let, $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be the pairs of n observations. Then the correlation coefficient between x and y is denoted by r_{xy} and defined as,

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \dots\dots\dots(1)$$

Suppose, $(x_i - \bar{x}) = X$ and $(y_i - \bar{y}) = Y$ therefore

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

Let us consider the following expression which is always positive.

$$\text{i.e, } \sum \left(\frac{X}{\sqrt{\sum X^2}} \pm \frac{Y}{\sqrt{\sum Y^2}} \right)^2 \geq 0$$

$$\text{or, } \sum \left(\frac{X^2}{\sum X^2} \pm 2 \frac{X}{\sqrt{\sum X^2}} \frac{Y}{\sqrt{\sum Y^2}} + \frac{Y^2}{\sum Y^2} \right) \geq 0$$

$$\text{or, } \left(\frac{\sum X^2}{\sum X^2} \pm 2 \frac{\sum XY}{\sqrt{\sum X^2 \sum Y^2}} + \frac{\sum Y^2}{\sum Y^2} \right) \geq 0$$

$$\text{or, } 1 \pm 2r + 1 \geq 0$$

$$\text{or, } 2(1 \pm r) \geq 0$$

$$\text{or, } (1 \pm r) \geq 0 \dots\dots(i)$$

From (i), $1+r \geq 0$ [considering +ve sign.]

$$\text{or, } r \geq -1$$

$$\text{or, } -1 \leq r \dots\dots(ii)$$

$$\text{and } 1-r \geq 0$$

$$\text{or, } 1 \geq r$$

or, $r \leq 1$ (iii)

From (ii) and (iii) we get, $-1 < r < 1$.

i.e, coefficient lies between -1 to +1.

Theorem: Show that for two independent variable correlation coefficient is zero.

Proof: Let, $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be the pairs of n observations. Then the arithmetic mean of x_i is \bar{x} and y_i is \bar{y} . Since x and y are independent therefore,

$$\text{Covariance, Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n} = 0$$

$$\text{or, } \sum (x_i - \bar{x})(y_i - \bar{y}) = 0$$

$$\begin{aligned} \text{We Know, } r_{xy} &= \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= \frac{0}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \\ &= 0 \text{ (proved)} \end{aligned}$$

Application Problem-1: If $y = mx + c$, then find the correlation coefficient between x and y.

Solution: Let, $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ be the pairs of n observations. Then the correlation coefficient between x and y is denoted by r_{xy} and defined as,

$$r_{xy} = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (y_i - \bar{y})^2}} \dots\dots\dots(1)$$

Now, $y = mx + c$(ii)

$$\begin{aligned}
 \text{Therefore, } r_{xy} &= \frac{\sum_{i=1}^n (x_i - \bar{x})(mx_i + c - m\bar{x} - c)}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (mx_i + c - m\bar{x} - c)^2}} \dots\dots\dots(1) \\
 &= \frac{\sum_{i=1}^n (x_i - \bar{x})(mx_i - m\bar{x})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (mx_i - m\bar{x})^2}} \\
 &= \frac{m \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})}{m \sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \sum_{i=1}^n (x_i - \bar{x})^2}} \\
 &= \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{\sum_{i=1}^n (x_i - \bar{x})^2} = 1
 \end{aligned}$$

Procedure for computing the correlation coefficient

Calculate the mean of the two series 'x' & 'y'

Calculate the deviations 'x' & 'y' in two series from their respective mean.

Square each deviation of 'x' & 'y' then obtain the sum of the squared deviation i.e. $\sum x^2$ & $\sum y^2$

Multiply each deviation under x with each deviation under y & obtain the product of 'xy'. Then obtain the sum of the product of x, y i.e. $\sum xy$

Substitute the value in the formula.

Application Problem-1: A research physician recorded the pulse rates and the temperatures of water submerging the faces of ten small children in cold water to control the abnormally rapid heartbeats. The results are presented in the following table. Calculate the correlation coefficient between temperature of water and reduction in pulse rate.

Temperature of water	68	65	70	62	60	55	58	65	69	63
Reduction in pulse rate.	2	5	1	10	9	13	10	3	4	6

Solution: Calculating table of correlation coefficient.

x_i	y_i	x_i^2	y_i^2	$x_i y_i$
68	2	4624	4	136
65	5	4225	25	325
70	1	4900	1	70
62	10	3844	100	620
60	9	3600	81	540
55	13	3025	169	715
58	10	3364	100	580
65	3	4225	9	195
69	4	4761	16	276
63	6	3969	36	378
$\sum x_i = 635$	$\sum y_i = 63$	$\sum x_i^2 = 40537$	$\sum y_i^2 = 541$	$\sum x_i y_i = 3835$

$$\sum x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n}$$

We know, $r_{xy} = \frac{\sum x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n}}{\sqrt{\left\{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}\right\}\left\{\sum_{i=1}^n y_i^2 - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n}\right\}}}$

$$= \frac{3835 - \frac{635 \times 63}{10}}{\sqrt{\left\{40537 - \frac{(635)^2}{10}\right\}\left\{541 - \frac{(63)^2}{10}\right\}}}$$

$$= -0.94$$

The result -0.94, indicates that the correlation coefficient between temperature of water and reduction in pulse rate is highly negatively correlated.

Assignment problem-1: Compute r for the for the following paired sets of values:

i.(x, y): (1,2) , (2, 3), (3, 5), (4, 4), (5, 7)

ii. (x, y): (1,1) , (2, 3), (3, 5), (4, 7), (5, 9)

iii.(x, y): (1,10) , (2, 8), (3, 6), (4, 4), (5, 2)

iv.(x, y): (2,9) , (3, 5), (4, 6), (5, 2), (6, 1)

v.(x, y): (-2,4) , (-1, 1), (0, 0), (1, 1), (2, 4)

Solution 1: (x, y): (1,2) , (2, 3), (3, 5), (4, 4), (5, 7)

The formula for finding correlation coefficient is

$$r_{xy} = \frac{\sum x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n}}{\sqrt{\left\{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}\right\}\left\{\sum_{i=1}^n y_i^2 - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n}\right\}}}$$

Let us make a table to calculate correlation coefficient.

x_i	y_i	x_i^2	y_i^2	$x_i y_i$
1	2	1	4	2
2	3	4	9	6
3	5	9	25	15
4	4	16	16	16
5	7	25	49	35

$\sum x_i = 15$	$\sum y_i = 21$	$\sum x_i^2 = 55$	$\sum y_i^2 = 103$	$\sum x_i y_i = 74$
-----------------	-----------------	-------------------	--------------------	---------------------

$$r_{xy} = \frac{\sum x_i y_i - \frac{\left(\sum_{i=1}^n x_i\right)\left(\sum_{i=1}^n y_i\right)}{n}}{\sqrt{\left\{\sum_{i=1}^n x_i^2 - \frac{\left(\sum_{i=1}^n x_i\right)^2}{n}\right\}\left\{\sum_{i=1}^n y_i^2 - \frac{\left(\sum_{i=1}^n y_i\right)^2}{n}\right\}}}$$

$$= \frac{74 - \frac{15 \times 21}{5}}{\sqrt{\left\{55 - \frac{(15)^2}{5}\right\}\left\{103 - \frac{(21)^2}{5}\right\}}}$$

$$= 0.90$$

Comment: There exists a strong positive relationship between x and y.

Problem: above ii-v (Assignment)

Assignment Problem-2: The following table gives the ages and blood pressure of 10 women:

Age in years x	56	42	36	47	49	42	72	63	55	60
Blood pressure y	147	125	118	128	125	140	155	160	149	150

Draw a scatter diagram

Find correlation coefficient between x and y and comment.

Ans: Try your-self

Assignment Problem-3: The scores of 12 students in their mathematics and physics classes are:

Mathematics	2	3	4	4	5	6	6	7	7	8	10	10
Physics	1	3	2	4	4	4	6	4	6	7	9	10

Find the correlation coefficient distribution and interpret it.

Comment on the followings:

(i) $r=0$ (ii) $r=-1$ (iii) $r=1$ (iv) $r \geq 1$ (v) $r < 1$

(i) $r=0$, indicates that the correlation coefficient between x and y is zero.

(ii) $r=-1$, indicates that the correlation coefficient between x and y is perfect negative.

(iii) $r=1$, indicates that the correlation coefficient between x and y is perfect positive.

(iv) $r \geq 1$ i.e, $r=1$ and $r > 1$ i.e, $r > 1$, is not possible, because the Correlation coefficient lies between -1 to +1.

(v) $r < 1$, not possible because, the Correlation coefficient lies between -1 to +1.

Uses of correlation coefficient.

1. To find the relationship between two variables.
2. To find the relationship between dependent variable and combined influence of a group of independent variables.
3. To solve many problem in biology.
4. In social studies like relationships between crime and educations, correlation analysis has got definite role to play.
5. In economies this is used specially.

RANK CORRELATION

Rank correlation: In some situation it is difficult to measure the values of the variables from bivariate distribution numerically, but they can be ranked. The correlation coefficient between these two ranks is usually called rank correlation coefficient, given by Spearman (1904). It is denoted by R. this is the only method for finding relationship between two qualitative variables like beauty, honesty, intelligence, efficiency and so on.

When there are no ties, the formula for computing the spearman's rank correlation coefficient

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

Here, R= rank correlation coefficient, n = number of pairs of observations being ranked.

d = difference between rank of x and rank of y.

Remarks:

(i) We always have $\sum d_i = \sum (R_1 - R_2) = 0$

(ii) Like simple correlation coefficient, rank correlation coefficient lies between -1 to +1.

Note: For finding rank correlation coefficient, we may have two types of data:

Actual observations are given

Actual ranks are given

Interpretation of Rank Correlation Coefficient (R)

The value of rank correlation coefficient, R ranges from -1 to +1

If $R = +1$, then there is complete agreement in the order of the ranks and the ranks are in the same direction

If $R = -1$, then there is complete agreement in the order of the ranks and the ranks are in the opposite direction

If $R = 0$, then there is no correlation

Application Problem-1: Obtain the rank correlation co-efficient for the following data:

A:	80	75	90	70	65	60
B:	65	70	60	75	85	80

Solution: Here ranks of the score are not given. Let us start ranking from the highest value for both the variables as shown in the table given below:

A	B	Rank of A (x)	Rank of B (y)	d = x-y	d ²
80	65	2	5	-3	9
75	70	3	4	-1	1
90	60	1	6	-5	25
70	75	4	3	1	1
65	85	5	1	4	16
60	80	6	2	4	16
Total				$\sum d_i = 0$	$\sum d_i^2 = 68$

$$R = 1 - \frac{6 \sum d^2}{n(n^2 - 1)} = 1 - \frac{6 \times 68}{6(6^2 - 1)} = -0.94$$

Conclusion: There exist strongly negative relationship between A and B.

Application Problem -2: Obtain the rank correlation co-efficient for the following data:

Examiner	A	B	C	D	E
I	1	2	3	4	5
II	2	4	1	5	4

Solution: Here ranks of the score are given:

Ranking by examiner-I: R ₁	Ranking by examiner-II: R ₂	d = R ₁ - R ₂	d ²
1	2	-1	1
2	3	-1	1
3	1	2	4
4	5	-1	1
5	4	1	1
Total		$\sum d_i = 0$	$\sum d_i^2 = 8$

$$R = 1 - \frac{6 \sum d^2}{n(n^2 - 1)} = 1 - \frac{6 \times 8}{5(5^2 - 1)} = 0.6$$

Comment: There is a positive rank correlation coefficient between the rankings of two examiners.

Repeated ranks or ties observations:

When ranks are repeated the following formula is used for finding rank correlation coefficient:

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

Problems of equal ranks or tie in ranks:

Application Problem -3: The following data refer to the marks obtained by 8 students in mathematics and statistics:

Marks in mathematics	20	80	40	12	28	20	15	60
Marks in statistics	30	60	20	30	50	30	40	20

Compute rank correlation coefficient and comment.

Solution: let the marks obtained by mathematics be x and the marks obtained by statistics be y.

Table for computation of rank correlation.

x	y	Rank of x (R ₁)	Rank of y (R ₂)	d = R ₁ - R ₂	d ²
---	---	-----------------------------	-----------------------------	-------------------------------------	----------------

20	30	3.5	4	-0.5	0.25
80	60	8	8	0	0
40	20	6	2	4	16
12	30	1	4	-3	9
28	50	5	7	-2	4
20	30	3.5	4	-0.5	0.25
15	40	2	6	-4	16
60	10	7	1	6	36
					$\sum d_i^2 = 81.5$

Here, $m_1 = 2$, $m_2 = 3$, $n = 8$

$$R = 1 - \frac{6 \left\{ 81.5 + \frac{1}{12}(2^3 - 2) + \frac{1}{12}(3^3 - 3) \right\}}{8(8^2 - 1)}$$

$$= 0$$

Merits Spearman's Rank Correlation

- This method is simpler to understand and easier to apply compared to Karl Pearson's correlation method.
- This method is useful where we can give the ranks and not the actual data. (qualitative term)
- This method is to use where the initial data in the form of ranks.

Limitation Spearman's Correlation

- Cannot be used for finding out correlation in a grouped frequency distribution.
- This method should be applied where N exceeds 30.

Assignment problem-4:

The following figures relate to advertisement expenditure and profit:

Profit (Tk.Crore):x	25	28	27	33	31	10	16	16	18	23
Adv. Exp.(Tk. Lakh):y	87	91	92	95	93	52	68	72	78	86

(i) Draw a scatter diagram and comment

(ii) Calculate Karl Pearson's and Spearman rank correlation coefficients and comment.

Assignment problem-5:

The following figures relate to advertisement expenditure and sales of a company:

Adv. Exp. (Tk. Lac)	62	67	73	78	85	78	91	92	96	98
Sales (Tk.Crore)	11	13	17	18	21	24	21	27	26	21

Calculate Karl Pearson's correlation coefficient and Spearman rank correlation

Coefficient and comment.

Website:

http://www.pindling.org/Math/Statistics/Textbook/Examples/Chapter3/chapter3_examples.htm