# Statistics for Computing

## (CSC 502 0.0 )
## MSc in Computer Science

# Prof. (Dr.) R.M. KAPILA RATHNAYAKA

*B.Sc. Special (Math. & Stat. ) (Ruhuna), M.Sc. (Industrial Mathematics) (USJ),*
*M.Sc. (Stat. ) (WHUT, China),*
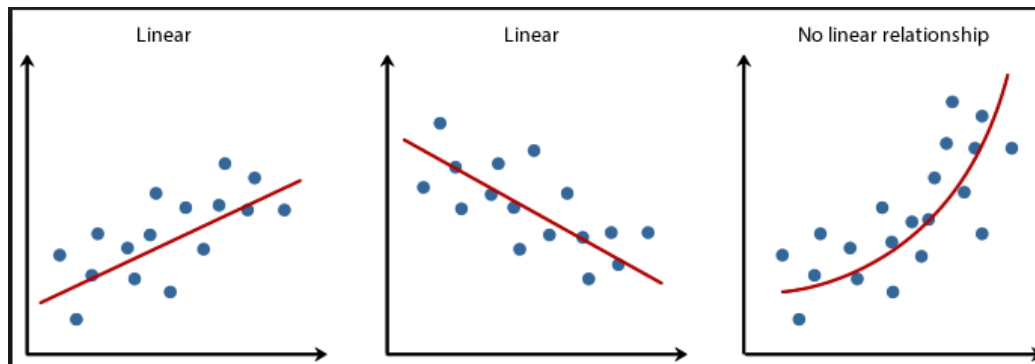*Ph.D. (Applied Statistics, WHUT)*

# Introduction to Correlation

# Chapter 04

# Correlation and Regression Analysis

- <u>Correlation and regression</u> are the two most commonly used techniques for **investigating the relationship** between <u>quantitative variables</u>.

- <u>Correlation is used to give the relationship between the variables</u> whereas linear regression uses an equation to express this relationship.

- In this section we will first discuss <u>correlation analysis</u>, which is used to <u>quantify the association between two continuous variables</u>

  ✓between an independent and a dependent variable
  ✓ between two independent variables

# Correlation Definition



- If an <u>increase (or decrease) in one variable causes a corresponding increase (or decrease) in another</u> then the two variables are said to be directly correlated.

- Similarly, if an <u>increase in one causes a decrease in another or vice versa</u>, then the variables are said to be indirectly correlated.

- <u>If a change in an independent variable does not cause a change in the dependent variable then they are uncorrelated.</u>

# Method 01: Karl Pearson's Coefficient of Correlation

- Correlation is a *bivariate analysis* that measures the *strength of association between two variables* and the *direction of the relationship*.

- Usually, in statistics, we measure four types of correlations:

  - Pearson correlation,

  - Kendall rank correlation,

  - Spearman correlation

# Method 01: Karl Pearson's Coefficient of Correlation

- **Assumptions**

- For the Pearson *r* correlation, <u>both variables should be normally distributed</u> (normally distributed variables have a bell-shaped curve).

- Other assumptions include <u>linearity and homoscedasticity</u>.

- <u>Linearity assumes a straight line relationship between each of the two variables</u> and <u>homoscedasticity assumes</u>

# Karl Pearson's Coefficient of Correlation
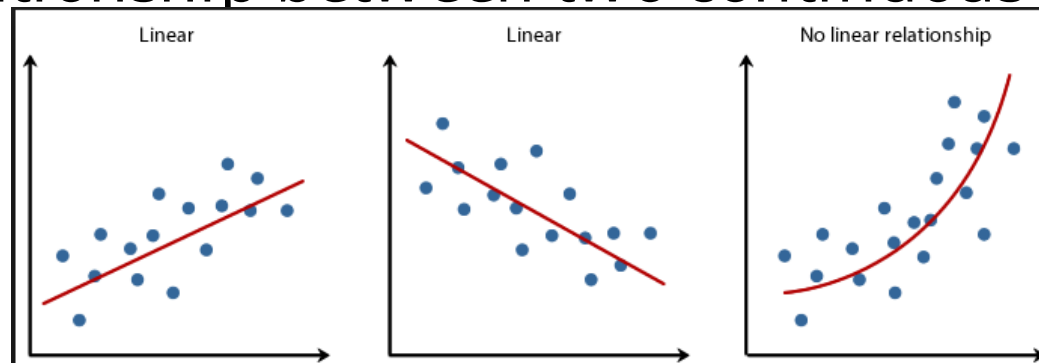
- **Definition:**

   **Karl Pearson's Coefficient of Correlation** is widely used mathematical method wherein the numerical expression is used to calculate ***the degree*** and ***direction*** of the relationship between <u>linear related variables</u>.

1. **Karl Pearson's Coefficient of Correlation-** evaluates the linear relationship between two continuous variables
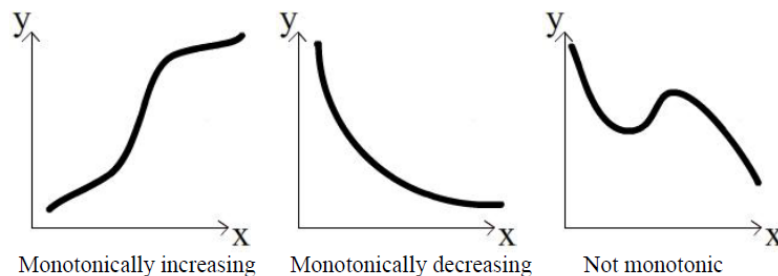
# **Correlation Definition**

- Correlation can be defined as a <u>measurement that is used to quantify the relationship between variables</u>.

1. **Karl Pearson's Coefficient of Correlation-** evaluates the linear relationship between two continuous variables



2. **Spearman rank-order correlation-** evaluates the monotonic relationship between two continuous or ordinal variables

# Method 01: Karl Pearson's Coefficient of Correlation

- The coefficient of correlation is denoted by symbol *'r'*.

$$r = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \sum (y_i - \bar{y})^2}}$$

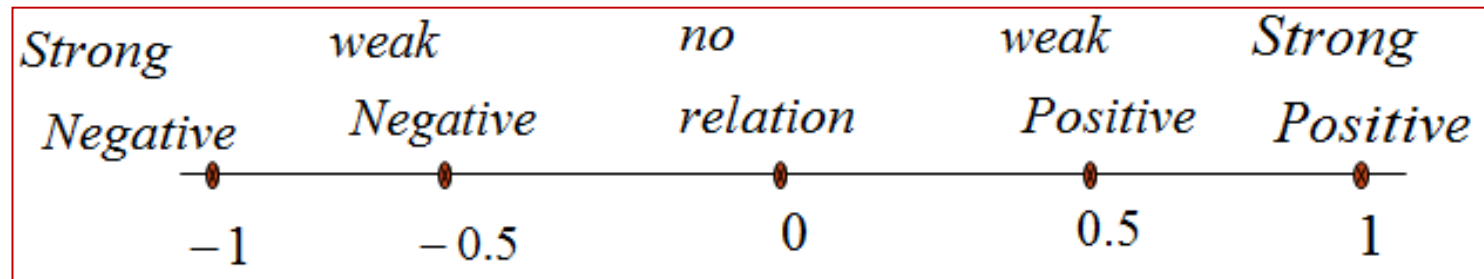$\bar{X}$ = mean of X variable
$\bar{Y}$ = mean of Y variable

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

# Correlation Analysis

- Correlation is one of the most common and most useful statistics

- It is a term that refers to the <u>strength of a relationship between two variables</u> *(single number that describes the degree of relationship between two variables)*.

- <u>A strong, or high, correlation means that two or more variables have a strong relationship with each other</u>.

- A weak or low correlation means that the variables are <u>hardly related.</u>

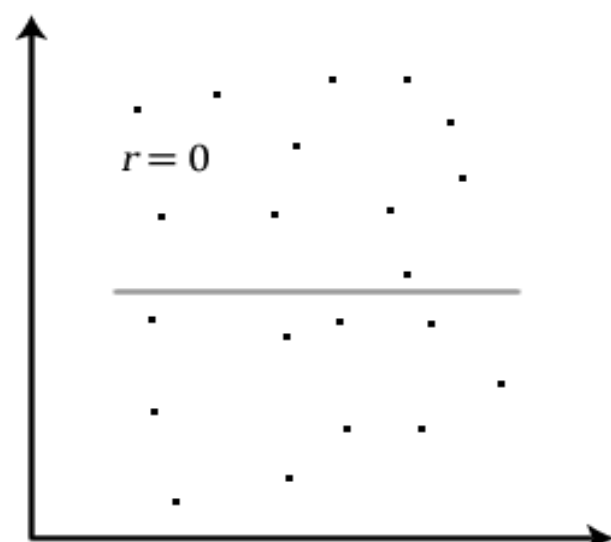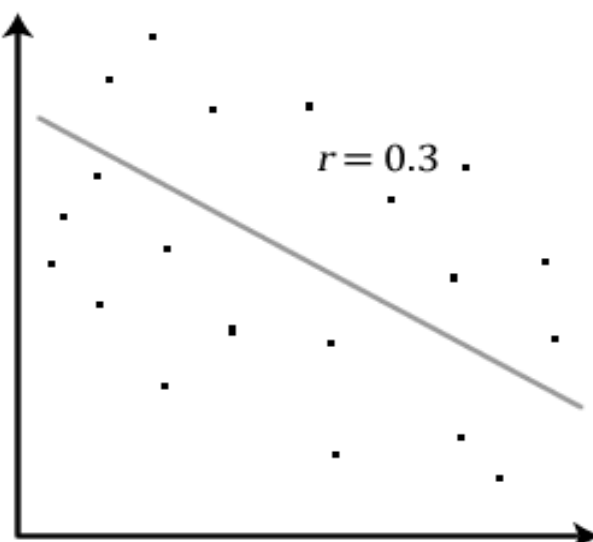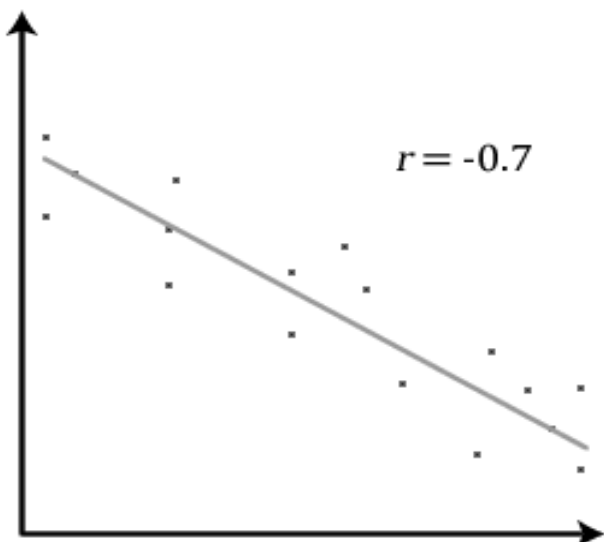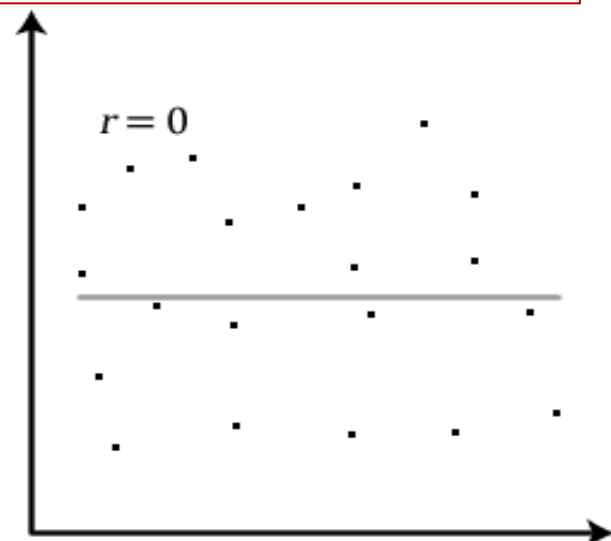# Properties of Coefficient of Correlation

1. The value of the coefficient of correlation (r) always **lies between ±1**.



| Strong Negative | weak Negative | no relation | weak Positive | Strong Positive |
|---|---|---|---|---|
| −1 | − 0.5 | 0 | 0.5 | 1 |

r=+1, perfect positive correlation

r=-1, perfect negative correlation

r=0, no correlation

$r = 0.7$

$r = 0.3$

$r = 0$

$r = -0.7$

$r = 0.3$

$r = 0$

**Example (01) :** Calculate the coefficient of correlation between X and Y from the following data.

| X | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| Y | 2 | 4 | 5 | 3 | 8 | 6 | 7 |

| X | 1 | 2 | 3 | 4 | 5 | 6 | 7 |
|---|---|---|---|---|---|---|---|
| Y | 2 | 4 | 5 | 3 | 8 | 6 | 7 |

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

$$\sum x = 28 \qquad \sum y = 35$$

$0.79 \Rightarrow$ *Strong positive*

## Example (02):

A sample of 6 children was selected, data about their age in years and weight in kilograms was recorded as shown in the following table . It is required to find the correlation between age and weight.

| Weight (Kg) | Age (years) | serial No |
|---|---|---|
| 12 | 7 | 1 |
| 8 | 6 | 2 |
| 12 | 8 | 3 |
| 10 | 5 | 4 |
| 11 | 6 | 5 |
| 13 | 9 | 6 |

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

$$r = \frac{461 - \dfrac{41 \times 66}{6}}{\sqrt{\left[291 - \dfrac{(41)^2}{6}\right]\left[742 - \dfrac{(66)^2}{6}\right]}}$$

r = 0.759

strong direct correlation