# QUANTITATIVE METHODS FOR ECONOMIC ANALYSIS-I

**III Semester**

**CORE COURSE**

# BA ECONOMICS

*(2014 Admission onwards)*
*(CU CBCSS)*



# UNIVERSITY OF CALICUT

## SCHOOL OF DISTANCE EDUCATION

**Calicut University P.O. Malappuram, Kerala, India 673 635**

**703**

# UNIVERSITY OF CALICUT

## SCHOOL OF DISTANCE EDUCATION

**STUDY MATERIAL**

**BA – ECONOMICS**
**(2014 Admission)**

**III Semester**

## ECO3 B03 - QUANTITATIVE METHODS FOR ECONOMIC ANALYSIS - I

**Prepared by**

| | |
|---|---|
| **Module I:** | **Shihabudheen  M. T.** |
| | **Assistant  Professor** |
| | **Department of Economics** |
| | **Farook College, Calicut** |
| | |
| **Module II:** | **Shabeer  K. P.** |
| | **Assistant Professor** |
| | **Department of Economics** |
| | **Governemnt College, Kodencherry** |
| | |
| **Module III, IV and V** | **Dr. Chacko Jose,** |
| | **Associate Professor** |
| | **Department of Economics** |
| | **Sacred Heart College, Chalakudy,T hrissur** |
| | |
| **Edited and Compiled By**: | **Dr. Yusuf Ali  P. P., Chairman, Board of Studies (UG)** |
| | **Associate Professor** |
| | **Department of Economics** |
| | **Farook College** |

# Contents

**Module I**       **Algebra**

**Module II**     **Basic Matrix Algebra**

**Module III**    **Functions and Graphs**

**Module IV**    **Meaning of Statistics and Description of of Data**

**Module V**     **Correlation and Regression Analysis**

## MODULE – I

## ALGEBRA

### Exponents

If we add the letter a, six times, we get a+ a+ a+ a+ a = 5a.  i.e., 5 x a.  If we multiply this, we get a $\times$ a $\times$ a $\times$ a $\times$ a = $a^5$ , i.e. a is raised to the power 5.  Here a is the factor.  a is called the base.  5 is called the exponent (Power or Index)

### 1.    Meaning of positive Integral power.

$a^n$ is defined only positive integral values.  If a is a positive integer, $a^n$ is defined as the product of *n* factors.  Each of which is *a*

$$a^n = a \times a \times a \ldots\ldots\ldots \text{ n times}$$

*Eg:*        $2^3 = 2 \times 2 \times 2 = 8$

### 2.    Meaning of zero exponent ( zero power)

If a≠0, $a^0 = 1$, i.e., any number (other than zero) raised to zero = 1

*Eg:*        $7^0 = 1$,        $\left(\dfrac{2}{3}\right)^0 = 1$

### 3.    Meaning of negative integral power (negative exponent)

If n is a positive integer and a≠0, $a^{-n}$ is the reciprocal of $a^n$ .

i.e.,    $a^{-n} = \dfrac{1}{a^n}$

*Eg:*        $3^{-2} = \dfrac{1}{3^2} = \dfrac{1}{9}$

### 4. Root of a number

**(a)**    Meaning of square root

If $a^2 = b$, then *a* is the square root of b and we write, $a = \sqrt{b}$ or $a = b^{\frac{1}{2}}$

*Eg:*  $\qquad 3^2 = 9, \qquad\qquad 3 = \sqrt{9} \qquad$ or $3 = 9^{\frac{1}{2}}$

**(b)** Meaning of cube root

If $a^3 = b$, then $a$ is the cube root of b and we write, $a = \sqrt[3]{b}$ , or $a = b^{\frac{1}{3}}$

*Eg:*  $\qquad 2^3 = 8, \qquad\qquad 2 = \sqrt[3]{8} \qquad$ or $2 = (8)^{\frac{1}{3}}$

**(c)** Meaning of $n^{th}$ root

If $a^n = b$ then $a$ is the $n^{th}$ root of $b$ and we write $a = \sqrt[n]{b}$ or $a = (b)^{\frac{1}{n}}$

Eg:  $\qquad 3^4 = 81, \qquad$ i.e $3 = \sqrt[4]{81}$ , or $3 = (81)^{\frac{1}{4}}$

**(d)** Meaning of positive fractional power.

If $m$ and $n$ are positive integers, then $a^{\frac{m}{n}}$ is defined as $n^{th}$ root of $m^{th}$ power of $a$.

i.e.,  $\qquad a^{\frac{m}{n}} = \sqrt[n]{a^m}$

*Eg:*  $\qquad (16)^{\frac{2}{4}} = \left(\sqrt[4]{16}\right)^2 = 2^2 = 4$

**(e)** Meaning of negative fractional powers.

If $m$ and $n$ are positive integers, $a^{\frac{-m}{n}}$ is defined as $\dfrac{1}{a^{\frac{m}{n}}}$ or $\dfrac{1}{\sqrt[n]{a^m}}$

*Eg.*  $\qquad (16)^{\frac{-3}{2}} = \dfrac{1}{(16)^{\frac{3}{2}}} = \dfrac{1}{\left(\sqrt{16}\right)^3}$

$\qquad\qquad\qquad = \dfrac{1}{(4)^3} = \dfrac{1}{64}$

**Laws of Indices**

1. Product rule:-
   When two powers of the same base are multiplied, indices or exponents are added.

   i.e.,  $\qquad a^m \times a^n = a^{m+n}$

*Eg.*  $\qquad 2^2 \times 2^3 = 2^{2+3} \qquad = 2^5 = 32$

$$3^1 \times 3^4 = 3^{1+4} = 3^5 = 243$$

$$3^{-2} \times 3^5 = 3^{-2+5} = 3^3 = 27$$

2. **Quotient rule:-**

When some power of *a* is divided by some other power of *a*, index of the denominator is subtracted from that of the numerator.

i.e,  $$a^m \div a^n = a^{m-n}$$

$$5^3 \div 5^2 = 5^{3-2} = 5$$

3. **Power rule:-**

When some power of a is raised by some other power, the indices are multiplied.

i.e.,  $$\left(a^m\right)^n = a^{mn}$$

Eg.  $$\left(3^2\right)^3 = 3^{2 \times 3} = 3^6 = 729$$

$$\left(4^3\right)^{\frac{1}{3}} = 4^{3 \times \frac{1}{3}} = 4^{\frac{3}{3}} = 4^1 = 4$$

4.  $$\left(ab\right)^n = a^n b^n$$

$$\left(2 \times 3\right)^2 = 2^2 \times 3^2 = 4 \times 9 = 36$$

$$3^3 \times 4^3 = 27 \times 64 = 1728$$

5.  $$\left(\frac{a}{b}\right)^n = \frac{a^n}{b^n}$$

$$\left(\frac{4}{3}\right)^2 = \frac{4^2}{3^2} = \frac{16}{9}$$

**Extension**

1.  $$a^m \times a^n \times a^P = a^{m+n+P}$$

2.  $$\frac{a^m \times a^n}{a^P} = a^{m+n-P}$$

3.  $$\left(abc\right)^n = a^n \times b^n \times c^n$$

4. $\left(\dfrac{ab}{cd}\right)^n = \dfrac{a^n \times b^n}{c^n \times d^n}$

Examples

1. $a^{\frac{2}{3}} \cdot b^{\frac{5}{3}} \cdot c^{\frac{2}{7}} \times a^{\frac{3}{2}} \cdot b^{\frac{3}{5}} \times c^{\frac{7}{2}}$

$a^{\frac{2}{3}+\frac{3}{2}} \cdot b^{\frac{5}{3}+\frac{3}{5}} \cdot c^{\frac{2}{7}+\frac{7}{2}}$

$= a^{\frac{13}{6}} \cdot b^{\frac{34}{15}} \cdot c^{\frac{53}{14}}$

2. $6a^2 b \times 8\, a^3 b^4 = \ ?$ $\qquad 6 \times 8 \times a^2 \times a^3 \times b \times b^4$

$= 48\, a^5 b^5$

3. $6a^{-1} \times b^{-5} \times 4a^{-2} b^3 = 24\, a^{-3} b^{-2} = \ ?$ $\qquad \dfrac{24}{a^3 b^2}$

4. $63\, x^8 y^5 \div 9 x^5 y^3$

$= \dfrac{63}{9} \cdot \dfrac{x^8}{x^5} \cdot \dfrac{y^5}{y^3} \qquad = \qquad 7\, x^{8-5} \cdot y^{5-3} \qquad = \qquad 7\, x^3 y^2$

5. $15\, x^7 y^3 \div \dfrac{5}{3} x^3 y^{-1}$

$= \dfrac{\frac{15}{5}}{\frac{5}{3}} \qquad \cdot \dfrac{x^7}{x^3} \cdot \dfrac{y^3}{y^{-1}}$

$= 15 \times \dfrac{3}{5} x^4 y^4$

$= \dfrac{45}{5} x^4 y^4 \qquad = \qquad 9\, x^4 y^4$

6. $\dfrac{\left(x^2 y\right)^5 \left(x^{-2} y^{-3}\right)^2}{\left(x^2\right)^3 \left(y^3\right)^4} \qquad = \qquad \dfrac{\left(x^{2\times5} y^5\right)\left(x^{-4} y^{-6}\right)}{x^6 y^{12}}$

$$= \frac{x^{10} y^5 x^{-4} y^{-6}}{x^6 y^{12}} \qquad = \frac{x^{10} \times x^{-4} y^5 \times y^{-6}}{x^6 y^{12}}$$

$$= \frac{x^6 y^{-1}}{x^6 y^{12}} \qquad = x^{6-6} y^{-1-12} = x^6 \times y^{-13}$$

$$= 1 \quad \times \quad \frac{1}{y^{13}} \qquad = \frac{1}{y^{13}}$$

7.    Multiply    $( \; x^{\frac{1}{2}} - x^{\frac{1}{4}} y^{\frac{1}{4}} + y^{\frac{1}{2}} \; )$ by $x^{\frac{1}{4}} + y^{\frac{1}{4}}$

$$= \left[ \left( x^{\frac{1}{4}} \right)^2 - x^{\frac{1}{4}} y^{\frac{1}{4}} + \left( y^{\frac{1}{4}} \right)^2 \right] \quad \left( x^{\frac{1}{4}} + y^{\frac{1}{4}} \right)$$

$$= \left( a^2 - ab + b^2 \right) \left( a + b \right)$$

$$= a^3 + b^3 \quad , \quad a = x^{\frac{1}{4}} \; , \quad b = y^{\frac{1}{4}}$$

$$= \left( x^{\frac{1}{4}} \right)^3 + \left( y^{\frac{1}{4}} \right)^3$$

$$= x^{\frac{3}{4}} + y^{\frac{3}{4}}$$

8.    Multiply    $\left( a^{\frac{2}{3}} - a^{\frac{1}{3}} b^{\frac{1}{3}} + b^{\frac{2}{3}} \right)$ by $\left( a^{\frac{1}{3}} + b^{\frac{1}{3}} \right)$

$$= \left( a^{\frac{1}{3}} \right)^2 - a^{\frac{1}{3}} b^{\frac{1}{3}} + \left( b^{\frac{1}{3}} \right)^2 \quad \left( a^{\frac{1}{3}} + b^{\frac{1}{3}} \right)$$

$$= \left( a^2 - ab + b^2 \right) \left( a + b \right)$$

$$= a^3 + b^3$$

$$= \left( a^{\frac{1}{3}} \right)^3 + \left( b^{\frac{1}{3}} \right)^3 \quad = a + b$$

9.    Prove that    $\dfrac{1}{1 + x^{a-b} + x^{a-c}} \; + \; \dfrac{1}{1 + x^{b-c} + x^{b-a}} \; + \; \dfrac{1}{1 + x^{c-a} + x^{c-b}} \; = 1$

$$= \cfrac{1}{1+\cfrac{x^a}{x^b}+\cfrac{x^a}{x^c}} \quad + \quad \cfrac{1}{1+\cfrac{x^b}{x^c}+\cfrac{x^b}{x^a}} \quad + \quad \cfrac{1}{1+\cfrac{x^c}{x^a}+\cfrac{x^c}{x^b}}$$

By multiplying each term by $x^{-a}, \quad x^{-b}, \quad x^{-c}$ respectively we get,

$$= \cfrac{x^{-a}}{x^{-a}\left(1+\cfrac{x^a}{x^b}+\cfrac{x^a}{x^c}\right)} \quad + \quad \cfrac{x^{-b}}{x^{-b}\left(1+\cfrac{x^b}{x^c}+\cfrac{x^b}{x^a}\right)} \quad + \quad \cfrac{x^{-c}}{x^{-c}\left(1+\cfrac{x^c}{x^a}+\cfrac{x^c}{x^b}\right)}$$

$$= \cfrac{x^{-a}}{x^{-a}\times 1+x^{-a}\left(\cfrac{x^a}{x^b}\right)+x^{-a}\left(\cfrac{x^a}{x^c}\right)} \quad + \quad \cfrac{x^{-b}}{x^{-b}\times 1+x^{-b}\left(\cfrac{x^b}{x^c}\right)+x^{-b}\left(\cfrac{x^b}{x^a}\right)}$$

$$+ \quad \cfrac{x^{-c}}{x^{-c}\times 1+x^{-c}\left(\cfrac{x^c}{x^a}\right)+x^{-c}\left(\cfrac{x^c}{x^b}\right)}$$

$$= \cfrac{x^{-a}}{x^{-a}+x^{-b}+x^{-c}} \quad + \quad \cfrac{x^{-b}}{x^{-b}+x^{-c}+x^{-a}} \quad + \quad \cfrac{x^{-c}}{x^{-c}+x^{-a}+x^{-b}}$$

$$= \cfrac{x^{-a}+x^{-b}+x^{-c}}{x^{-a}+x^{-b}+x^{-c}} \quad = 1$$

10. If $a = b^c$, $b = c^a$ and $c = a^b$, prove hat $abc = 1$

$a = b^c$

$a = \left(c^a\right)^c$          ($\because$ b = c$^a$)

$a = c^{ac}$

$a = \left(a^b\right)^{ac}$          ($\because$ c = a$^b$)

$a = a^{abc}$

$\therefore$    $abc = 1$    (Since the bases are same)

11. If $x^a = y$, $y^b = z$ and $z^c = x$, prove that $abc = 1$

$$y = x^a$$

$$y = \left(z^c\right)^a \qquad (\because x = z^c)$$

$$y = z^{ca}$$

$$y = \left(y^b\right)^{ca} \qquad (\because z = y^b)$$

$$y = y^{abc}$$

$$\therefore \quad 1 = abc \quad \text{(Since the bases are same)}$$

12. Simplify $\dfrac{3 \cdot 2^{n+1} + 2^n}{2^{n+2} - 2^{n-1}}$

$$= \frac{3 \times 2^n \times 2^1 + 2^n}{2^n \times 2^2 - \dfrac{2^n}{2^1}}$$

$$= \frac{2^n(6+1)}{2^n\left(4 - \dfrac{1}{2}\right)} \quad = \quad \frac{7}{\dfrac{8-1}{2}} \quad \frac{7}{\dfrac{7}{2}} \quad = \quad \frac{14}{7} \quad = 2$$

13. Solve $2^{x+7} = 4^{x+2}$

$$= \qquad 2^{x+7} = 2^{2(x+2)}$$

$$= \qquad 2^{x+7} = 2^{2x+4}$$

$$= \qquad x+7 = 2x+4$$

$$= \qquad 2x - x = 7 - 4$$

$$= \qquad x = 3$$

## LOGARITHMS

Logarithm of a positive number to a given base is the power to which the base must be raised to get the number.

For eg:- $4^2 = 16$ logarithm of 16 to be the base 4 is 2. It can be written as $\log_4 16 = 2$

Eg:- $\log_7 49 = 2$

Eg: $10^4 = 10{,}000$, $\log_{10} 10{,}000 = 4$

# LAWS OF LOGARITHMS

### 1. Product rule

The logarithm of a product is equal to the sum of the logarithms of its factors.

i.e., $\log_a mn = \log_a m + \log_a n$

Eg: $\log_2 2 \times 3 = \log_2 2 + \log_2 3$

### 2. Quotient rule

The logarithm of a Quotient is the logarithm of the numerator minus the logarithm of the denominator

$$\log_a \frac{m}{n} = \log_a m - \log_a n$$

Eg: $\log_3 \frac{5}{2} = \log_3 5 - \log_3 2$

### 3. Power rule

The logarithm of a number raised to a power is equal to the product of the power and the logarithm of the number.

$$\log_a m^n = n \times \log_a m$$

Eg: $\log_{10} 2^3 = 3 \times \log_{10} 2$

### 4. The logarithm of unity to any base is zero

$\log_a 1 = 0$

Eg: $\log_{10} 1 = 0$

### 5. The logarithm of any number to the same base is unity

i.e., $\log_a a$ = **1**

Eg: $\log_2 2$ = 1

### 6. Base changing rule

The logarithm of a number to a given base is equal to the logarithm of the number to a new base multiplied by the logarithm of the new base to the given base.

$$\log_a m = \log_a m \cdot \log_a b$$

**7.** The logarithm obtained by interchanging the number and the base of a logarithm is the reciprocal of the original logarithm.

$$\log_m a = \frac{1}{\log_a m}$$

Eg:- Find logarithm of 10,000 to the base 10

$$10^4 = 10,000$$

$$\log_{10} 10,000 = 4$$

Eg:- Find logarithm of 125 to the base 5

$$5^3 = 125$$

$$\log_5 125 = 3$$

Eg:- $7^3 = 343$, $\log_7 343 = 3$

Eg:- Find logarithm of

(1) $\log 12 = \log 3 \times 4 = \log 3 \times 2^2$

$$= \log 3 + 2 \log 2$$

**LOGARITHM TABLES**

The logarithm of a number consists of two parts, the integral parts called the characteristics and the decimal part called the mantissa

**Characteristic**

The characteristic of the logarithm of any number greater than 1 is positive and is one less than the number of digits to the left of the decimal point in the given number. The characteristic of the logarithm of any number less then 1 is negative and it is numerically one more than the number of zeros to the right to the decimal point.

| Number | Character | |
|--------|-----------|---|
| 75.3 | 1 | |
| 2400.0 | 3 | One less than the number of digits |
| 144.0 | 2 | |
| 3.2 | 0 | |
| .5 | -1 | |
| .0902 | -2 | One more than the number of zero immediately after the decimal point |
| .0032 | -3 | |
| .0007 | -4 | |

## Antilogarithm

If the logarithm of a number '$a$' is b, then the antilogarithm of 'b' is a

For example if log 61720 = 4.7904, then antilog  4.7904 = 61720

## EQUATIONS

An equation is a statement of equality between two expressions. In other words, an equation sets two expressions, which involves one or more than one variable, equal to each other.

For example, (a) $2x = 10$,    (b) $3x + 2 = 20$,   (c) $x^2 - 5x + 6 = 0$.

An equation consists of one or more unknown variables. In the above example first and second equation (a and c) contain only one unknown variable ( x) and equation 2 contains two unknowns (x and y)

The value (or values) of unknown for which the equation is true are called solution of equations.

Eg:-  In the equation  4x = 2, the value of x is:  x = 2/4 = ½

## Difference between an equation and an Identity

An equation is true for only certain values of the unknown.  But an identity is true for all real values of the unknown.

Eg:-  $x^2 + 2x - 3$  = 0 is true for x = -3 or x = 1.  So it is an equation.  But $(x+1)^2 = x^2 + 2x + 1$  is true for all real values of x.  So it is an identity.

### Solutions of the Equation

An equation is true for some particular value or values of the unknown. The value of the unknown for which equation is true is called solutions of the equation. It is also known as root of the equation

$$x = \frac{10}{2} = 5$$

For example, (a) 2x = 10, so                 , Thus this equation is true for the value x = 5

### Linear and Non-linear Equations

The highest degree of the variables in an equation determines the nature of the equation. If the equation is of first degree, then it is known as linear equation otherwise it is known as non-linear.

$$5x + y = 20$$

For example:                 is a linear equation. It is a linear equation because there is no term involving  $x^2$ ,  $y^2$ ,  $x \times y$ , or  any higher powers of *x* and *y*.

$x^2 - 7x + 12 = 0$                 is a non-linear equation. It is non-linear because the highest degree of the unknown variable in the equation is two.

## Variables

A variable is a symbol or letter used to denote a quantity whose value changes over a period of time. In other words, a variable is a quantity which can assume any one of the values from a range of possible values.

Example: income of the consumer is a variable, since it assumes different values at different time.

## Dependent and Independent Variable

$$y = f(x)$$

If  x  and  y  are two variables such that                 ,for any value of the x there is a corresponding y value, then x is independent variable and y is dependent variable. The value of y depends on the value of x.

$$c = f(y)$$

Example: Consider the consumption function . Here consumption c depends on income. For each value of income there corresponds a value of consumption. Thus c is dependent variable and y is independent variable.

Parameters are similar to variables –that is, letters that stand for numbers– but have a different meaning. We use parameters to describe a set of similar things. Parameters can take on different values, with each value of the parameter specifying a member of this set of similar objects.

**Solution of Simple Linear Equations**

A simple linear equation is an equation which consists of only one unknown and its exponent is one.

**Steps for Solving a Linear Equation in One Variable**

1. Simplify both sides of the equation.
2. Use the addition or subtraction properties of equality to collect the variable terms on one side of the equation and the constant terms on the other.
3. Use the multiplication or division properties of equality to make the coefficient of the variable term equal to 1.

Note: In order to isolate the variable, perform operations on both sides of the equation.

1- **Use of Inverse Operation**
   a) **Use subtraction to undo addition.**

   If $a = b$

   $a - c = b - c$

Example (1): Solve $x + 5 = 15$

   **Solution:**

   $$x + 5 = 15$$

   Subtract 5      $5 = 5$

   $$x = 10$$

   OR

   $x + 5 = 15$        $x = 15 - 5 = 10$

Example (2)    $y + 6 = 2y$

   $$y + 6 = 2y$$

   **Solution:**

   $$y = y$$

   Subtract y

   $$6 = y$$

   OR

   $y + 6 = 2y$    $6 = 2y - y$    $y = 6$

   b) **Use Addition to Undo Subtraction**

   $a = b$

   If

   $a + c = b + c$

   then

For example, solve $x - 4 = 6$

**Solution:**

$$x - 4 = 6$$

Add $\qquad 4 = 4$

$$x = 10$$

OR

$$x - 4 = 6 \qquad x = 6 + 4 \qquad x = 10$$

c) **Use Division to Undo Multiplications**

If $a = b$

then $\dfrac{a}{c} = \dfrac{b}{c}$

**Example:** $3x = 18$

Solution:

$$3x = 18$$

$$\frac{3x}{3} = \frac{18}{3}$$

$$x = 6$$

Answer:

OR

$$x = \frac{18}{3} = 6$$

d) **Use Multiplication to Undo Division**

If $a = b$

then $ac = bc$

Example: $\dfrac{x}{4} = 6$

Solution:

$$4 \cdot \left( \frac{x}{4} \right) = 4 \cdot 6$$

$$x = 24$$

Answer

OR

$$x = 4.6 = 24$$

## 2. Equation having Fractional Coefficient

The coefficient of x also be a rational number. This section discusses how to solve the equation having only one fraction and equation having different fractions.

### a) Equation having Only One Fraction

To clear fractions, multiply both sides of the equation by the denominator of the fractions or by the reciprocal of the fraction

$$\frac{1}{7}x = 5$$

Example (1) :

**Solution**

$$7 \cdot \frac{1}{7}x = 5.7$$

$$x = 35$$

**Answer**:

$$\frac{2}{6}x = 15$$

Example (2):

**Solution**:

$$\frac{6}{2} \cdot \left(\frac{2}{6}x\right) = 15 \cdot \left(\frac{6}{2}\right) \quad x = \frac{90}{2} = 45$$

,

$$x = 45$$

**Answer:**

### b) Equation Containing Fractions having Different Denominator

To clear fractions, multiply both sides of the equation by the LCD of all the fractions. The Lowest Common Denominator (L.C.D) of two or more fractions is the smallest number divisible by their denominators without reminder

$$\frac{x}{3} + \frac{x}{4} = 14$$

For example: solve

**Solution:** Here L.C.D is 12

$$12 \times \left(\frac{x}{3} + \frac{x}{4}\right) = 12 \times 14$$

$$4x + 3x = 168 \qquad 7x = 168$$

$$x = 24$$

**Answer:**

### 3. Equations Containing Parentheses

Follow the following steps to solve the equation which contains parenthesis
a) Remove the parenthesis
b) Solve the resulting equation

$$10 + 3(x - 6) = 16$$

For example: solve

$$10 + 3x - 18 = 16$$

**Solution:**

$$3x - 8 = 16$$

$$3x = 16 + 8$$

$$x = \frac{24}{3} = 8$$

Examples:

1. 4x = 2,     x = 2/4 = ½

2. X -3 = 2,   x = 2 + 3 = 5

3.     Find two numbers of which sum is 25 and the difference is 5

Let one number be x so, that the other is 25-x.

Since the difference is 5, (25-x)-x = 5

$$25\text{-}2x = 5$$

$$\text{-}2 x = \text{-}20$$

$$x = \text{-}20/\text{-}2 = 10$$

So, one number is 10, and other is 25-10 = 15.

### Simultaneous Equations

Simultaneous equations are set of two or more equations, each containing two or more variables whose values can simultaneously satisfy both or all equations in the set. The number of variables will be equal to or less than the number of equations in the set.

#### Simultaneous Equation in Two Unknowns (First Degree)

The simultaneous equation can be solved by the following methods.

a. Elimination method
b. Substitution method
c. Cross multiplication method

#### (A) Elimination method
i.   Multiply the equations with suitable non-zero constants, so that the coefficients of one variable in both equations become equal.
ii.  Subtract one equation from another, to eliminate the variable with equal coefficients. Solve for the remaining variable.

iii. Substitute the obtained value of the variable in one of the equations and solve for the second variable.

**Example**

1. Solve
$$2x + 2y = 40$$
$$3x + 4y = 65$$

**Solution:**

$$2x + 2y = 40$$
...........................(1)

$$3x + 4y = 65$$
.......................... (2)

Multiply equation (1) by 2, we will get
$$4x + 4y = 80$$
...........................(3)

Subtract equation (2) from equation (3)
$$4x + 4y = 80$$
$$3x + 4y = 65$$
$$-$$
$$x = 15$$

Substitute x = 15 either in equation (1) or in equation (2)
Substituting in equation (1), we get
$$2(15) + 2y = 40$$

$$2y = 40 - 30$$

$$y = \frac{10}{2} = 5$$

Checking answers by substituting the obtained value into the original equation.
$$2(15) + 2(5) = 40$$

$$30 + 10 = 40$$

Both sides are equal (L.H.S=R.H.S)

**So the answers x = 15 and y = 5**

2. Solve   4x + 3y = 6

8x + 4y = 18

**Solution:**   4x + 3y = 6   ………………….. (1)

8x + 4y = 18  …………………. (2)

Multiplying first equation by 4, and second equation by 3.

$$16x + 12y = 24 \quad \rightarrow (3)$$

$$24x + 12y = 54 \quad \rightarrow (4)$$

$$-8x = -30$$

$$x = 30/8 = 15/4$$

Substituting x = 15/4 in equation (1), we get

$$4x + 3y = 6$$

$$4 \times \frac{15}{4} + 3y = 6$$

$$3y = -9$$

$$Y = -9/3 = -3$$

The solution is   x = $\frac{15}{4}$   and   y = -3

3.    Solve   5x - 2y = 4

x - 3y = 6

by multiplying first equation by 1 and second equation by 5, we get

$$5x - 2y = 4 \rightarrow (1)$$

$$5x - 15y = 30 \rightarrow (2)$$

$$13y = -26$$

$$y = -26/13 = -2$$

by substituting it in equation (2)

$$x - 3 \times -2 = 6$$

$$x + 6 = 6, \quad x = 0$$

So , the solution is   x = 0  and  y = -2

4.  Find the equilibrium price and the quantity exchanged at the equation price, if supply and dd functions are given by s = 20 + 3p and D = 160 – 2p, where p is the price charged.

Ans:    s = 20 + 3p

D = 160 – 2p

For Equation   s = D

$$20 + 3p = 160 - 2p$$

$$3p + 2p = 160 - 20$$

$$5p = 140, \quad P = 140/5 = 28$$

Equation price = Rs. 28

Quantity exchanged

$$20 + 3p = 20 + (3 \times 28)$$

$$= 20 + 84 = 104$$

## B. Substitution Method

The substitution method is very useful when one of the equations can easily be solved for one variable. Here we reduce one equation in to the form of $y = f(x)$ or $x = f(y)$. That is expressing the equation either in terms of x or in terms of y. Then substitute this reduced equation in the non-reduced equation and find the values of both unknowns.

### Steps involved in Substitution Method

i.   Choose one equation and isolate one variable; this equation will be considered the first equation.

ii.  Substitute the transformed equation into the second equation and solve for the variable in the equation.

iii. Using the value obtained in step ii, substitute it into the first equation and solve for the second variable.

iv.  Check the obtained values for both variables into both equations.

Solve
$$4x + 2y = 6$$
$$5x + y = 6$$

**Solution:**
$$4x + 2y = 6$$
........................... (1)
$$5x + y = 6$$
........................... (2)

Express equation (2) in terms of x, we will get
$$y = 6 - 5x$$
........................... (3)

Substitute equation (3) in equation (2), we will get
$$4x + 2(6 - 5x) = 6$$

$$4x + 12 - 10x = 6$$

$$-6x = 6 - 12$$

$$x = \frac{-6}{-6} = 1$$

Substitute x = 1 in equation (1)
$$4(1) + 2y = 6$$

$$2y = 6 - 4 \qquad y = \frac{2}{2} = 1$$

,

Checking answers by substituting the obtained value into the original equation.
$$4(1) + 2(1) = 6$$

$$4 + 2 = 6$$

Both sides are equal (L.H.S=R.H.S)

**So the answers are x = 1 and y = 1**

### C. Cross Multiplication

This method is very useful for solving the linear equation in two variables.Let us consider

$$a_1 x + b_1 y + c_1 = 0 \qquad a_2 x + b_2 y + c_2 = 0$$

the general form of two linear equations , and . To solve this pair of equations for x and y using cross-multiplication, we will arrange the variables, coefficients, and the constants as follows.

| X | | Y | | 1 | |
|---|---|---|---|---|---|
| coefficient of y terms | constant | constant terms of x | coefficient | coefficient of x | coefficient of y |
| $b_1$ | $c_1$ | $c_1$ | $a_1$ | $a_1$ | $b_1$ |
| $b_2$ | $c_2$ | $c_2$ | $a_2$ | $a_2$ | $b_2$ |

**That is**

$$x = \frac{b_1 c_2 - b_2 c_1}{a_1 b_2 - a_2 b_1} \qquad\qquad y = \frac{c_1 a_2 - c_2 a_1}{a_1 b_2 - a_2 b_1}$$

Example:      Solve

$$2x + 2y = 40$$
$$3x + 4y = 65$$

**Solution**

On transposition, we get

$$2x + 2y - 40 = 0$$
$$3x + 4y - 65 = 0$$

| X | | Y | | 1 | |
|---|---|---|---|---|---|
| coefficient of y terms | constant | constant terms of x | coefficient | coefficient of x | coefficient of y |
| 2 | -40 | -40 | 2 | 2 | 2 |
| 4 | -65 | -65 | 3 | 3 | 4 |

(2×-65) - (4×-40) = (-130) – (-160) = 30
(-40×3) – (-65×2) = (-120) - (-130) = 10
(2×4) – (3×2) = (8) – (6) = 2

$$x = \frac{30}{2} = 15 \qquad y = \frac{10}{2} = 5$$

S0

,   and

The same answer that we got in the first problem

## Simultaneous Equation in Three Unknowns (First Degree)

Steps
1. Take any two equation form the given equations and eliminate any one of the unknowns.
2. Take the remaining equation and eliminate the same unknown
3. Follow the rules of simultaneous equation in two unknowns

$$9x + 3y - 4z = 35$$

**Examples:**   1.  Solve

$$x + y - z = 4$$

$$2x - 5y - 4z = -48$$

**Solution:**

$$9x + 3y - 4z = 35 .................(1)$$

$$x + y - z = 4 ...............(2)$$

$$2x - 5y - 4z = -48 ...............(3)$$

Take equation (1) and (2)
Multiply equation (2) by 4,we will get

$$4x + 4y - 4z = 16 ...............(4)$$

Subtract it from equation (1), we will get
$$5x - y = 19 ...............(5)$$

Take equation (2) and (3)
Multiply equation (2) by 4,we will get

$$4x + 4y - 4z = 16 ...............(4)$$

Subtract equation (3) from (4), we will get
$$2x + 9y = 64 ...............(6)$$

Take equation (5) and multiply it by 9
$$45x - 45y = 171 ............(7)$$

Add equation (6) from equation (7)

$$45x - 45y = 171 \dots\dots\dots(7)$$

$$2x + 9y = 64 \dots\dots\dots(6)$$

$$47x = 235 \qquad x = \frac{235}{47} = 5$$

Substitute *x=5* in equation (5),
$$5x - y = 19$$

$$5(5) - y = 19$$

$$-y = 19 - 25 = -6$$

So y = 6

Substitute *x=5 and y=6 in equation (1)*
$$9x + 3y - 4z = 35$$

$$9(5) + 3(6) - 4z = 35$$

$$45 + 18 - 4z = 35$$

$$-4z = 35 - 63 = -28$$

$$z = \frac{-28}{-4} = 7$$

**Answer**: *x = 5, y = 6, and z = 7*

2.      Solve   $9x + 3y - 4z = 35$

$$x + y - z = 4$$

$$2x - 5y - 4z + 48 = 0$$

Solution:  $9x + 3y - 4z = 35$ $\qquad \rightarrow \quad$ (1)

$\qquad x + y - z = 4$ $\qquad \rightarrow \quad$ (2)

$\qquad 2x - 5y - 4z + 48 = 0$ $\quad \rightarrow \quad$ (3)

(1) is $\qquad 9x + 3y - 4z = 35$

(2) $\overset{\times}{} \ 9$ $\qquad 9x + 9y - 9z = 36$

$\qquad\qquad$ -6y +5z = -1 $\qquad \rightarrow \qquad$ (4)

(2) $\times$ 2    2x + 2y – 2z = 8

(3) is      3x – 5y – 4z = -48

$$7y + 2z = 56 \quad \rightarrow \quad (5)$$

$$7y + 2z = 56 \quad \rightarrow \quad (5)$$
$$-6y + 5z = -1 \quad \rightarrow \quad (4)$$

(5) $\times$ 5    35y +10z = 280

(4) $\times$ 2    -12y + 10z = -2

$$47y = 282$$

$$y = 282/47 = 6$$

Substituting 6 in equation (4)

$$-6 \times 6 + 5z = -1$$

$$-36 + 5z = -1, \quad 5z = 35, \quad z = 35/5 = 7$$

Substituting y = 6, z = 7 in equ. (2)

$$x +6-7 = 4$$

$$x – 1 = 4, \quad x = 4 + 1 = 5$$

3.    Solve    7x – 4y – 20z = 0
       10x – 13y – 14z = 0
       3x + 4y – 9z = 11

Solution:    7x – 4y – 20z = 0      →      (1)
       10x – 13y – 14z = 0      →      (2)
       3x + 4y – 9z = 11      →      (3)

       7x – 4y – 20z = 0      →    (1)
       10x – 13y – 14z = 0      →    (2)

(1) $\times$ 13    91x – 52y -260z = 0

(2) $\times$ 4    40x – 52y – 56z = 0

$$51x – 204z = 0 \quad \rightarrow \quad (4)$$

7x – 4y – 20z = 0
3x +4y – 9z = 11

$$10x -29z = 11 \quad \rightarrow \quad (5)$$

$$51x - 204z = 0$$
$$10x - 29z = 11$$

(4) $\times$ 10 $\quad 510x - 2040z = 0$

(5) $\times$ 51 $\quad 510x - 1479z = 561$

$$-561z = 561$$

$$z = -561/-561 = 1$$

Substituting z = 1 in equ. (5)

$$10x = 29 \times 1 = 11$$
$$10x - 29 = 11$$
$$10x = 11 + 29, 10x = 40, \qquad x = 40/10 = 4$$

Substituting x = 4, z = 1 in equ. (1)

$$7x - 4y - 20z = 0$$
$$7 \times 4 - 4y - 20 \times 1 = 0$$
$$28 - 4y - 20 = 0$$
$$8 - 4y = 0, \qquad 8 = 4y, y = 8/4 = 2$$

$\therefore$ the solutions are x = 4, y = 2, z = 1

## DEMAND AND SUPPLY FOR A GOOD

Now we can apply simple linear equation and simultaneous linear equations in the analysis of demand and supply. Here we use both demand function and supply function. Demand function depicts the negative relationship between quantity demanded and price. The linear demand function can be written as $q = a - bp$ .where q denotes quantity demanded and p denotes price.

For example: $q = 80 - 2p$ . This equation can be written as $p = 40 - \frac{1}{2}q$ . This called inverse demand function.

Supply function depicts the positive relationship between quantity demanded and price. The linear supply function can be written as $q = a + bp$ .where q denotes quantity supplied and p denotes price.

For example: $q = 40 + 2p$ . This equation can be written as $p = 20 + \frac{1}{2}q$ . This called inverse supply function.

The equilibrium quantity and equilibrium price is determined by the interaction of both demand supply curve. At equilibrium point the demand will be equal to supply. The price that equates demand and supply is called equilibrium price. If current price exceeds the equilibrium price, there will be an excess supply. This situation will compel the producer to reduce the price of the product so that they can sell unsold goods. The reduction in the price will continue until it reaches equilibrium point ($q^d = q^s$) . On the other hand, if current price is below the equilibrium price there is an excess demand for the product. This shortage leads buyers to bid the price up. The increase in the price will continue until it reaches the equilibrium point ($q^d = q^s$).

Now we are able to find the equilibrium price and quantity by using the system of two linear equations; demand function and supply function. Consider the following equations.

$$p = 20 + \frac{1}{2}q$$

$$p = 40 - \frac{1}{2}q$$

This set of equation is system of two linear equations in the variable p and q. We have to find the values of both p and q that satisfy both equations simultaneously.

Example: Find the equilibrium price of the following demand and supply function

$$q^s = 20 + 3p$$

$$q^d = 160 - 2p$$

**Solution:**

At equilibrium demand is equal to supply

$$q^s = 20 + 3p = q^d = 160 - 2p$$

Collect all $p$ values on left side and the constants on right side

$$3p + 2p = 160 - 20$$

$$5p = 140$$

$$p = \frac{140}{5} = 28$$

Now substitute p=28 in either $q^d$ or $q^s$

$$q^s = 20 + 3p$$

$$q^s = 20 + 3(28)$$

$$q^s = 104$$

Check the answer with the $q^d$ equation,

$$q^d = 160 - 2p$$

$$q^d = 160 - 2(28)$$

$$q^d = 160 - 56 = 104$$

Thus, $q^d = q^s$. Here equilibrium price is Rupees 28 and the equilibrium quantity is 104.

## QUADRATIC EQUATIONS

A quadratic function is one which involves at most the second power of the independent

variable in the equation $ax^2 + bx + c$ where a and b are coefficients and c is constant. The graph

of a quadratic function is parabola.

Equation of degree two is known as quadratic equation. This is one of the non-linear

equations. The general format of this equation can be written as $ax^2 + bx + c = 0$. Where *a, b* and

*c* are real numbers and *a* is not equal to zero. The numbers *b* and *c* can also be zero. The number

*a* is the coefficient of $x^2$, *b* is the coefficient of *x*, and *c* is the constant term. These numbers can

be positive or negative.

Solving the quadratic equation, we get the two values for x. These two values are known as the roots of the quadratic equation. It may be pure or general

### Pure quadratic equation

If in the equation $ax^2 + bx + c = 0$, b is zero, then the equation becomes $ax^2 + c = 0$, this is called pure quadratic equation.

### General Quadratic Equation

$ax^2 + bx + c = 0$ is the general form of the quadratic equation.

The general quadratic equation may be solved by one of the following methods.

### Methods to Find the Roots of the Quadratic Equation:

The general quadratic equation $ax^2 + bx + c = 0$ can be solved by one of the following methods

1) By factorization method
2) By quadratic formula
3) By completing the square method

1. **By Factorization Method**

The factorization is an inverse process of multiplication. When an algebraic expression is the product of two or more quantities, each these quantities is called factor. Consider this example, if (x+3) be multiplied by (x+2) the product is $x^2 + 5x + 6$ .The two expressions $x^2 + bx + c$

A. **Procedures to Factorise the Quadratic Equation**

1. Factor the first term ( $x^2$ is the product of x and x)
2. Find two numbers that their sum becomes equal to b (the coefficient of x) and the product becomes equal to c (the constant term)
3. Equate these two expressions with zer0.
4. Apply Zero Property: if we have two expressions multiplied together resulting in zero, then one or both of these must be zero. In other words, if m and n are complex numbers, then m × n= 0, iff m=0 or n=0

Example: Find the roots of $x^2 - 5x + 6 = 0$

Factors of $x^2$ are x and x. Next find two numbers whose sum is -5 and the product is six. The numbers are -2 and -3

$$(x - 3)\ (x - 2) = 0$$

Thus either $(x - 3)$ or $(x - 2)$ should be equal to zero

$(x - 3) = 0$ , $x=3$

$(x - 2) = 0$ $x=2$

$$x^2 - 5x + 6 = 0$$

OR

This equation can rewrite as

$$x^2 - 3x - 2x + 6 = 0$$

-5 broken into two numbers

$$x(x - 3) - 2(x - 3) = 0$$

by factorising the first two terms and last two terms

$$(x - 3)(x - 2) = 0$$

by noting the common factor of $x + 3$

$(x - 3) = 0$ $(x - 2) = 0$

or

*So x=3 or x=2*

2. **Quadratic Formula**

$$ax^2 + bx + c = 0$$

The roots of a quadratic equation can be solved by the following quadratic formula

$$x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

We can split this formula into two parts as

$$\alpha = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$$

, and

$$\beta = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$$

$$\alpha + \beta = -\frac{b}{a} \ and$$

Accordingly, sum of roots:

$$\alpha \times \beta = \frac{c}{a}$$

Product of roots

$$6x^2 - 10x + 4 = 0$$

Example: Find the roots of

Here $a=6$, $b=-10$, and $c=4$

$$\alpha = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$$

$$x = \frac{-(-10) + \sqrt{(-10)^2 - 4 \times 6 \times 4}}{2 \times 6}$$

$$= \frac{10 + \sqrt{100 - 96}}{2 \times 6}$$

$$= \frac{10 + \sqrt{4}}{2 \times 6} \qquad = \frac{10 + 2}{12} = 1$$

$$\beta = \frac{-b - \sqrt{b^2 - 4ac}}{2a}$$

$$x = \frac{-(-10) - \sqrt{(-10)^2 - 4 \times 6 \times 4}}{2 \times 6}$$

$$= \frac{10 - \sqrt{100 - 96}}{2 \times 6}$$

$$= \frac{10 - \sqrt{4}}{2 \times 6} \qquad = \frac{10 - 2}{12} = \frac{8}{12} = \frac{2}{3}$$

$$= \frac{2}{3}$$

Answer: $x=1$ or $x=$

3. **Completing the Square**

This is based on the idea that a perfect square trinomial is the square of a binomial. Consider the following examples:

$$x^2 + 10x + 25$$

is a perfect trinomial because this can be written in the square of a binomial as

$(x+5)^2$ $\qquad x^2 - 6x + 9$ $\qquad\qquad (x-3)^2$

, this equation can be written as

. Consider

Now look at the constant terms of the above two equations, it is the square of half of the coefficient of x equals the constant term;

$$\left(\frac{1}{2} \times 10\right)^2 = 25 \qquad \left(\frac{1}{2} \times (-6)\right)^2 = 9$$

, and . Thus we use this idea in the completing the square method.

**Steps under Completing the Square Method**

1) Rewrite the equation $x^2 + bx - c$ in to $x^2 + bx = c$

2) Add $\left(\frac{1}{2}b\right)^2$ to each side of the equation
3) Factor the perfect-square trinomial
4) Take the square root of both sides of the equation
5) Solve for $x$

**Example:** Solve $x^2 + 6x - 4 = 0$ by completing the square method.

**Solution:** First rewrite the equation as $x^2 + bx = c$

$$x^2 + 6x = 4$$

Add $\left(\frac{1}{2}b\right)^2$ on both sides. Here b = 6 and $\left(\frac{1}{2}b\right)^2 = 3^2 = 9$

$$x^2 + 6x + 9 = 4 + 9$$

$$(x+3)^2 = 13$$

Now take the square root of both sides

$$\sqrt{(x+3)^2} = \sqrt{13}$$

$$(x+3) = \pm\sqrt{13}$$

$$(x = -3 \pm \sqrt{13}$$

**So x=** $-3 + \sqrt{13}$ **or** $-3 - \sqrt{13}$

**OR**

Rewrite the equation so that it becomes complete square. To rewrite the equation take the half of the coefficient of x, add or subtract (depends on the sign of coefficient of x) with the x and square it. Here, $\left(\frac{1}{2}b\right) = 3$

$$(x+3)^2 = x^2 + 6x + 9$$

$$\Rightarrow x^2 + 6x - 4 = (x+3)^2 - 9 - 4$$

Deduct 9 from the expression

$$(x+3)^2 - 13 = 0$$

Take 13 to right side and put square root on both sides

$$\Rightarrow \sqrt{(x+3)^2} = \sqrt{13}$$

$$(x+3) = \pm\sqrt{13}$$

$$(x = -3 \pm \sqrt{13}$$

**So x =** $\dfrac{-3+\sqrt{13}}{}$ **or** $\dfrac{-3-\sqrt{13}}{}$

## SIMULTANEOUS QUADRATIC EQUATIONS

In the second module you have learned simultaneous equations where both equations are linear. In this section we would learn how to solve simultaneous quadratic equation. We start with simultaneous equations where one equation is linear and other is quadratic. This will give you a quadratic equation to solve.

**Example**: solve simultaneous equations

$$y = x^2 - 1$$

$$y = 5 - x$$

**Solution:**

$$y = x^2 - 1 ........................(1)$$

$$y = 5 - x ..........................(2)$$

Subtract equation (2) from (1)

$$(y = x^2 - 1) \quad (y = 5 - x) \quad x^2 - 1 - 5 + x$$
$$- \qquad = \qquad\qquad y \text{ will be cancelled}$$
$$x^2 + x - 6 = 0$$

Now solve this quadratic equation either by factorisation method or by quadratic formula.

$$(x+3)\ (x-2) = 0$$

By factorization

$$x + 3 = 0 \qquad x - 2 = 0$$

So $\qquad$ or $\qquad$ Therefore, $\qquad$ _x=-3 or x= 2_

OR

Substitute equation (2) in equation (1)

$$\Rightarrow x^2 - 1 = 5 - x$$

$$x^2 + 1 - 5 + x \quad x^2 + x - 6 = 0$$
$$=$$
$$(x+3)\ (x-2) = 0$$

By factorisation

$$x + 3 = 0 \qquad x - 2 = 0$$

So $\qquad$ or $\qquad$ Therefore, $\qquad$ _x=-3 or x= 2_

Now we can move to simultaneous quadratic equations

Solve simultaneous quadratic equations

$$y = 2x^2 + 3x + 2$$

$$y = x^2 + 2x + 8$$

**Solution:**

$$y = 2x^2 + 3x + 2..............(1)$$

$$y = x^2 + 2x + 8...............(2)$$

Now equate equation (1) and equation (2)

$$2x^2 + 3x + 2 = x^2 + 2x + 8$$

$$2x^2 + 3x + 2 - x^2 - 2x - 8 = 0$$

$$x^2 + x - 6 = 0$$

$$(x + 3)(x - 2) = 0$$

By factorization

$$x + 3 = 0 \qquad x - 2 = 0$$

So $\qquad$ or $\qquad$ Therefore, $\qquad$ <u>*x=-3 or x= 2*</u>

## ECONOMIC APPLICATION

The quadratic equation has application in the field of economics. Here we discuss two important Economics application of quadratic equation.

**Supply and Demand**

The quadratic equation can be used to represent supply and demand function. Market equilibrium occurs when the quantity demanded equals the quantity supplied. If we solve the system of quadratic equations for quantity and price we get equilibrium quantity and price.

$$p = q^2 + 50$$

For example: The supply function for a commodity is given by $\qquad$ and the demand

$$p = -10q + 650$$

function is given by $\qquad$ find the point of equilibrium.

Solution: $\qquad$ At the equilibrium demand is equal to supply

$$q^2 + 50 = -10q + 650$$

$$q^2 + 50 + 10q - 650 = 0$$

$$q^2 + 10q - 600 = 0$$

$$(q + 30)(q - 20) = 0$$

By factorization

So q=-30 or 20

Since negative quantity is not possible we take positive value as quantity. Thus the equilibrium quantity is 20. Put q=20 in either demand function or supply function.

$$p = q^2 + 50$$

Supply function

$$p = (20)^2 + 50$$

P=450

## Cost and Revenue

The cost and revenue function can be represented by the quadratic equation. The total cost is composed of two parts, fixed cost and variable cost. The fixed cost remains the same regardless of the number of units produced. It does not depend on the quantity produced. Rent on building and machinery is an example for the fixed cost. The variable cost is directly related to the number of unit produced. Cost on raw material is an example for the variable cost. Thus, TC=FC+VC
The revenue of the firm depends on the number of unit sold and its price.
TR= P×Q. Where TR denotes total revenue, P shows price, and Q denotes quantity.

## BREAK-EVEN POINT

Firm's break-even point occurs when total revenue is equal to total cost.
**Steps:** 1- Find the profit function
2- Equate profit function with zero and solve for q.
If we deduct total cost function from total revenue function we get profit function.

$$TC = 10.75q^2 + 5q + 125$$

**Example**: A firm has the total cost function

$$p = 180 - 0.5q$$

and demand function                                  Find revenue function, profit function, and break-even

point .

**Solution:**

Total revenue function= price × quantity (TR = p × q)
$$p \times q = (180 - 0.5q)q$$

$$= 175q - 0.5q^2$$

$$(\pi = TR - TC)$$

Profit function= Total revenue- total cost

$$= 175q - 0.5q^2 - 10.75q^2 + 5q + 125$$

$$= 180q - 11.25q^2 - 125$$

$$11.25q^2 - 180q - 125$$

Break –even point
$$11.25q^2 - 180q - 125 = 0$$

Use quadratic formula

$$q = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

Here a=11.25, b=-180, and c= -125

$$q = \frac{-(-180) \pm \sqrt{(-180)^2 - 4 \times 11.25 \times -125}}{2 \times 11.25}$$

$$= \frac{180 \pm \sqrt{32400 + 5625}}{22.5}$$

$$= \frac{180 \pm \sqrt{38025}}{22.5}$$

$$= \frac{180 \pm 195}{22.5}$$

$$= \frac{180 + 195}{22.5} = 16.66 \approx 17$$

$$= \frac{180 - 195}{22.5} = -0.67$$

Since negative quantity is not possible we take positive value as quantity. Thus the break-even point is 17.

## MODULE  II

## BASIC MATRIX ALGEBRA

### MATRICES : DEFINITION AND TERMS

A matrix is defined as a rectangular array of numbers, parameters or variables.  Each of which has a carefully ordered place within the matix.  The members of the array are referred to as "elements" of the matrix and are usually enclosed in brackets, as shown below.

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$$

The members in the horizontal line are called rows and members in the vertical line are called columns. The number of rows and the number of columns together define the dimension or order of the matrix. If a matrix contains 'm' rows and 'n' columns, it is said to be of dimension m x n (read as '). The row number precedes the column number. In that sense the above matrix is of dimension 3 x 3. Similarly

$$B = \begin{bmatrix} 3 & 5 & 1 \\ 2 & 7 & 4 \end{bmatrix}_{2 \times 3}$$

$$C = \begin{bmatrix} 7 \\ 8 \\ 10 \end{bmatrix}_{3 \times 1}$$

$$D = \begin{bmatrix} 10 & 2 \end{bmatrix}_{1 \times 2}$$

$$E = \begin{bmatrix} 2 & 0 \\ 1 & 4 \end{bmatrix}_{2 \times 2}$$

## TYPES OF MATRICES

### 1. Square Matrix
A matrix with equal number of rows and colums is called a square matrix. Thus, it is a special case where m=n. For example

$\begin{bmatrix} 2 & 1 \\ 3 & 4 \end{bmatrix}$ is a square matrix of order 2

$\begin{bmatrix} 2 & 1 & 3 \\ 4 & 0 & 6 \\ 9 & 7 & 5 \end{bmatrix}$ is a square matrix of order 3

### 2. Row matrix or Row Vector
A matrix having only one row is called row vector of row matrix. The row vector will have a dimension of 1×0. For example

$\begin{bmatrix} 2 & 5 & 0\,1 \end{bmatrix}_{1 \times 4}$

$\begin{bmatrix} 2 & 1 \end{bmatrix}_{1 \times 2}$

$\begin{bmatrix} 0 & 2 & 3 \end{bmatrix}_{1 \times 3}$

### 3. Column matrix or Column Vector

A matrix having only one column is called column vector or column matrix. The column vector will have a dimension of m 1 . For example

$$\begin{bmatrix} 5 \\ 8 \end{bmatrix}_{2 \times 1} \qquad \begin{bmatrix} ¿ \\ 8 \\ 9 \\ 21 \\ 4 \end{bmatrix}_{4 \times 1}$$

$$\begin{bmatrix} 0 \\ 2 \\ 5 \end{bmatrix}_{3 \times 1}$$

## 4. Diagonal Matrix

In a matrix the elements lie on the diagonal from left top to the right bottom are called diagonal elements. For instance, in the matrix $\begin{bmatrix} 2 & 5 \\ 4 & 6 \end{bmatrix}$ the element 2 and 6 are diagonal elements. A square matrix in which all elements except those in diagonal are zero are called diagonal matrix. For example

$$\begin{bmatrix} 2 & 0 \\ 0 & 6 \end{bmatrix} \quad ¿2 \times 2 \qquad \begin{bmatrix} 4 & 0 & 0 \\ 0 & 9 & 0 \\ 0 & 0 & 2 \end{bmatrix}_{3 \times 3}$$

## 5. Identity matrix or Unit Matrix

A diagonal matrix in which each of the diagonal elements is unity is said to be unit matrix and denoted by I. The identity matrix is similar to the number one in algebra since multiplication of a matrix by an identity matrix leaves the original matrix unchanged. That is, AI = I A =A

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \quad 2 \times 2 \qquad \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \quad 3 \times 3 \quad \text{are examples of identity matrix}$$

## 6. Null Matrix or Zero Matrix

A matrix in which every element is zero is called null matrix or zero matrix. It is not necessarily square. Addition or subtraction of the null matrix leaves the original matrix unchanged and multiplication by a null matrix produces a null matrix.

$$\begin{bmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix}_{2 \times 3} \qquad \begin{bmatrix} 0 & 0 \\ 0 & 0 \end{bmatrix} \quad 2 \times 2 \qquad \begin{bmatrix} 0 & 0 \\ 0 & 0 \\ 0 & 0 \end{bmatrix}_{3 \times 2} ¿ \quad \text{are examples of null matrix}$$

## 7. Triangular Matrix

If every element above or below the leading diagonal is zero, the matrix is called a triangular matrix. Triangular matrix may be upper triangular or lower triangular. In the upper triangular matrix, all elements below the leading diagonal are zero, like

$$A = \begin{bmatrix} 1 & 9 & 2 \\ 0 & 3 & 7 \\ 0 & 0 & 4 \end{bmatrix}$$

In the lower triangular matrix, all elements above leading diagonal are zero like

$$B = \begin{bmatrix} 4 & 0 & 0 \\ 2 & 9 & 0 \\ 5 & 6 & 3 \end{bmatrix}$$

## 8. Idempotent Matrix

A square matrix A is said to be idempotent if $A = A^2$.

## TRANSPOSE OF A MATRIX

A matrix obtained from any given matrix A by interchanging its rows and columns is called its transpose and is denoted by or A'. If A is $m \times n$ matrix A' will be $n \times m$ dimension. For example

$$A = \begin{bmatrix} 6 & 7 & 9 \\ 2 & 8 & 4 \end{bmatrix} \quad \dot{c} 2 \times 3 \qquad\qquad A^t = \begin{bmatrix} 6 & 2 \\ 7 & 8 \\ 9 & 4 \end{bmatrix}_{3 \times 2}$$

$$B = \begin{bmatrix} 1 & 23 & 4 \\ 2 & 34 & 1 \\ 3 & 42 & 5 \end{bmatrix}_{3 \times 4} \qquad B^t = \begin{bmatrix} 123 \\ 234 \\ 342 \\ 415 \end{bmatrix}_{4 \times 3}$$

$$C = \begin{bmatrix} 12 \\ 19 \\ 25 \end{bmatrix}_{3 \times 1} \qquad C^t = \begin{bmatrix} 12 & 19 & 25 \end{bmatrix}_{1 \times 3}$$

$$D = \begin{bmatrix} 21 \\ 78 \\ 30 \\ 95 \end{bmatrix} \qquad D^t = \dot{c} \begin{bmatrix} 2 & 73 & 9 \\ 1 & 80 & 5 \end{bmatrix}$$

## Symmetric and skew Symmetric Matrix

Any square matrix A is said to be symmetric if it is equal to its transpose. That is, A is symmetric if $A = A^t$ Consider the following examples

$$A = \begin{bmatrix} 1 & 5 \\ 5 & 3 \end{bmatrix} \qquad A^t = \begin{bmatrix} 1 & 5 \\ 5 & 3 \end{bmatrix} \quad A = A^t \text{ , hence A is symmetric}$$

$$B = \begin{bmatrix} 5 & 2 & 6 \\ 2 & 3 & 9 \\ 6 & 9 & 7 \end{bmatrix} \qquad B^t = \begin{bmatrix} 5 & 2 & 6 \\ 2 & 3 & 9 \\ 6 & 9 & 7 \end{bmatrix} \quad B = B^t \quad \therefore \text{ B is symmetric}$$

At the same time, any square matrix A is said to be skew symmetric if it is equal to its negative transpose. That is $A = -A^t$, then A is skew symmetric consider the following examples

$$A = \begin{bmatrix} 0 & 4 \\ -4 & 0 \end{bmatrix} \qquad A^t = \begin{bmatrix} 0 & 4 \\ -4 & 0 \end{bmatrix} \qquad -A^t = \begin{bmatrix} 0 & 4 \\ -4 & 0 \end{bmatrix}$$

$$A = -A^t \qquad \therefore \quad \text{A is skew symmetric}$$

$$B = \begin{bmatrix} 0 & 3 & 5 \\ -3 & 0 & -2 \\ -5 & 2 & 0 \end{bmatrix} \qquad B^t = \begin{bmatrix} 0 & -3 & -5 \\ 3 & 0 & 2 \\ 5 & -2 & 0 \end{bmatrix} \qquad -B^t = \begin{bmatrix} 0 & 3 & 5 \\ -3 & 0 & -2 \\ -5 & 2 & 0 \end{bmatrix}$$

## OPERATION OF MATRICES

### 1. Addition and subtraction of Matrices

Two matrixes can be added or subtracted if and only if they have the same dimension. That is, given two matrixes A and B, their addition or subtraction that is, A + B and A − B requires that A and B have the same dimension. When this dimensional requirement is met, the matrices are said to be "conformable for addition or subtraction". Then, each element of one matrix is added to (or subtracted from) the corresponding element of the other matrix.

For example, if $A = \begin{bmatrix} 4 & 9 \\ 2 & 1 \end{bmatrix}_{2\times 2}$ and $B = \begin{bmatrix} 6 & 3 \\ 7 & 0 \end{bmatrix}_{2\times 2}$

$$A + B = \begin{bmatrix} 4 & 9 \\ 2 & 1 \end{bmatrix} + \begin{bmatrix} 6 & 3 \\ 7 & 0 \end{bmatrix} = \begin{bmatrix} 4+6 & 9+3 \\ 2+7 & 1+0 \end{bmatrix} = \begin{bmatrix} 10 & 12 \\ 9 & 1 \end{bmatrix}$$

$$A - B = \begin{bmatrix} 4 & 9 \\ 2 & 1 \end{bmatrix} - \begin{bmatrix} 6 & 3 \\ 7 & 0 \end{bmatrix} = \begin{bmatrix} 4-6 & 9-3 \\ 2-7 & 1-0 \end{bmatrix} = \begin{bmatrix} -2 & 6 \\ -5 & 1 \end{bmatrix}$$

**Example :2**

If $A = \begin{bmatrix} 8 & 9 & 7 \\ 3 & 6 & 2 \\ 4 & 5 & 10 \end{bmatrix}$ $B = \begin{bmatrix} 1 & 3 & 6 \\ 5 & 2 & 4 \\ 7 & 9 & 2 \end{bmatrix}$ $A+B = \begin{bmatrix} 9 & 12 & 13 \\ 8 & 8 & 6 \\ 11 & 14 & 12 \end{bmatrix}$

**Example : 3**

$A = \begin{bmatrix} 3 & 7 & 11 \\ 12 & 9 & 2 \end{bmatrix}$ $B = \begin{bmatrix} 6 & 8 & 1 \\ 9 & 5 & 8 \end{bmatrix}$ $A - B = \begin{bmatrix} -3 & -1 & 10 \\ 3 & 4 & -6 \end{bmatrix}$

**Example : 4**

$$A = \begin{bmatrix} 2 & 2 & 2 \\ 1 & 1 & -3 \\ 1 & 0 & 4 \end{bmatrix} \quad B = \begin{bmatrix} 3 & 3 & 3 \\ 3 & 0 & 5 \\ 6 & 9 & -1 \end{bmatrix} \quad C = \begin{bmatrix} 4 & 4 & 4 \\ 5 & -1 & 0 \\ 2 & 3 & 1 \end{bmatrix}$$

$$A + B - C = \begin{bmatrix} 1 & 1 & 1 \\ -1 & 2 & 2 \\ 5 & 6 & 2 \end{bmatrix}$$

## Example 5

$$A = \begin{bmatrix} 12 & 16 & 27 & 8 \end{bmatrix} \quad B = \begin{bmatrix} 0 & 19 & 5 & 6 \end{bmatrix} \quad A + B = \begin{bmatrix} 12 & 17 & 11 & 12 & 14 \end{bmatrix}$$

## 2. Scalar Multiplication

In the matrix algebra, a simple number such as 1,2, -1, -2 ……. is called a scalar. Multiplication of matrix by a scalar or number involves multiplication of every element of the matrix by the number.  The process is called scalar multiplication.

Let 'A' be any matrix and 'k' any scalar, then the matrix obtained by multiplying every element of A by K is said to be the scalar multiple of A by K, because it scales the matrix up or down according to the size of the scalar.

## Example 1

If $A = \begin{bmatrix} 3 & -1 \\ 0 & 5 \end{bmatrix}$ and scalar k = 7 then $KA = 7 \begin{vmatrix} 3 & -1 \\ 0 & 5 \end{vmatrix} = \begin{bmatrix} 21 & -7 \\ 0 & 35 \end{bmatrix}$

## Example 2

Determine KA if K = 4 and $A = \begin{bmatrix} 3 & 2 \\ 9 & 5 \\ 6 & 7 \end{bmatrix} \quad KA = \begin{bmatrix} 12 & 8 \\ 36 & 20 \\ 24 & 28 \end{bmatrix}$

## Example 3

K = -2 and $A = \begin{bmatrix} 7 & -3 & 2 \\ -5 & 6 & 8 \\ 2 & -7 & -9 \end{bmatrix} \quad KA = \begin{bmatrix} -14 & 6 & -4 \\ 10 & -12 & -16 \\ -4 & 14 & 18 \end{bmatrix}$

## Example 4

If $A = \begin{bmatrix} 2 & 3 & 1 \\ 0 & -1 & 5 \end{bmatrix} \quad B = \begin{bmatrix} 1 & 2 & -1 \\ 0 & -1 & 3 \end{bmatrix}$ Find 2A -3B

$2A = \begin{bmatrix} 4 & 6 & 2 \\ 0 & -2 & 10 \end{bmatrix} \quad 3B = \begin{bmatrix} 3 & 6 & -3 \\ 0 & -3 & 9 \end{bmatrix} \quad 2A - 3B = \begin{bmatrix} 1 & 0 & 5 \\ 0 & 1 & 1 \end{bmatrix}$

## 3. Vector Multiplication

Multiplication of a row vector 'A' by a column vector 'B' requires that each vector has precisely the same number of elements. The product is found by multiplying the individual elements of the row vector by their corresponding elements in the column vector and summing the product. For example

If A = [a b c]  B = $\begin{bmatrix} d \\ e \\ f \end{bmatrix}$

AB = [ $ad + be + cf$ ]

Thus the product of row – column multiplication will be a single number or scalar. Row –column vector multiplication is very important because it serves the basis for all matrix multiplication.

**Example 1:**     A = [ 4 7 2 9 ] B =     AB = (4 $\times$ 12) + (7 $\times$ 1) + (2 $\times$ 5) + (9×0) = 119

**Example 2 :**     C = [ 3 6 8 ] D = $\begin{bmatrix} 2 \\ 4 \\ 5 \end{bmatrix}$   CD = 70

**Example 3:**     A = [ 12 -5 6 11 ] B = AB = 44

**Example 4:**     A = [ 9 6 2 0 -5 ] B AB = 101

### 4. Matrix Multiplication

The matrices A and B are conformable for multiplication if and only if the number of columns in the matrix A is equal to the number of rows in the matrix B. That is, to find the product AB, conformity condition for multiplication requires that the column dimension of A (the lead matrix in the expression AB) must be equal to the row dimension of B (the lag matrix)

In general, if A is of the order $m \times n$ then B should be of the order $n \times p$ and dimension of AB will be $m \times p$. That is, if dimension of A is and $1 \times 2$ and dimension of B is $2 \times 3$, then AB will be of $1 \times 3$ dimension. For multiplication, the procedure is that take each row and multiply with all column. For example if

$$A = \begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix} \text{ and B} = \begin{bmatrix} b_{11} & b_{12} & b_{13} \\ b_{21} & b_{22} & b_{23} \\ b_{31} & b_{32} & b_{33} \end{bmatrix}$$

$$AB = \begin{bmatrix} a_{11}b_{11}+a_{12}b_{21}+a_{13}b_{31} & a_{11}b_{12}+a_{12}b_{22}+a_{13}b_{32} & a_{11}b_{13}+a_{11}b_{13}+a_{13}b_{33} \\ a_{21}b_{11}+a_{22}b_{21}+a_{23}b_{31} & a_{21}b_{12}+a_{22}b_{22}+a_{23}b_{32} & a_{21}b_{13}+a_{22}b_{23}+a_{23}b_{33} \\ a_{31}b_{11}+a_{32}b_{21}+a_{33}b_{31} & a_{31}b_{12}+a_{32}b_{22}+a_{33}b_{32} & a_{31}b_{13}+a_{32}b_{23}+a_{33}b_{33} \end{bmatrix}$$

Similarly if A = $\begin{bmatrix} 3 & 6 & 7 \\ 12 & 9 & 11 \end{bmatrix}$ B = $\begin{bmatrix} 6 & 12 \\ 5 & 10 \\ 13 & 2 \end{bmatrix}$

Since A is of 2 × 3 dimension and B is of 3 × 2 dimension the matrices are conformable for multiplication and the product AB will be of 2 × 2 dimension. Then

AB = $\begin{bmatrix} 3 \times 6 + 6 \times 5 + 7 \times 13 & 3 \times 12 + 6 \times 10 + 7 \times 2 \\ 12 \times 6 + 9 \times 5 + 11 \times 13 & 12 \times 12 + 9 \times 10 + 11 \times 2 \end{bmatrix}$

AB = $\begin{bmatrix} 139 & 110 \\ 260 & 256 \end{bmatrix}$

**Example 1**

A = $\begin{bmatrix} 3 & 5 \\ 4 & 6 \end{bmatrix}$ B = $\begin{bmatrix} -1 & 0 \\ 4 & 7 \end{bmatrix}$ AB = $\begin{bmatrix} 17 & 35 \\ 20 & 42 \end{bmatrix}$

**Example 2 :** A = $\begin{bmatrix} 1 & 3 \\ 2 & 8 \\ 4 & 0 \end{bmatrix}$ B = $\begin{bmatrix} 5 \\ 9 \end{bmatrix}$ AB = $\begin{bmatrix} 32 \\ 82 \\ 20 \end{bmatrix}$

**Example 3:** A = $\begin{bmatrix} 7 & 11 \\ 2 & 9 \\ 10 & 6 \end{bmatrix}$ B = $\begin{bmatrix} 12 & 4 & 5 \\ 3 & 6 & 1 \end{bmatrix}$ AB = $\begin{bmatrix} 117 & 94 & 46 \\ 51 & 62 & 19 \\ 138 & 76 & 56 \end{bmatrix}$

**Example 4**

A = B =[ 2 6 5 3 ] B = $\begin{bmatrix} 2 & 65 & 3 \end{bmatrix}$ AB = $\begin{bmatrix} 6 & 18 & 15 & 9 \\ 2 & 6 & 5 & 3 \\ 8 & 24 & 20 & 12 \\ 10 & 30 & 25 & 15 \end{bmatrix}$

**Matrix Expression of a System of Linear Equations**

Matrix algebra permits the concise expression of a system of linear equations. For example, the following system of linear equation

$$a_{11}x_1 + a_{12}x_2 = b_1$$

$$a_{21}x_1 + a_{22}x_2 = b_2$$

Can be expressed in matrix form as

$A X$ =B where,

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix} \quad X = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad \text{and } B = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

Here, A is the coefficient matrix, x is the solution vector an B is the vector of constant terms. $X$ and B will always be column vector. Since A is 2x2 matrix and x is 2x1 vector, they we conformable for multiplication, and the product matrix will be 2 x 1.

**Example 1** :

$$7x_1 + 3x_1 = 45$$

$$4x_1 + 5x_2 = 29$$

In matrix from $AX=B$

$$\begin{bmatrix} 7 & 3 \\ 4 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 45 \\ 29 \end{bmatrix}$$

**Example 2:**

$$7x_1 + 8x_2 = 120$$

$$6x_1 + 9x_2 = 92$$

In matrix form $AX=B$

$$\begin{bmatrix} 7 & 8 \\ 6 & 9 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 120 \\ 92 \end{bmatrix}$$

**Example 3:**    $2x_1 + 4x_2 + 9x_2 = 143$
$2x_1 + 8x_2 + 7x_3 = 204$
$5x_1 + 6x_2 + 3x_3 = {}^-168$

In matrix form

$AX = B$

$$\begin{bmatrix} 2 & 4 & 9 \\ 2 & 8 & 7 \\ 5 & 6 & 3 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 143 \\ 204 \\ -168 \end{bmatrix}$$

**Example 4**

$8w + 12x - 7y + 22 = 139$

$3w - 13x + 4y + 92 = 242$

In matrix from

$AX=B$

$$\begin{bmatrix} 8 & 12 & -7 & 2 \\ 3 & -13 & 4 & 9 \end{bmatrix} \begin{bmatrix} w \\ x \\ y \\ z \end{bmatrix} = \begin{bmatrix} 139 \\ 242 \end{bmatrix}$$

**Concept of Determinants**

The determinant is a single number or scalar associated with a square matrix. Determinants are defined only for square matrix. In other words, determinant denoted as $|A|$, is a uniquely defined number or scalar associated with that matrix

If A = $[a_{11}]$ is a 1×1 matrix, then the determinant of A, ie $|A|$ is the number $a_{11}$ itself. If A is a 2 × 2 matrix then the determinant of such matrix, like

$$A = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$$

called the second order determinant is derived by taking the product of two elements on the principal diagonal and subtracting from it the product of two elements off the principal diagonal. That is,

$$|A| = a_{11}a_{22} - a_{21}a_{12}$$

Thus $|A|$ is obtained by cross multiplication of the elements. If the determinant is equal to zero, the determinant is said to vanish and the matrix is termed as singular matrix. That is, a singular matrix is one in which there exists linear dependence between at least two rows or columns. If $|A| \neq 0$, matrix A is non-singular and all its rows and columns are linearly independent.

**Example:** If A = $\begin{bmatrix} 10 & 4 \\ 8 & 5 \end{bmatrix}$ $|A|$ = (10 × 5) − (4 × 8) = 18

**Example 2:** B = $\begin{bmatrix} 2 & 1 \\ -3 & 2 \end{bmatrix}$ $|B|$ = 7

**Example 3:** C = $\begin{bmatrix} 6 & 4 \\ 7 & 9 \end{bmatrix}$ $|C|$ = 26

**Example 4:** D = $\begin{bmatrix} 4 & 6 \\ 6 & 9 \end{bmatrix}$ $|D|$ = 0

**RANK OF MATRIX**

The rank (P) of a matrix is defined as the maximum number of linearly independent rows and columns in the matrix. For example, if

$$A = \begin{bmatrix} 2 & 3 \\ 3 & 6 \end{bmatrix}$$

$|A|$ = 3 and the matrix A is non singular and its rows and columns are linearly independent and the rank of the matrix A, ie, P(A) = 2.  If

$$B = \begin{bmatrix} 4 & 2 \\ 8 & 4 \end{bmatrix}$$

$|B|$ = 0 and matrix B is singular and a Linear dependence exists between its rows and columns.  Hence the rank of the matrix P(B) = 1

**Third order Determinants**

A determinant of order three is associated with a 3 × 3 matrix.  Given.

$$A = \begin{bmatrix} a11 & a12 & a13 \\ a21 & a22 & a23 \\ a31 & a32 & a33 \end{bmatrix}$$

Then

$$|A| = a_{11}\begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix} - a_{12}\begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix} + a_{13}\begin{vmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{vmatrix}$$

$$|A| = a_{11}(a_{22}a_{33} - a_{32}a_{23}) - a_{12}(a_{21}a_{33} - a_{31}a_{23}) + a_{13}(a_{21}a_{32} - a_{31}a_{22})$$

$|A|$ = a scalar

$|A|$ is called a third order determinant and is the summation of three products to desire three products.

1. Take the first element of the first row, ie, $a_{11}$ and mentally delete the row and column in which it appears.  Then multiply $a_{11}$ by the determinant of the remaining elements.

2. Take the second element of the first row, ie, a12 and mentally delete the row and column in which it appears.  Then multiply $a_{12}$ by -1 time the determinant of the remaining element.

3. Take the third element of the first row, ie, $a_{13}$ and mentally delete the row and column in which it appears.  Then multiply by the determinant of the remaining elements.

In the like manner, the determinant of a 4 × 4 matrix is the sum of four products.  The determinant of a 5 × 5 matrix is the sum of five products and so on.

**Example 1**

$$A = \begin{vmatrix} 8 & 3 & 2 \\ 6 & 4 & 7 \\ 5 & 1 & 3 \end{vmatrix}$$

$$|A| = 8 \begin{vmatrix} 4 & 7 \\ 1 & 3 \end{vmatrix} - 3 \begin{vmatrix} 6 & 7 \\ 5 & 3 \end{vmatrix} + 2 \begin{vmatrix} 6 & 4 \\ 5 & 1 \end{vmatrix}$$

$$|A| = (8 \times 5) - (3 \times {}^-17) + 2({}^-4)$$

$$|A| = 63$$

**Example 2**

$$A = \begin{bmatrix} -3 & 6 & 2 \\ 2 & 1 & 8 \\ 7 & 9 & 1 \end{bmatrix} \qquad |A| = 3 \begin{vmatrix} 1 & 8 \\ 9 & 1 \end{vmatrix} - 6 \begin{vmatrix} 2 & 8 \\ 7 & 1 \end{vmatrix} + 5 \begin{vmatrix} 2 & 1 \\ 7 & 9 \end{vmatrix}$$

$$|A| = 166$$

**Example 3**

$$B = \begin{bmatrix} -3 & 6 & 2 \\ 1 & 5 & 4 \\ 4 & -8 & 2 \end{bmatrix} \qquad |B| = -3 \begin{vmatrix} 5 & 4 \\ -8 & 2 \end{vmatrix} - 6 \begin{bmatrix} 1 & 4 \\ 4 & 2 \end{bmatrix} + 2 \begin{vmatrix} 1 & 5 \\ 4 & -8 \end{vmatrix}$$

$$|B| = 98$$

**Example 4**

$$C = \begin{bmatrix} 5 & 7 & 2 \\ 2 & 3 & 1 \\ 4 & 6 & 2 \end{bmatrix} \qquad |C| = 5 \begin{vmatrix} 3 & 1 \\ 6 & 2 \end{vmatrix} - 7 \begin{vmatrix} 2 & 1 \\ 4 & 2 \end{vmatrix} + 2 \begin{vmatrix} 2 & 3 \\ 4 & 6 \end{vmatrix}$$

$$|C| = 0$$

**PROPERTIES OF A DETERMINANT**

1. The value of the determinant does not change if the rows and columns of it are interchanged. That is, the determinant of a matrix equals the determinant of its transpose. That is  .For Example

$$A = \begin{bmatrix} 4 & 3 \\ 5 & 6 \end{bmatrix} \qquad |A| = 9 \qquad A^t = \begin{bmatrix} 4 & 5 \\ 3 & 6 \end{bmatrix} \qquad |A^t| = ¿ \; 9$$

2. The interchange of any two rows or any two columns will alter the sign,  but not the numerical value of the determinant.  For example, if

$$A = \begin{bmatrix} 3 & 1 & 0 \\ 7 & 5 & 2 \\ 1 & 0 & 3 \end{bmatrix} \quad |A| = 3 \begin{vmatrix} 5 & 2 \\ 0 & 3 \end{vmatrix} - 1 \begin{vmatrix} 7 & 2 \\ 1 & 3 \end{vmatrix} + 0 \begin{vmatrix} 7 & 5 \\ 1 & 0 \end{vmatrix} = 26$$

Now if we interchange first and third column,

$$\begin{bmatrix} 0 & 1 & 3 \\ 2 & 5 & 7 \\ 3 & 0 & 1 \end{bmatrix} \quad 0 \begin{vmatrix} 5 & 7 \\ 0 & 1 \end{vmatrix} - 1 \begin{vmatrix} 2 & 7 \\ 3 & 1 \end{vmatrix} + 3 \begin{vmatrix} 2 & 5 \\ 3 & 0 \end{vmatrix} = -26 = - |A|$$

3. If any two rows or columns of a matrix are identical or proportional, ie linearly dependent, the determinant is zero. For Example

$$A = \begin{bmatrix} 2 & 3 & 1 \\ 4 & 1 & 0 \\ 2 & 3 & 1 \end{bmatrix} \quad |A| = 2 \begin{vmatrix} 1 & 0 \\ 3 & 1 \end{vmatrix} - 3 \begin{vmatrix} 4 & 0 \\ 2 & 1 \end{vmatrix} + 1 \begin{vmatrix} 4 & 1 \\ 2 & 3 \end{vmatrix}$$

$|A|$ =0, since first and third row are identical

4. The multiplication of any one row or one column by a scalar or constant 'k' will change the value of the determinant k . For example

$$\text{If } A = \begin{bmatrix} 3 & 5 & 7 \\ 2 & 1 & 4 \\ 4 & 2 & 3 \end{bmatrix} \quad |A| = 35$$

Now forming a new matrix B by multiplying the first row of A by 2, then

$$B = \begin{bmatrix} 6 & 10 & 14 \\ 2 & 1 & 4 \\ 4 & 2 & 3 \end{bmatrix} \quad |B| = 70, \text{ ie, } 2 \times |A|$$

Thus, multiplying a single row or column of a matrix by a scalar will cause the value of determinant to be multiplied by the scalar.

5. The determinant of triangular matrix is equal to the product of elements on the principal diagonal For example, for the following lower triangular matrix

$$A = \begin{bmatrix} -3 & 0 & 0 \\ 2 & -5 & 0 \\ 6 & 1 & 4 \end{bmatrix} \quad |A| = 60, \text{ ie } ^-3 \times {}^-5 \times 4$$

6. If all the elements of any row or column are zero the determinant is zero. For example

$$A = \begin{bmatrix} 12 & 16 & 13 \\ 0 & 0 & 0 \\ -15 & 20 & -9 \end{bmatrix}$$

$|A|$ = 0 Since all elements of second row is zero

7. If every element in a row or column of a matrix is sum of two numbers, then the given determinant can be expressed as the sum of two determinants.

$$|A| = \begin{bmatrix} 2+3 & 1 \\ 4+1 & 5 \end{bmatrix} = 20$$

$$\text{ie } \begin{vmatrix} 2 & 1 \\ 4 & 5 \end{vmatrix} + \begin{vmatrix} 3 & 1 \\ 1 & 5 \end{vmatrix} = 6 + 14 = 20$$

8. Addition or subtraction of a non zero multiple of any one row or column from another row or column does not change the value of determinant. For example

$$A = \begin{bmatrix} 20 & 3 \\ 10 & 2 \end{bmatrix} \qquad |A| = 10$$

Now subtract tow times of second column from first column and for a new matrix.

$$B = \begin{vmatrix} 14 & 3 \\ 6 & 2 \end{vmatrix} \qquad |B| = 10$$

## MINORS AND COFACTORS

Every element of a square matrix has a minor.  It is the value of the determinant formed with the elements obtained when the row and the column in which the element lies are deleted. Thus, a minor, denoted as is the determinant of the sub matrix formed by deleting the i$^{th}$ row and j$^{th}$ column of the matrix.

For example, if A = $\begin{bmatrix} a_{11} & a_{12} & a_{13} \\ a_{21} & a_{22} & a_{23} \\ a_{31} & a_{32} & a_{33} \end{bmatrix}$

Minor of $a_{11} = \begin{vmatrix} a_{22} & a_{23} \\ a_{32} & a_{33} \end{vmatrix}$   Minor of $a_{12} = \begin{vmatrix} a_{21} & a_{23} \\ a_{31} & a_{33} \end{vmatrix}$   Minor of $a_{13} = \begin{bmatrix} a_{21} & a_{22} \\ a_{31} & a_{32} \end{bmatrix}$

Minor of $a_{21} = \begin{vmatrix} a_{12} & a_{13} \\ a_{32} & a_{33} \end{vmatrix}$   Minor of $a_{22} = \begin{vmatrix} a_{11} & a_{13} \\ a_{31} & a_{33} \end{vmatrix}$   Minor of $a_{23} = \begin{bmatrix} a_{11} & a_{12} \\ a_{31} & a_{32} \end{bmatrix}$

Minor of $a_{31} = \begin{vmatrix} a_{12} & a_{13} \\ a_{22} & a_{23} \end{vmatrix}$   Minor of $a_{32} = \begin{vmatrix} a_{11} & a_{13} \\ a_{21} & a_{23} \end{vmatrix}$   Minor of $a_{33} = \begin{bmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{bmatrix}$

A cofactor (cij) is a minor with a prescribed sign.  Cofactor of an element is obtained by multiplying the minor of the element with where i is the number of row and j is the number of column.

That is $|cij| = (-1)^{i+j} Mij$

A cofactor matrix is a matrix in which every element is replaced with its cofactor cij.

**Example 1:**   A = $\begin{bmatrix} 7 & 12 \\ 4 & 3 \end{bmatrix}$ Matrix of cofactors C$_{ij}$ = $\begin{bmatrix} 3 & -4 \\ -12 & 7 \end{bmatrix}$

**Example 2:**   B = $\begin{bmatrix} -2 & 5 \\ 13 & 6 \end{bmatrix}$ Matrix of cofactors C$_{ij}$ = $\begin{bmatrix} 6 & 13 \\ -5 & -2 \end{bmatrix}$

**Example 3:**   C $\begin{bmatrix} 2 & 3 & 1 \\ 4 & 1 & 2 \\ 5 & 3 & 4 \end{bmatrix}$ = Matrix of Cofactors C$_{ij}$ = $\begin{bmatrix} -2 & -6 & 7 \\ -9 & 3 & 9 \\ 5 & 0 & -10 \end{bmatrix}$

**Example 4:**   D = $\begin{bmatrix} 6 & 2 & 7 \\ 5 & 4 & 9 \\ 3 & 3 & 1 \end{bmatrix}$ Matrix of Cofactors C$_{ij}$ = $\begin{bmatrix} -23 & 22 & 3 \\ 19 & -15 & -12 \\ -10 & -19 & 14 \end{bmatrix}$

## ADJOINT MATRIX

An adjoint matrix is transpose of a cofactor matrix that is adjoint of a given square matrix is the transpose of the matrix formed by cofactors of the elements of a given square matrix taken in order.

**Example 1:** $\qquad A = \begin{bmatrix} 13 & 17 \\ 19 & 15 \end{bmatrix}$

Matrix of Cofactors $C_{ij} = \begin{bmatrix} 15 & -19 \\ -17 & 13 \end{bmatrix}$

Adjoint of A = $\left[ c_{ij} \right]^t = \begin{bmatrix} 15 & -17 \\ -19 & 13 \end{bmatrix}$

**Example 2:** $\qquad A = \begin{bmatrix} 6 & 7 \\ 12 & 9 \end{bmatrix} \quad C_{ij} = \begin{bmatrix} 9 & -12 \\ -7 & 6 \end{bmatrix}$

Adj A = $\left[ c_{ij} \right]^t = \begin{bmatrix} 9 & -7 \\ -12 & 6 \end{bmatrix}$

**Example 3:**

$B = \begin{bmatrix} 0 & 1 & 2 \\ 1 & 2 & 3 \\ 3 & 1 & 1 \end{bmatrix} \qquad c_{ij} = \begin{bmatrix} -1 & 8 & -5 \\ 1 & -6 & 3 \\ -1 & 2 & -1 \end{bmatrix}$

Adj A = $\left[ c_{ij} \right]^t = \begin{bmatrix} -1 & 1 & -1 \\ 8 & -6 & 2 \\ -1 & 2 & -1 \end{bmatrix}$

**Example 4:**

$A = \begin{bmatrix} 13 & 2 & 8 \\ -9 & 6 & -4 \\ -3 & 2 & -1 \end{bmatrix} \quad C_{ij} = \begin{bmatrix} 2 & 3 & 0 \\ 14 & 11 & -20 \\ -40 & -20 & 60 \end{bmatrix}$

A $C_{ij}$ A = $C_{ij}{}^t = \begin{bmatrix} 2 & 14 & -40 \\ 3 & 11 & -20 \\ 0 & -20 & 60 \end{bmatrix}$

## **INVERSE  MATRIX**

For a square matrix A, if there exists a square matrix B such that AB=BA=1, then B is called the inverse matrix of A and is denoted as $A^{-1} = \dfrac{A d_j A}{IAI}$

**Example 1:** $\qquad A = \begin{bmatrix} 3 & 2 \\ 1 & 0 \end{bmatrix}$

$|A| = {}^-2$

$$\text{Adj A} = \begin{bmatrix} 0 & -2 \\ -1 & 3 \end{bmatrix}$$

$$A^{-1} = \frac{Adj A}{|A|} = \begin{bmatrix} \dfrac{0}{-2} & \dfrac{-2}{-2} \\ \dfrac{-1}{-2} & \dfrac{3}{-2} \end{bmatrix}$$

$$A^{-1} = \begin{bmatrix} 0 & 1 \\ \dfrac{1}{2} & -2 \end{bmatrix}$$

**Example 2:**  $A = \begin{bmatrix} 7 & 9 \\ 6 & 12 \end{bmatrix}$

$|A| = 30$

$$A^{-1} = \begin{bmatrix} \dfrac{Adj A}{|A|} \end{bmatrix} = \dfrac{\begin{matrix} 2 & -9 \\ -6 & 9 \end{matrix}}{30}$$

$$A^{-1} = \begin{bmatrix} \dfrac{2}{5} & \dfrac{-3}{10} \\ \dfrac{-1}{5} & \dfrac{7}{30} \end{bmatrix}$$

**Example 3:**  $A = \begin{bmatrix} 1 & 2 & 3 \\ 5 & 7 & 4 \\ 2 & 1 & 3 \end{bmatrix}$

$|A| = {}^-24$

$$A_{dj} A = \begin{bmatrix} 17 & -3 & -13 \\ -7 & -3 & 11 \\ -9 & 3 & -3 \end{bmatrix}$$

$$A^{-1} = \frac{Adj}{|A|} = \begin{bmatrix} \dfrac{17}{-24} & \dfrac{-3}{-24} & \dfrac{-13}{-24} \\ \dfrac{-7}{-24} & \dfrac{-3}{-24} & \dfrac{11}{-24} \\ \dfrac{-9}{-24} & \dfrac{3}{-24} & \dfrac{-3}{-24} \end{bmatrix}$$

**Example 4:**  $A = \begin{bmatrix} 4 & 2 & 5 \\ 3 & 1 & 8 \\ 9 & 6 & 7 \end{bmatrix}$

$|A| = {}^-17$

$$AdjA = \begin{bmatrix} -41 & 16 & 11 \\ 51 & -17 & -17 \\ 9 & -6 & -2 \end{bmatrix}$$

$$A^{-1} = \frac{AdjA}{|A|} = \begin{bmatrix} \frac{41}{17} & \frac{-16}{17} & \frac{-11}{17} \\ -3 & 1 & 1 \\ \frac{-9}{17} & \frac{6}{17} & \frac{2}{17} \end{bmatrix}$$

## CRAMERS RULE FOR MATRIX SOLUTIONS

Cramers rule provides a simplified method of solving a system of linear equations through the use of determinants. Cramer's rule states

$$x_i = \frac{|Ai|}{|A|}$$

where $x_i$ is the unknown variable, $|A|$ is the determinant of the coefficient matrix and $|Ai|$ is the determinant of special matrix formed from the original coefficient matrix by replacing the column of coefficient of $x_i$ with the column vector of constants

**Example 1 :**          To solve          $6x_1 + 5x_2 = 49$

$$3x_1 + 4x_2 = 32$$

In matrix from A×=B

$$\begin{bmatrix} 6 & 5 \\ 3 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 49 \\ 32 \end{bmatrix}$$

$$|A| = 9$$

To solve for $x_i$ , replace the first column of A, that is coefficient of $x_i$ with vector of constants B, forming a new matrix $A_1$ ,

$$A_1 = \begin{bmatrix} 49 & 5 \\ 32 & 4 \end{bmatrix}$$

$$|A1| = 36$$

$$x_1 = \frac{|A_1|}{|A|} = \frac{36}{9} = 4$$

Similarly, to solve for $x_2$, replace the second column of A, that is coefficient of $x_2$, with vector of constants B, forming a new matrix $A_2$,

$$A_2 = \begin{bmatrix} 6 & 49 \\ 3 & 32 \end{bmatrix}$$

$$|A_2| = 45$$

$$x_2 = \frac{|A_2|}{|A|} = \frac{45}{9} = 5$$

$\therefore$ the solution is $x_1 = 4$ and $x_2 = 5$

**Example 2:**

$$2x_1 + 6x_2 = 22$$

$$-x_1 + 5x_2 = 22$$

$$Ax = B$$

$$\begin{bmatrix} 2 & 6 \\ -1 & 5 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} = \begin{bmatrix} 22 \\ 53 \end{bmatrix}$$

$$|A| = 16, \qquad A_1 = \begin{bmatrix} 22 & 6 \\ 53 & 5 \end{bmatrix}$$

$$|A2| = {}^-208$$

$$x_1 = \frac{|A_1|}{|A|} = \frac{-208}{16} = {}^-13$$

$$A_2 = \begin{bmatrix} 2 & 22 \\ -1 & 53 \end{bmatrix}$$

$$|A_2| = 128$$

$$x_2 = \frac{|A_2|}{|A|} = \frac{128}{16} = 8$$

$\therefore$ the solution is $x_1 = -13$ and $x_2 = 8$

**Example 3:**

$$7x_1 - x_2 - x_3 = 0$$

$$10x_1 - 2x_2 + x_3 = 8$$

$$6x_1 + 3x_2 - 2x_3 = 7$$

$$AX = B$$

$$\begin{bmatrix} 7 & -1 & -1 \\ 10 & -2 & 1 \\ 6 & 3 & -2 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 8 \\ 7 \end{bmatrix}$$

$|A| = -61$

$$|A_1| = \begin{bmatrix} 0 & -1 & -1 \\ 8 & -2 & 1 \\ 7 & 3 & -2 \end{bmatrix}$$

$|A_1| = -61$

$$x_1 = \frac{|A_1|}{|A|} = \frac{-61}{-61} = 1$$

$$|A_2| = \begin{bmatrix} 7 & 0 & -1 \\ 10 & 8 & 1 \\ 6 & 7 & -2 \end{bmatrix}$$

$|A_2| = -183$

$$x_2 = \frac{|A_2|}{|A|} = \frac{-183}{-61} = 3$$

$$A_3 = \begin{bmatrix} 7 & -1 & 0 \\ 10 & -2 & 8 \\ 6 & 3 & 7 \end{bmatrix}$$

$|A_3| = -244$

$$x_3 = \frac{|A_3|}{|A|} = \frac{-244}{-61} = 4$$

$\therefore$ the solution is $x_1 = 1$

$x_2 = 3$

$x_3 = 4$

**Example 4:**    $5x_1 - 2x_2 + 3x_3 = 16$

$2x_1 + 3x_2 - 5x_3 = 2$

$4x_1 - 5x_2 + 6x_3 = 7$

$$AX = B \; \Box \; A = \begin{bmatrix} 5 & -2 & 3 \\ 2 & 3 & -5 \\ 4 & -5 & 6 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix} = \begin{bmatrix} 16 \\ 2 \\ 17 \end{bmatrix}$$

$|A| = -37$

$$A_1 = \begin{bmatrix} 16 & -2 & 3 \\ 2 & 3 & -5 \\ 7 & -5 & 6 \end{bmatrix}$$

$$|A_1| = {}^-111 \qquad x_3 = \frac{|A_1|}{|A|} = \frac{-111}{-37} = 3$$

$$A_2 = \begin{bmatrix} 5 & 16 & 3 \\ 2 & 2 & -5 \\ 4 & 7 & 6 \end{bmatrix}$$

$$|A_2| = {}^-259 \qquad x_2 = \frac{|A_2|}{|A|} = \frac{-259}{-37} = 7$$

$$A_3 = \begin{bmatrix} 5 & -2 & 16 \\ 2 & 3 & 2 \\ 4 & -5 & 7 \end{bmatrix}$$

$$|A_3| = {}^-185 \qquad x_3 = \frac{|A_3|}{|A|} = \frac{-185}{-37} = 5$$

$\therefore$ solution is $x_1 = 3$, $x_2 = 7$ and $x_3 = 5$

# MODULE III
# FUNCTIONS AND GRAPHS
## Part A
## FUNCTIONS

Suppose you worked in a shop part time in the evening. You are paid on an hourly basis and you earn Rs. 100 for an hour. The more you work the more you are paid. That means if you work for one hour you get Rs. 100, if you work for two hours, you get Rs. 200 and so on. This implies that the amount of money you earn depends on the time you work. If this sentence is written in mathematical format, we can write the amount of money you earn is a function of the time you work. Here the amount of money you earn is a dependent on the time you work. But the time you work is independent. This is the crux of a functional relation. For example, if we represent the amount you earn by 'y' and 'the time you work by x' and write in mathematical form, it can be written as $y = f(x)$.

Before seeing the formal definition of a function, let us first see understand the concept of a variable better.

Variable: A variable is a value that may change within the scope of a given problem or set of operations. Thus a variable is a symbol for a number we don't know yet. It is usually a letter like x or y. We call these letters 'variables' because the numbers they represent can vary- that is, we can substitute one or more numbers for the letters in the expression.

In the above example $y = f(x)$, both y and x are variables. But there is one difference. The amount you earn (y) is depend on the time you work (x). But the length of time you work (x) is independent of the amount you earn(y). Here 'y' is called a dependent variable and 'x' is called an independent variable. Sometimes the independent variable is called the 'input' to the function and the dependent variable the 'output' of the function.

Thus dependent variable is a variable whose value depends on those of others; it represents a response, behaviour, or outcome that the researcher wishes to predict or explain. Independent variable is any variable whose value determines that of others. Sometimes we also see a extraneous variable, which is a factor that is not itself under study but affects the measurement of the study variables or the examination of their relationships.

Constant: In contrast to a variable, a constant is a value that remains unchanged, though often unknown or undetermined. In Algebra, a constant is a number on its own, or sometimes a letter such as a, b or c to stand for a fixed number.

Coefficients: Coefficients are the number part of the terms with variables. In $3x^2 + 2y + 7xy + 5$, the coefficient of the first term is 3. The coefficient of the second term is 2, and the coefficient of the third term is 7. Note that if a term consists of only variables, its coefficient is 1.

Expressions consisting of a real number or of a coefficient times one or more variables raised to the power of a positive integer are called monomials. Monomials can be added or subtracted to form polynomials.

Variable Expression:    A variable expression is a combination of numbers (or constants), operations, and variables. For example in a variable expressions 5a + 3b, 'a' and 'b' are variables, 5 and 3 are constants and + is an operator.

Function: A mathematical function relates one variable to another. There are lots of other definitions for a function and most of them involve the terms rule, relation, or correspondence. While these are more technically accurate than the definition that we are going to use is

A function is an equation for which any x that can be plugged into the equation will yield exactly one y out of the equation. Let us explore further.

When one value is completely determined by another, or several others, then it is called a function of the other value or values. In this case the value of the function is a dependent variable and the other values are independent variables. The notation f(x) is used for the value of the function f with x representing the independent variable. Similarly, notation such as f(x, y, z) may be used when there are several independent variables.

Here we used 'f' to represent a function. We may also represent a function by using symbols like '$g$' or '$h$' or the Greek letter $\phi$ *phi*. Just as your name signifies all of the many things that make 'you,' a symbol like '$f$' serves as a shorthand for what may be a long or complicated rule expressing a particular relationship between variable quantities. Then a function $y = f(x)$ shows that $f$ is a rule which assigns a unique value of the variable quantity $y$ to values of the variable quantity $x$.

**Example :** (1) y = 3x – 2          (2) h = 5x + 4y

In other words, a function can be defined as a set of mathematical operations performed on one or more inputs (variables) that results in an output. Consider a simple function y = x + 1. In this case, x in the input value (independent variable) and y is the output (dependent variable). By putting any number in for x, we calculate a corresponding output y by simply adding one. The set of possible input values is known as the domain, while the set of possible outputs is known as the range.

| x | 1 | 2 | 3 | 4 | 5 | 6 | Domain |
|---|---|---|---|---|---|---|--------|
| y | 2 | 3 | 4 | 5 | 6 | 7 | Range |

Note that the above table shows that as we give different values to the independent variable x, the function (y) assumes different values.

Now consider the function y = 3x – 2. Here the variable y represents the function of whatever inputs appear on the other side of the equation. In other words, y is a function of the variable x in y = 3x – 2.  Because of that, we sometimes see the function written in this form f(x) = 3x – 2. Here f(x) means just the same as "y =" in front of an equation. Since there is really no significance to y, and it is just an arbitrary letter that represents the output of the function, sometimes it will be written as f(x) to indicate that the expression is a function of x. As we have discussed above, a function may also be written as g(x), h(x), and so forth, but f(x) is the most common because function starts with the letter f.

Thus the key point to remember is that all of the following are same, they are just different ways of expressing a function.

$$y = 3x - 2,\ y = f(3x - 2),\ f(x) = 3x - 2,\ h(x) = 3x - 2,\ g(x) = 3x - 2$$

Evaluate a function: To evaluate a function means to pick different values for the independent variable x (often named the input) in order to find the dependent variable y (often named output). In terms of evaluation, for every choice of x that we pick, only one corresponding value of y will be the end result. For example if you are asked to evaluate the function y = 3x – 2 at x = 5, substitute the value at the place of x. Then you get y = 3(5) – 2 = 13. Note that since a function is a unique mapping from the domain (the inputs) to the range (the outputs), there can only be one output for any input. However, there can, be many inputs which give the same output.

Thus, strictly speaking, a function is a rule that produces one and only one value of y for any given value of x. Some equations, such as y = $\sqrt{x}$ , are not functions because they produce more than one value of y for a given value of x.

To plot a graph of a function, as a matter of convention, we denote the independent variable as x and plot it on the horizontal axis of a graph, and we denote the dependent variable as y and plot it on the vertical axis.

**Types of functions:** (Classification of Functions)

Functions take a variety of forms, but to begin with, we will concern ourselves with the three broad categories of functions mentioned earlier: linear, exponential, and power.

Linear functions take the form y = a + bx where a is called the y-intercept and b is called the slope. The reasons for these terms will become clear in the course of this exercise.

Exponential functions take the form $y = a + q^x$

Power functions take the form $y = a + kx^p$

Note the difference between exponential functions and power functions. Exponential functions have a constant base (q) raised to a variable power (x); power functions have a variable base (x) raised to a constant power (p). The base is multiplied by a constant (k) after raising it to the power (p).

Here are some more functions with examples.

Linear function: Linear function is a first-degree polynomial function of one variable. These functions are known as 'linear' because they are precisely the functions whose graph is a straight line. Example: When drawn on a common (x, y) graph it is usually expressed as $f(x) = mx + b$. This is a very common way to express a linear function is named the 'slope-intercept form' of a linear function.

Equations whose graphs are not straight lines are called nonlinear functions. Some nonlinear functions have specific names. A quadratic function is nonlinear and has an equation in the form of $y = ax^2 + bx + c$ where $a \neq 0$. Another nonlinear function is a cubic function. A cubic function has an equation in the form of $y = ax^3 + bx^2 + cx + d$, where $a \neq 0$.

Or, in a formal function definition it can be written as $f(x) = mx + b$. Basically, this function describes a set, or locus, of (x, y) points, and these points all lie along a straight line. The variable m holds the slope of this line. The variable b holds the y-coordinate for the spot where the line crosses the y-axis. This point is called the 'y-intercept'.

Quadratic function: A polynomial of the second degree, represented as $f(x) = ax^2 + bx + c$, $a \neq 0$. Where a, b, c are constant and x is a variable. Example, value of $f(x) = 2x^2 + x - 1$ at $x = 2$. Put $x = 2$, $f(2) = 2.2^2 + 2 - 1 = 9$

Single valued function: If for every value of x, there correspond only one value of y, then y is called a single valued function. Example: $f(x) = 2x + 3$, when $x=3$, $f(3) = 2(3)+3 = 6+3 = 9$

Many valued function: If for each value of x, there corresponds more than one value of y, then y is called many or multi valued function.

Explicit function: If the functional relation between the two variables x and y is expressed in the form $y = f(x)$, y is called an explicit function of x. Example: $y = 4x - 5$

Implicit function: If the relation between two variables x and y is expressed in the form $f(x,y)=0$, where x cannot be expressed as a function of y, or y cannot be expressed as a function of x, is called an implicit function. Example: $y - 4x - 5 = 0$

Odd and Even Functions: When there is no change in the sign of f(x) when x is changed to –x, then that function is called an even function. (i.e) f(-x) = f(x). The graph of an even function is such that the two ends of the graph will be directed towards the same side.(Figure 1)



Figure 1

When the sign of f(x) is changed when x is changed to –x, then it is called an odd function.

(i.e) f(-x) = - f(x). The following graph (Figure 2) shows the odd function, f(x)=x3 , and its reflection about the y-axis, which is f(-x) = −x3.



Figure 2

Inverse functions: From every function y = f(x), we may be able to deduce a function x=g(y). This means that the composition of the function and its inverse is an identity function. In an inverse function we may be able to express the independent variable in terms of the dependent variable. Function g(x) is inverse function of f(x) if, for all x, g(f(x)) = f(g(x)) = x.  A function f(x) has an inverse function if and only if f(x) is one-to-one.  For example, the inverse of f(x) = x + 1 is g(x) = x - 1

Polynomial functions of degree n: A polynomial is an expression of finite length constructed from variables (also called indeterminate) and constants, using only the operations of addition, subtraction, multiplication, and non-negative integer exponents. For example, $x2 − x/4$

+ 7 is a polynomial, but $x2 - 4/x + 7x3/2$ is not, because its second term involves division by the variable x (4/x), and also because its third term contains an exponent that is not an integer (3/2).

Example $f(x) = a_n x^n + a_{n-1} x^{n-1} + \ldots + a_0$ (n = nonnegative integer, $a_n \neq 0$)

Rational function: The term rational comes from 'ratio'

$$f(x) = \frac{g(x)}{h(x)}$$

Where g(x) and h(x) are both polynomials and h(x) $\neq$ 0

Algebraic Functions: A function which consists of finite number of terms involving powers and roots of independent variable x and the four fundamental operations of addition, subtraction, multiplication and division is called an algebraic function. Polynomials, rational functions and irrational functions are all the examples of algebraic functions.

Transcendental functions: Functions which are not algebraic are called transcendental functions. Trigonometric functions, Inverse trigonometric functions, exponential functions, logarithmic functions are all transcendental functions.



Figure 3

Example, f(x) = sinx, g(x) = log(x), h(x)= ex , k(x) = tan−1 (x)

 Modulus functions: Modulus functions are defined as follows.

$y = |x| = \{$ x, if x $\geq$ 0

- x , if x < 0.

For example, | 3 | = 3, and | -4 |= -(-4) = 4 , since -4 < 0.

The graph of the modulus function y = |x| is shown below. (Figure 4)

(Figure 4)

Onto functions: A function f(x) is one-to-one (or injective),if f be a function with domain D and range R. A function g with domain R and range D is an inverse function for f if, for all x in D, y = f(x) if and only if x = g(y).

One-to-one functions: A function f(x) is one-to-one (or injective) if, for every value of f, there is only one value of x that corresponds to that value of 'f'. Eg. f(x) = x + 3 is is one-to-one, because, for every possible value of f(x), there is exactly one corresponding value of x.

Identity functions: A polynomial of the first degree, represented as f(x) = x, example, values of f(x) = x, at x = 1,2 are f(1) = 1 and f(2) = 2

Constant functions: It is a polynomial of the 'zeroth' degree where f(x) = cx0 = c(1) = c. It disregards the input and the result is always c. Its graph is a horizontal line. For example f(x) = 2, whatever the value of x result is always 2.

Linear functions: It is a polynomial of the first degree, the input should be multiplied by m and it adds to c. It is represented as f(x) = m x + c such as f(x) = 2x + 1 at x = 1.f(1) = 2 . 1 + 1 = 3 that is f(1) = 3

Trigonometrical functions: Trigonometric functions are often useful in modeling cyclical trends such as the seasonal variation of demand for certain items, or the cyclical nature of recession and prosperity. There are six trigonometric functions: sin(x), cos(x), tan(x), csc(x), sec(x) and cot(x)

Analytic functions: All polynomials and all power series in the interior of their circle of convergence are analytic functions. Arithmetic operations of analytic function are differentiated according to the elementary rules of the calculation, and, hence analytic function.

Differentiable functions: A differentiable function is a function whose derivative exists at each point in its domain. If x0 is a point in the domain of a function f, then f is said to be differentiable at x0 if the derivative f '(x0) is defined. The graph of differential functions is always smooth.

Smooth functions: A smooth function is a function that has continuous derivatives over some domain or we can say that, it is a function on a Cartesian space Rn with values in the real line R if its derivatives exist at all points

Even and odd functions: A function for every x in the domain of f is an even function if f(-x) = f(x) and odd function if f(-x) = - f(x) for example, f(x) = x2 is an even function because f(-x) = (-x)2 = x2 = f(x) and f(x) = x is an even function because f(x) = (-x) = -x = - f(x)

Curves functions: A space curve C is the set of all ordered triples (f(t),g(t),h(t)) together with their definition parametric equations x= f(t) y = g(t) and z = h(t) where f, g, h are continuous functions of t on an interval I.

Composite functions: There is one particular way to combine functions. The value of a function f depends upon the value of another variable x and that variable could be equal to another function g, so its value depends on the value of a third variable. If this is the case, then the first variable which is a function h, is called as the composition of two functions ( f and g). It denoted as f o g = (f o g) x = f(g(x)). For example, let f(x) = x+1 and g(x) = 2x then h(x) = f(g(x)) = f(2x) = 2x + 1.

Monotonic functions: A monotonic function is a function that preserves the given order. These are the functions that tend to move in only one direction as x increases. A monotonic increasing function always increases as x increases, that is f(a) > f(b) for all a>b. A monotonic decreasing function always decreases as x increases, that is . f(a) < f(b) for all a>b

Periodic functions: Functions which repeat after a same period are called as a periodic function, such as trigonometric functions like sine and cosine are periodic functions with the period $2\prod$.

**Evaluation of Functions**

To evaluate a function, substitute the value of the variable. If a is any value of x, the value of the function f(x) for x = a is denoted by f(a).

Examples

**1.** Given $y = f(x) = x^2 + 4x - 5$ find $f(2)$ and $f(-4)$

To find $f(2)$, substitute x = 2

$f(2) = 2^2 + 4(2) - 5 = 7$

To find $f(-4)$, substitute x = -4

$f(-4) = -4^2 + 4(-4) - 5 = 16 - 16 - 5 = -5$

**2.** Given $y = f(x) = 3x^3 - 4x^2 + 4x - 10$, find value of the function at $f(2)$ and $f(-3)$

$f(2) = 3(2)^3 - 4(2)^2 + 4(2) - 10 = 6$

$f(-3) = 3(-3)^3 - 4(-3)^2 + 4(-3) - 10 = -139$

**3.** Given $y = f(x) = \sqrt{x+2}$ find value of the function when x = 0, x = 7, x= -2

$$f(0) = \sqrt{0+2} = \sqrt{2}$$

$$f(17) = \sqrt{7+2} = \sqrt{9} = \pm 3$$

$$f(-2) = \sqrt{-2+2} = 0$$

**4.** Given $f(x) = x^2 + 4x - 11$ find f(x+1)

This means we have to evaluate the function when x = x + 1

$$f(x + 1) = (x + 1)^2 + 4(x + 1) - 11$$

Simplifying using $(a+b)^2$

$$f(x + 1) = x^2 + 2x + 1 + 4x + 4 - 11 = x^2 + 6x - 6$$

**5.** $f(x) = \dfrac{x^3-5}{x+3}$     find $f(2)$

$$f(x) = \dfrac{2^3-3}{2+3} = f(x) = \dfrac{8-3}{5} = 1$$

**6.** Given $y = x^2 + 1$, find f(0), f(–1) and f(3). Also find $\dfrac{f(0)+f(-1)}{f(3)}$

f(0) = 1, f(–1) = 2, f(3) = 10

$$\dfrac{f(0)+f(-1)}{f(3)} = \dfrac{1+2}{10} = \dfrac{3}{10}$$

**7.** Given a function f(x) = 2x + 7, find f(4), f(1) and f(–2)

f(4) = 2(4) + 7, so f(4) = 8 + 7 = 15

f(1) = 2(1) + 7 = 10

f(–2) = 2(–2) + 7 = –4 + 7 = 3

**RECTANGULAR CO-ORDINATE SYSTEM AND GRAPHS OF FUNCTIONS**

Equations can be graphed on a set of coordinate axes. The location of every point on a graph can be determined by two coordinates, written as an ordered pair, (x,y). These are also known as Cartesian coordinates, after the French mathematician Rene Descartes, who is credited with their invention.

The rectangular coordinate system, also called the Cartesian coordinate system or the x-y coordinate system consists of 4 quadrants, a horizontal axis, a vertical axis, and the origin. The horizontal axis is usually called the x-axis, and the vertical axis is usually called the y-axis. The origin is the point where the two axes cross. The coordinates of the origin are (0,0). This notation is called an ordered pair. The first coordinate (or abscissa) is known as the x-coordinate, while the second coordinate (or ordinate) is the y-coordinate. These tell how far and in what direction we move from the origin.

The Rectangular Coordinate System



In the above figure please note the following **features of the Quadrants**

The x-axis and y-axis separate the coordinate plane into four regions called **quadrants.** The upper right quadrant is quadrant I, the upper left quadrant is quadrant II, the lower left quadrant is quadrant III, and the lower right quadrant is quadrant IV. Notice that, as shown in Figure 1,

- in quadrant I, x is always positive and y is always positive (+,+)
- in quadrant II, x is always negative and y is always positive (–,+)
- in quadrant III, x is always negative and y is always negative (–,–)
- in quadrant IV, x is always positive and y is always negative (+,–)

**Examples**

1. Plot the ordered pair (−3, 5) and determine the quadrant in which it lies.

The coordinates x=−3 and y=5 indicate a point 3 units to the left of and 5 units above the origin. The point is plotted in quadrant II (QII) because the x-coordinate is negative and the y-coordinate is positive.

**Example 2:** Plot ordered pairs: (4, 0), (−6, 0), (0, 3), (−2, 6), (−4, −6)

Ordered pairs with 0 as one of the coordinates do not lie in a quadrant; these points are on one axis or the other (or the point is the origin if both coordinates are 0). Also, the scale indicated on the *x*-axis may be different from the scale indicated on the *y*-axis. Choose a scale that is convenient for the given situation.



**Distance Formula**

Frequently you need to calculate the distance between two points in a plane. To do this, form a right triangle using the two points as vertices of the triangle and then apply the Pythagorean Theorem. The Distance Formula is a variant of the Pythagorean Theorem.

The distance formula can be obtained by creating a triangle and using the Pythagorean Theorem to find the length of the hypotenuse. The hypotenuse of the0 triangle will be the distance between the two points.

$$d=\sqrt{\left(x_2-x_1\right)^2+\left(y_2-y_1\right)^2}$$

x2 and y2 are the x,y coordinates for one point.
x1and y1 are the x,y coordinates for the second point.

d is the distance between the two points.

For example, consider the two points A (1,4) and B (4,0). Find the distance between them.
so: $x_1 = 1$, $y_1 = 4$, $x_2 = 4$, and $y_2 = 0$.
Substituting into the distance formula we have:

$$d=\sqrt{(4-1)^2+(0-4)^2}$$

$$d=\sqrt{(3)^2+(-4)^2}$$

$$d=\sqrt{9+16}$$

$$d=\sqrt{25}$$

$$d=5$$

Example 1:  Find the distance between (−1, 2) and (3, 5).



so: $x_1 = -1$, $y_1 = 2$, $x_2 = 3$, and $y_2 = 5$.
Substituting into the distance formula we have:
$$d=\sqrt{\left(x_2-x_1\right)^2+\left(y_2-y_1\right)^2}$$

$$d=\sqrt{(3--1)^2+(5-2)^2}$$

$$d=\sqrt{(4)^2+(3)^2}$$

$$d=\sqrt{16+9}$$

$$d=\sqrt{25}$$

$$d=5$$

**Midpoint Formula**

The midpoint is an ordered pair formed by finding the average of the *x*-values and the average of the *y*-values of the given points. The point that bisects the line segment formed by two points, (x1, y1) and (x2, y2), is called the midpoint and is given by the following formula:

$$\left(\frac{x_1+x_2}{2}, \frac{y_1+y_2}{2}\right)$$

Example1: Find the midpoint of the line segment joining P(–3, 8)and Q(4, –4)

Use the midpoint formula:

$$\left(\frac{-3+4}{2}, \frac{8\pm4}{2}\right)$$

$$\left(\frac{1}{2}, \frac{4}{2}\right)$$

$$\left(\frac{1}{2}, 2\right)$$

Example 2: Find the midpoint of the line segment joining P(-5, -4)and Q(3, 7)

Use the midpoint formula:

$$\left(\frac{-5+3}{2}, \frac{-4+7}{2}\right)$$

$$\left(\frac{-2}{2}, \frac{3}{2}\right)$$

$$\left(-1, \frac{3}{2}\right)$$

**Graphs**

A graph is a picture of one or more functions on some part of a coordinate plane. Our paper usually is of limited size. The coordinate plane is infinite. Therefore, we need to pick an area of the coordinate plane that would fit on the piece of paper. That means specifying minimum and maximum values of x and y coordinates.

To draw the graph of a function y = $f$(x), the values of the independent variable, x, is marked on the horizontal axis. The values of the dependent variable, y, is placed on the vertical axis. To graph a function we should know the slope of the function.

Slope: Slope is used to find how steep a particular line is, or it can be used to show how much something has changed over time. The slope of a line measures the change in y $(\Delta y)$ divided by change in x $(\Delta x)$. We calculate slope by using the following definition. In Algebra, slope is defined as the rise over the run. This is written as a fraction like this:
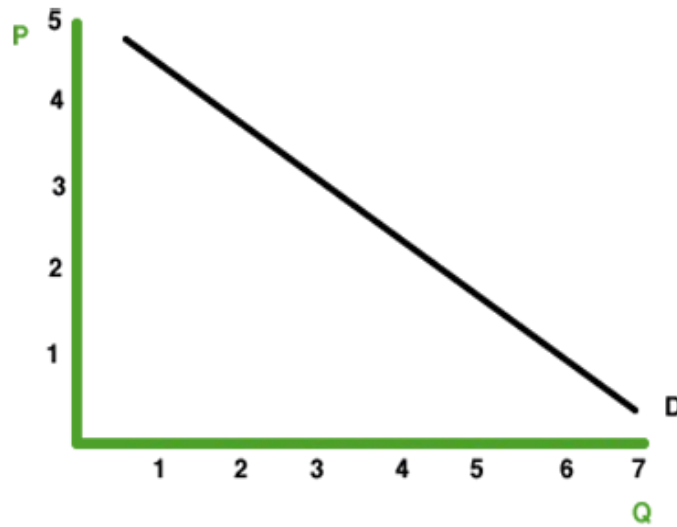
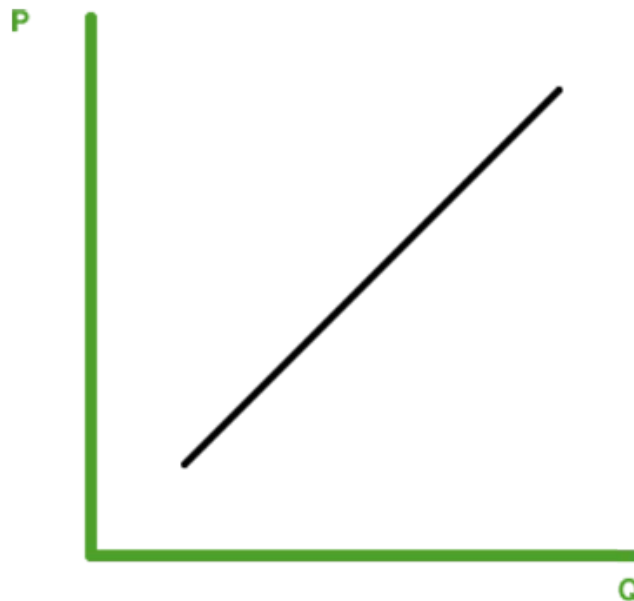$$Slope = \frac{rise}{run} = \frac{\Delta y}{\Delta x}$$





$$Slope = Slope = \frac{difference\ of\ y\ coordinates}{difference\ of\ x\ coordinates} = \frac{1-2}{3-1} = \frac{-1}{2}$$

The value of slope indicates the steepness and direction of a line. The greater the absolute value of the slope, the steeper the line. A positively slopped line moves up from left to right (for example a supply curve). A negatively slopped line moves down (for example a demand curve). The slope of a horizontal line (for example a perfectly elastic demand curve), for which

$\Delta y = 0,$ is zero. The slope of a vertical line (for example a perfectly inelastic demand curve),

for which $\Delta x = 0$ , is undefined. That is here slope does not exist since dividing by zero is not possible.

Negative Slope

Positive slope

Rules for Calculating the Slope of a Line

     1. Find two points on the line. Every straight line has a consistent slope. In other words, the slope of a line never changes. This fundamental idea means that you can choose ANY two points on a line to find the slope. This should intuitively make sense with your own understanding of a straight line. After all, if the slope of a line could change, then it would be a zigzag line and not a straight line.

     2. Count the rise (How many units do you count up or down to get from one point to the next?) Record this number as your numerator.

     3. Count the run (How many units do you count left or right to get to the point?) Record this number as your denominator.

    4. Simplify your fraction if possible.

    Important note: If you count up or right your number is positive. If you count down or left your number is negative.

Another approach to finding slope

Find the difference of both the y and x coordinates and place them in a ratio. Pick two points on the line. Because you are using two sets of ordered pairs both having x and y values, a subscript must be used to distinguish between the two values. Consider the following graph.
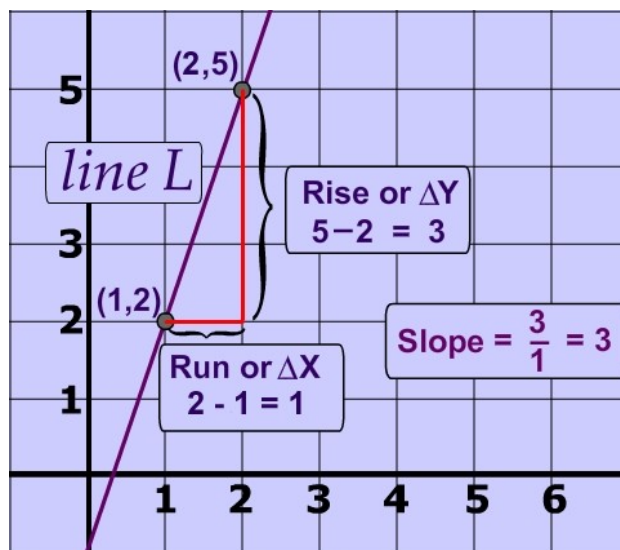


    Choose any two points. The two points we choose are $(x_1, y_1) = (3,2)$. and $(x_2, y_2) = (-1,-1)$ This is a simpler formula for finding the slope of a line. $m = \dfrac{y_2 - y_1}{x_2 - x_1}$ where m is the variable used for the slope.

Substituting

$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{-1-2}{-1-3} = \frac{-3}{-4} = \frac{3}{4}$$

Example 1: Find slope



    The slope of a line through the points (1, 2) and (2, 5) is 3 because every time that the line moves up three (the change in y or the rise) the line moves to the right (the run) by 1.

Example 2 : Find slope



rise to its run. The rise is the vertical change between two points on a line. The difference between the y-coordinates creates the rise.
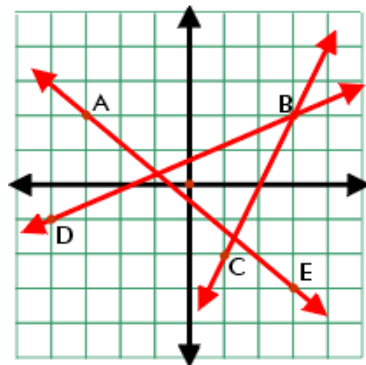
In this example, the difference between the y-coordinates of points A and B is $3 - (-3) = 6$.

The run is the horizontal difference between two points or the difference between the x-coordinates. In this example the difference is $-3 - 1 = -4$.

So here slope of the line = $Slope = \dfrac{\Delta y}{\Delta x} = \dfrac{6}{-4} = \dfrac{3}{-2}$

Example 3
Find Slope of the lines in the graph based on the information given



Line AE, coordinates $(-3,2) = (x1, y1)$ $(3,-3) = (x2, y2)$ , slope =
$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{-3 - 2}{3 - -3} = \frac{-5}{6}$$

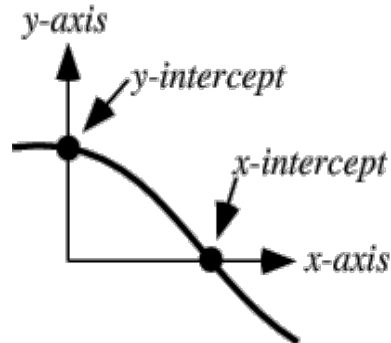Line DB, coordinates $(-4,-1) = (x1, y1)$, $(3,2) = (x2, y2)$ , slope =
$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{2 - -1}{3 - -4} = \frac{3}{7}$$

Line BC, coordinates (3,2) = (x1, y1), (1,-2) = (x2, y2) , slope =

$$m = \frac{y_2 - y_1}{x_2 - x_1} = \frac{-2 - 2}{1 - 3} = \frac{-4}{-2} = 2$$
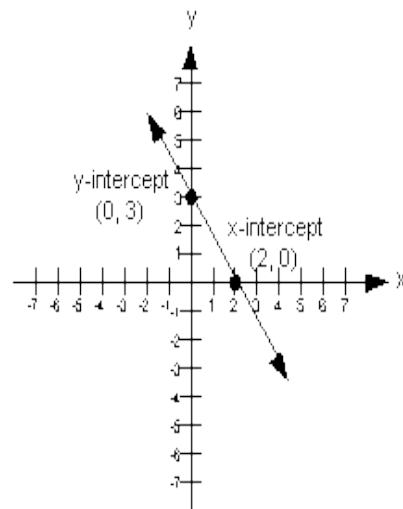
**Intercept:**

We are going to talk about x and y intercepts. An x intercept is the point where your line crosses the x-axis. The y intercept is the point where your line crosses the y-axis.



Algebraically,

An 'x-intercept' is a point on the graph where y is zero, that is, an x-intercept is a point in the equation where the y-value is zero.

A 'y-intercept' is a point on the graph where x is zero, that is. a 'y-intercept' is a point in the equation where the x-value is zero.



In the above illustration, the x-intercept is the point (2, 0) and the y-intercept is the point (0, 3).

Example 1: Find the x and y intercepts of the equation 3x + 4y = 12.

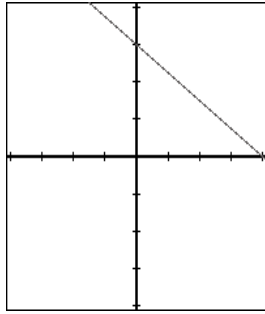To find the *x*-intercept, set *y* = 0 and solve for *x*.

$$3x + 4(\,0\,) = 12$$
$$3x + 0 = 12$$
$$3x = 12$$
$$x = 12/3$$
$$x = 4$$

To find the *y*-intercept, set *x* = 0 and solve for *y*.

$$3( 0 ) + 4y = 12$$
$$0 + 4y = 12$$
$$4y = 12$$
$$y = 12/4$$
$$y = 3$$

Therefore, the *x*-intercept is ( 4, 0 ) and the *y*-intercept is ( 0, 3 ).

The graph of the line looks like this:



Example 2:  Find the x and y intercepts of the equation 2x - 3y = -6.
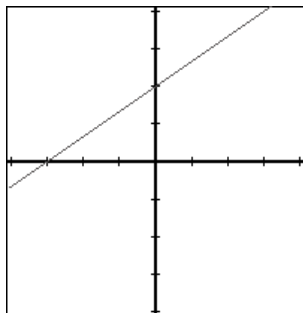
To find the *x*-intercept, set $y = 0$ and solve for *x*.

$$2x - 3( 0 ) = -6$$
$$2x - 0 = -6$$
$$2x = -6$$
$$x = -6/2$$
$$x = -3$$

To find the *y*-intercept, set $x = 0$ and solve for *y*.

$$2( 0 ) - 3y = -6$$
$$0 - 3y = -6$$
$$-3y = -6$$
$$y = -6/(-3)$$
$$y = 2$$

Therefore, the *x*-intercept is ( –3, 0 ) and the *y*-intercept  is ( 0, 2 ).

The graph of the line looks like this:

Example 3:  Find the x and y intercepts of the equation ( -3/4 )x + 12y = 9.
To find the *x*-intercept, set *y* = 0 and solve for *x*.

$$( -3/4 )x + 12( 0 ) = 9$$
$$( -3/4 )x + 0 = 9$$
$$( -3/4 )x = 9$$
$$x = 9/( -3/4 )$$
$$x = -12$$

To find the *y*-intercept, set *x* = 0 and solve for *y*.
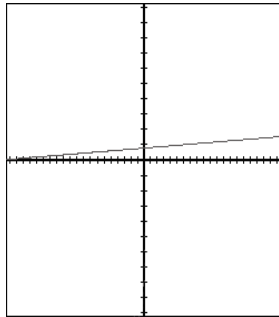
$$( -3/4 )( 0 ) + 12y = 9$$
$$0 + 12y = 9$$
$$12y = 9$$
$$y = 9/12$$
$$y = 3/4$$

Therefore, the *x*-intercept is ( −12, 0 ) and the *y*-intercept  is ( 0, 3/4 ).
The graph of the line looks like this:



Example 4:  Find the x and y intercepts of the equation -2x + ( 1/2 )y = -3.
To find the *x*-intercept, set *y* = 0 and solve for *x*.

$$-2x + ( 1/2 )( 0 ) = -3$$
$$-2x + 0 = -3$$
$$-2x = -3$$
$$x = -3/( -2 )$$
$$x = 3/2$$

To find the *y*-intercept, set *x* = 0 and solve for *y*.

$$-2( 0 ) + ( 1/2 )y = -3$$
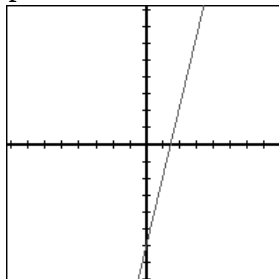$$0 + ( 1/2 )y = -3$$
$$( 1/2 )y = -3$$
$$y = -3/( 1/2 )$$
$$y = -6$$

Therefore, the *x*-intercept is ( 3/2, 0 ) and the *y*-intercept is ( 0, -6 ).
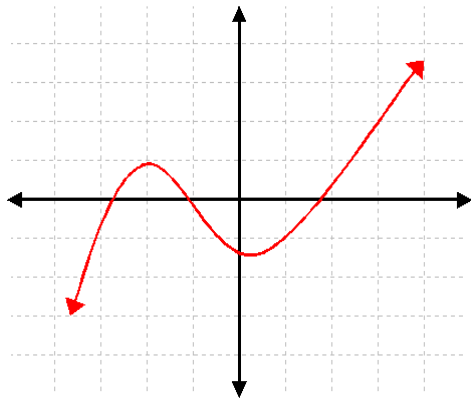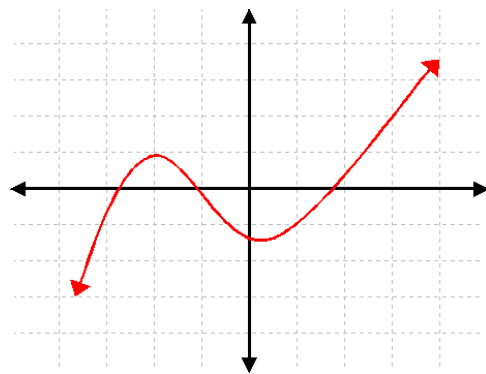The graph of the line looks like this:



**Graphing of Functions**

A function is a relation (usually an equation) in which no two ordered pairs have the same x-coordinate when graphed.

One way to tell if a graph is a function is the vertical line test, which says if it is possible for a vertical line to meet a graph more than once, the graph is not a function. The figure below is an example of a function.

Functions are usually denoted by letters such as f or g. If the first coordinate of an ordered pair is represented by x, the second coordinate (the y coordinate) can be represented by f(x). In the figure below, f(1) = -1 and f(3) = 2.

When a function is an equation, the domain is the set of numbers that are replacements for x that give a value for f(x) that is on the graph. Sometimes, certain replacements do not work, such as 0 in the function: f(x) = 4/x (because we cannot divide by 0).

**Straight-Line Equations: Slope Intercept form of a graph**

The slope intercept form of a line is

y = mx + b.

Where m is the slope and b is the y-intercept.

y = m x + b
↑ ↑
slope y-intercept


y = 2 x + 3
↑ ↑
slope y-intercept
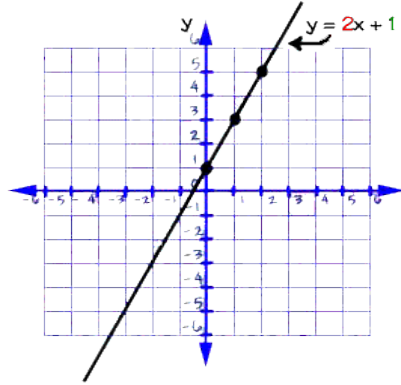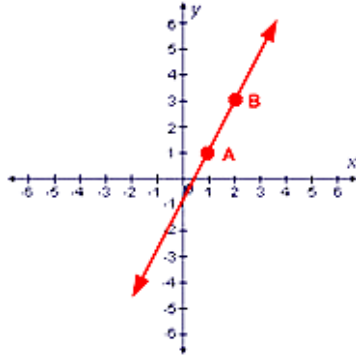
2/1 is the slope

(0,3) is the y intercept

Every straight line can be represented by an equation: y = mx + b. The coordinates of every point on the line will solve the equation if you substitute them in the equation for x and y.

The equation of any straight line, called a linear equation, can be written as: y = mx + b, where m is the slope of the line and b is the y-intercept.

The y-intercept of this line is the value of y at the point where the line crosses the y axis

Given the slope intercept form,  = mx + b, let's find the equation for this line. Pick any two points, in this diagram, A = (1, 1) and B = (2, 3).

We found that the slope m for this line is 2. By looking at the graph, we can see that it intersects the y-axis at the point (0, –1), so –1 is the value of b, the y-intercept. Substituting these values into the equation formula, we get:

**y = 2x –1**

The line shows the solution to the equation: that is, it shows all the values that satisfy the equation. If we substitute the x and y values of a point on the line into the equation, you will get a true statement.

To graph the equation of a line, we plot at least two points whose coordinates satisfy the equation, and then connect the points with a line. We call these equations "linear" because the graph of these equations is a straight line.

If the equation is given in the form y = mx + b, then the constant term, which is b, is the y intercept, and the coefficient of x, which is m, is the slope of the straight line.

The easiest way to graph such a line, is to plot the y-intercept first. Then, write the slope m in the form of a fraction, like rise over run, and from the y-intercept, count up (or down) for the rise, over (right or left) for the run, and put the next point. Then connect the two points and this is your line.

Example1: y = 3x + 2

Here Y-intercept = 2, slope = $\dfrac{3}{1} = 1$

Start by graphing the y-intercept by going up 2 units on the y-axis.

From this point go UP (rise) another 3 units, then 1 unit to the RIGHT (run), and put another point. This is the second point. Connect the points and it should look like this:

Example 2: y = –3x + 2

Y-intercept = 2, slope = $\dfrac{-3}{1}=-3$

Start by graphing the y-intercept by going up 2 units on the y-axis.

From this point go down (rise) 3 units, then 1 unit to the RIGHT (run), and put another point. This is the second point. Connect the points and it should look like this

Steps for Graphing a Line With a Given Slope

Plot a point on the y-axis. (In the next lesson, Graphing with Slope Intercept Form, you will learn the exact point that needs to be plotted first. For right now, we are only focusing on slope!)

Look at the numerator of the slope. Count the rise from the point that you plotted. If the slope is positive, count up and if the slope is negative, count down.

Look at the denominator of the slope. Count the run to the right.

Plot your point.

Repeat the above steps from your second point to plot a third point if you wish.

Draw a straight line through your points.

The trickiest part about graphing slope is knowing which way to rise and run if the slope is negative!

If the slope is negative, then only one - either the numerator or denominator is negative (NOT Both!) Remember your rules for dividing integers? If the signs are different then the answer is negative!

If the slope is negative you can plot your next point by going down and right OR up and left.

If the slope is positive you can plot your next point by going up and right OR down and left.

Example 1

This example shows how to graph a line with a slope of 2/3.

Example 1:



Slope of 2/3        Slope $= \dfrac{2 = rise}{3 = run}$

1. Plot 1st point. (I plotted this point at (0,-2)

2. Count the rise. Since the rise is positive 2, I counted up 2.

3. Count the run. Since the run is positive 3, I counted to the right 3.

4. Plot your second point. This point is (3, 0).

5. Draw a straight line through your points.

Example 2

Graph the linear function f given by f (x) = 2x + 4

We need only two points to graph a linear function. These points may be chosen as the x and y intercepts of the graph for example.

Determine the x intercept, set f(x) = 0 and solve for x.

2x + 4 = 0

x = -2

Determine the y intercept, set x = 0 to find f(0).

f(0) = 4

The graph of the above function is a line passing through the points (-2 , 0) and (0 , 4) as shown below.

Example 6

Graph the linear function f given by f (x) = -(1 / 3)x - 1 / 2

Determine the x intercept, set f(x) = 0 and solve for x.

-(1 / 3)x - 1 / 2 = 0

x = − 3 / 2 = − 1.5

Determine the y intercept, set x = 0 to find f(0).

f(0) = –1 / 2 = –0.5

The graph of the above function is a line passing through the points (-3 / 2 , 0) and (0 , -1 / 2) as shown below.

**Straight-Line Equations: Point-Slope Form**

Point-slope refers to a method for graphing a linear equation on an x-y axis. When graphing a linear equation, the whole idea is to take pairs of x's and y's and plot them on the graph. While you could plot several points by just plugging in values of x, the point-slope form makes the whole process simpler. Point-slope form is also used to take a graph and find the equation of that particular line.

Point slope form gets its name because it uses a single point on the graph and the slope of the line.

The standard point-slope equation looks like this:
$$y - y_1 = m(x - x_1)$$

Using this formula, If you know:
- one point on the line
- and the slope of the line,
  you can find other points on the line.

Example 1 :

You are given the point (4,3) and a slope of 2. Find the equation for this line in point slope form.

Just substitute the given values into your point-slope formula above $y - y_1 = m(x-x_1)$. Your point (4,3) is in the form of $(x_1, y_1)$. That means where you see $y_1$, use 3. Where you see $x_1$, use 4. Slope is given, so where you see m, use 2. So we get  Y − 3 = 2(x–4)

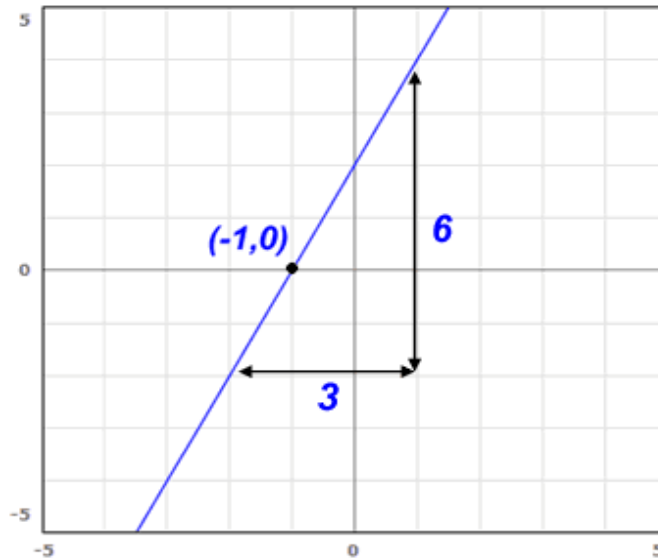Example 2 : You are given the point (– 1,5)and a slope of  1/2. Find the equation for this line in point slope form.

Substituting in the formula $y - y_1 = m(x–x_1)$ we get $y - 5 = 0.5(x- -1)$. Ie. $y - 5 = 0.5(x+1)$.

Point-slope form is about having a single point and a direction (slope) and converting that between an algebraic equation and a graph. In the example above, we took a given set of point and slope and made an equation. Now let's take an equation and find out the point and slope so we can graph it.

Example:     Find the equation (in point-slope form) for the line shown in this graph:



To write the equation, we need a point, and a slope. To find a point because we just need <u>any</u> point on the line. This means you can select any point as per your convenience. The point indicated in the figure is (–1,0). We have chosen this point for our convenience since it is the easiest one to find. We choose it also because  it is useful to pick a point on the axis, because one of the values will be zero.

Now let us find slope. Just count the number of lines on the graph paper going in each direction of a triangle, like shown in the figure. Remember that slope is rise over run, or y/x. Therefore the slope of this line is 2. Putting it all together, our point is (–1,0) and our slope is 2, the point-slope form is $y - 0 = 2(x + 1)$

Example:        Write the equation of the line passing through the points (9,-17) and (-4,4).

Calculate the slope:

$$m = \frac{(-17-4)}{(9+4)}$$

$$m = -21 / 13$$

Use the point-slope formula:

$$y - 4 = (-21/13)(x + 4)$$
$$y - 4 = (-21/13)x - 84/13$$
$$y = (-21/13)x - 32/13$$

**Straight-Line Equations: Two Point Form**

Two Point Slope Form is used to generate the Equation of a straight line passing through the two given points. The two-point form of a line in the Cartesian plane passing through the points $(x_1,y_1)$ and $(x_2,y_2)$ is given by

$$y - y_1 = \frac{y_2 - y_1}{x_2 - x_1}(x - x_1)$$

or equivalently,

$$y - y_2 = \frac{y_2 - y_1}{x_2 - x_1}(x - x_2)$$

Or alternatively

$$\frac{x - x_1}{x_2 - x_1} = \frac{y - y_1}{y_2 - y_1}$$

Two point form diagram is as follows.



Example 1:  Determine the equation of the line that passes through the points A = (1, 2) and B = (−2, 5).

$$\frac{x - 1}{-2 - 1} = \frac{y - 2}{5 - 2}$$

Example 2:  Find the equation of the line passing through (3, − 7) and (− 2,− 5).

      Solution : The equation of a line passing through two points $(x_1, y_1)$ and $(x_2, y_2)$ is given by

$$y - y_1 = \frac{y_2 - y_1}{x_2 - x_1}(x - x_1)$$

Since $x_1 = 3$, $y_1 = − 7$ and $x_2 = − 2$, and $y_2 = − 5$, equation becomes,

$$y - -7 = \frac{-5 - -7}{-2 - 3}(x - 3)$$

$$y + 7 = \frac{-5 + 7}{-2 - 3}(x - 3)$$

Or

$$y + 7 = \frac{2}{-5}(x - 3)$$

Or

$$2x + 5y + 29 = 0$$

**Straight-Line Equations: Intercept Form**

      The intercept form of the line is the equation of the line segment based on the intercepts with both axes. The intercept form of a line in the Cartesian plane with x intercept a and y intercept b is given by

$$\frac{x}{a}+\frac{y}{b}=1$$



$$\frac{x}{a}+\frac{y}{b}=1$$

a is the x-intercept. b is the y-intercept. a and b must be nonzero.

The values of a and b can be obtained from the general form equation.

If y = 0, x = a.

If x = 0, y = b.

A line does not have an intercept form equation in the following cases:

1.A line parallel to the x-axis, which has the equation y = k.

2.A line parallel to the x-axis, which has the equation x = k.

3.A line that passes through the origin, which has equation y = mx.

Examples

1. A line has an x-intercept of 5 and a y-intercept of 3. Find its equation.

$$\frac{x}{a}+\frac{y}{b}=1$$

$$\frac{x}{5}+\frac{y}{3}=1$$

2. The line x − y + 4 = 0 forms a triangle with the axes. Determine the area of the triangle.

The line forms a right triangle with the origin and its legs are the axes.

If y = 0 $\Rightarrow$ x = −4 = **a**.

If x = 0 $\Rightarrow$ y = 2 = **b**.

$$\frac{x}{-4} + \frac{y}{2} = 1$$

The area is:

$$s = \frac{1}{2}\left|(-4).2\right| = 4u^2$$

3. Find the equation of a line which cuts off intercepts 5 and –3 on x and y axes respectively.

The intercepts are 5 and –3 on x and y axes respectively. i.e., a = 5, b = – 3

The required equation of the line is

$$\frac{x}{5} + \frac{y}{-3} = 1$$

3 x – 5y –15 = 0

4. Find the equation of a line which passes through the point (3, 4) and makes intercepts on the axes equl in magnitude but opposite in sign.

Solution : Let the x-intercept and y-intercept be a and –a respectively

∴ The equation of the line is

$$\frac{x}{a} + \frac{y}{-a} = 1$$

x – y = a … (i)

Since (i) passes through (3, 4)

∴ 3 – 4 = a or

a = –1

Thus, the required equation of the line is

x – y = – 1

or  x – y + 1 = 0

4. Determine the equation of the line through the point (– 1,1) and parallel to   x – axis
Since the line is parallel to x-axis its slope ia zero. Therefore from the point slope
form of the equation, we get

y – 1 = 0 [ x – (– 1)]

y – 1 = 0

which is the required equation of the given line.

**Straight-Line Equations: Standard Form**

In the Standard Form of the equation of a straight line, the equation is expressed as:

ax + by = c where a and b are not both equal to zero.

1.  7x + 4y = 6

2.  2x – 2y = –2

3. –4x + 17y = –432

**Exercise:**

**1.** Graph a linear function $y=\dfrac{-1}{4}x+3$

   To draw graph of this function we should first find two points (the x intercept and the y intercept) which satisfy the equation. Then we join these two points using a straight line. This is because what we are given is the equation of linear function. For a linear function, all the points satisfying the equation must lie on the line.

   To find the y intercept, let us use our knowledge of the definition of intercept, ie, the y intercept is the point where the line touches the y axis. This means at this point x = 0. So to find the y intercept, set x equal to zero.

$y=\dfrac{-1}{4}(0)+3$

$y=3$

So the y intercept is (x,y) = (0,3)
Similarly to find the x intercept, put y = 0.

$0=\dfrac{-1}{4}x+3$

$\dfrac{1}{4}x=3$

$x=12$

So the x intercept is (x,y) = (10,0)
Now plot the two intercept points (0,3) and (12,0) and connect them with a straight line.

**2.** Graph the linear function 2y + 10x = 20

   To find the y intercept, put x = 0

   2y + 10(0) = 20, 2y = 20, y = 10

   So the y intercept is (x,y) = (0,10)

   To find the x intercept, put y = 0

   2(0) + 10x = 20, 10x = 20, x = 2

   So the x intercept is (x,y) = (2,0)

   Now plot the two intercept points (0,10) and (2,0) and connect them with a straight line.

Alternatively, you may first of all rewrite the equation 2y + 10x = 20 by solving for Q,

2y = 20 – 10x,   $y=\dfrac{20-10x}{2}, y=10-5x, that\ is\ y=-5x+10$

Now the function is in y = mx + c from. From this form, we can easily find the slope (m)of the function which is – 5.

**3.** Graph the linear function 2y -10x = 20

To find the y intercept, put x = 0

2y + 10(0) = 20, 2y = 20, y = 10

So the y intercept is (x,y) = (0,10)

To find the x intercept, put y = 0

2(0) – 10x = 20, –10x = 20, x = –2

So the x intercept is (x,y) = (–2,0)

Now plot the two intercept points (0,10) and (–2,0) and connect them with a straight line.

**4.** Graph the linear function 2y – 10x + 20  = 0
To find the y intercept, put x = 0
2y – 10(0) – 20 = 0, 2y = –20, y = –10
So the y intercept is (x,y) = (0, –10)
To find the x intercept, put y = 0
2(0) – 10x + 20 = 0, –10x = –20, x = 2
So the x intercept is (x,y) = (2,0)
Now plot the two intercept points (0, –10) and (2,0) and connect them with a straight line.

**5.** Given a linear function 6y + 3x – 18 = 0. Also find its slope.
Rewrite in the y = mx + c form
$$6y + 3x - 18 = 0$$
$$6y = -\ 3x + 18$$
$$y = \frac{-1}{2}x + 3$$

Slope = $\frac{-1}{2}x$

To find the y intercept, put x = 0
$$y = \frac{-1}{2}(0) + 3 \quad , y = 3$$

So the y intercept is (x,y) = (0, 3)
To find the x intercept, put y = 0
$$0 = \frac{-1}{2}x + 3$$

$$\frac{1}{2}x = 3$$

$$x = \frac{2}{1} \times 3 = 2 \times 3 = 6$$

So the x intercept is (x,y) = (6,0)
Now plot the two intercept points (4, 12) and (8,2) and connect them with a straight line.

**6.** Find the slope *m* of a linear function passing through (0,10), (2,0)
Substituting in the formula

$$m=\frac{y_2-y_1}{x_2-x_1}=\frac{2-12}{8-4}=\frac{-10}{4}$$

## Graphing of Non Linear Functions

So far we have been seeing graphing of linear functions which gives a straight line. Now let us see the graphing of some non linear functions.

Graphs of Quadratic Functions

Quadratic Function

The term quadratic comes from the word quadrate meaning square or rectangular. Similarly, one of the definitions of the term quadratic is a square. In an algebraic sense, the definition of something quadratic involves the square and no higher power of an unknown quantity; second degree. So, for our purposes, we will be working with quadratic equations which mean that the highest degree we'll be encountering is a square. Normally, we see the standard quadratic equation written as the sum of three terms set equal to zero. Simply, the three terms include one that has an $x^2$, one has an x, and one term is "by itself" with no $x^2$ or x.

A quadratic function in its normal form is written in the form    $f(x)=ax^2+bx+c$

where a, b, and c are constants and a is not equal to zero. If a = 0, the $x^2$ term would disappear and we would have a linear equation. Note that in a quadratic function there is a power of two on the independent variable and that is the highest power.

The graph of a quadratic function is called a **parabola.** It is basically a curved shape opening up or down. When we have a quadratic function in either form,   $f(x)=ax^2+bx+c$

if a > 0, then the parabola opens up $\cup$ , and , if a < 0, then the parabola opens down $\cap$ .

To make things simple, let us consider a simple quadratic function where a = 1, b = 0 and c = 0. So we get the normal quadratic equation as y = $1x^2$ or y = $x^2$. To graph this function, let us try substituting values in for x and solving for y as shown in the following table.

| x | y = $x^2$ | y = $x^2$ | (x, y) |
|---|---|---|---|
| $-3$ | (-3)2 | 9 | (-3, 9) |
| $-2$ | (-2)2 | 4 | (-2, 4) |
| $-1$ | (-1)2 | 1 | (-1, 1) |
| 0 | (0)2 | 0 | (0, 0) |
| 1 | (1)2 | 1 | (1, 1) |
| 2 | (2)2 | 4 | (2, 4) |
| 3 | (3)2 | 9 | (3, 9) |

Plot the graph on a graph paper and we will get the following graph.

As stated earlier, If a>0, then the parabola has a minimum point and it opens upwards (U-shaped) For example see the following graph for the function $y=x^2+2x-3$



Similarly, If a<0, then the parabola has a maximum point and it opens downwards (n-shaped) For example see below the graph of the function $y=-2x^2+5x+3$

## MODULE IV
## Meaning of Statistics and Description of Data

### 1. DEFINITION, SCOPE AND LIMITATIONS OF STATISTICS

The word 'Statistics' is derived from the Latin word Status, means a political state. The theory of statistics as a distinct branch of scientific method is of comparatively recent growth. Research particularly into the mathematical theory of statistics is rapidly proceeding and fresh discoveries are being made all over the world.

Statistics is concerned with scientific methods for collecting, organizing, summarizing, presenting and analyzing data well as deriving valid conclusions and making reasonable decisions on the basis of this analysis. Statistics is concerned with the systematic collection of numerical data and its interpretation.

The word 'statistic' is used to refer to

1. Numerical facts, such as the number of people living in particular area.

2. The study of ways of collecting, analyzing and interpreting the facts.

### 1.1 Definitions

Statistics is defined differently by different authors over period of time.  In the olden days statistics was confined to only state affairs but in modern days it embraces almost every sphere of human activity. Therefore a number of old definitions, which was confined to narrow field of enquiry were replaced by more definitions, which are much more comprehensive and exhaustive. Secondly, statistics has been defined in two different ways–Statistical data and statistical methods.

1.  Statistics can be defined as the collection presentation and interpretation of numerical data- Croxton and Cowden.

2.  Statistics are numerical statement of facts in any department of enquiry placed interrelation to each other.- Bowley.

3.  Statistics are measurement, enumerations or estimates of natural or social phenomena systematically arrangement to exhibit their inner relation.- Conner.

4. By Statistics we mean quantitative data affected to a marked extend by a multiplicity of causes. – Youle and Kendal.

5.  The science of Statistics is essentially a branch of applied mathematics and can be regarded as a mathematics applied to observation data. - R.A fisher.

Statistics can be defined in two senses i.e. singular and plural. In singular sense it may be defined as the various methods and techniques for attaining and analyzing the numerical

information. Different economists have different view about statistics. According to Boddingtons Statistics is, "the science of estimates and probabilities". The techniques and method means the collection of data, organization, presentation, analysis and interpretation of numerical data. The above definition covers the following aspects of statistics.

1. Collection of data: The collection of data is the first step of statistical investigation. It must be collected very carefully. So, the data must be covered, if not the conclusion will not be reliable.

2. Organization: The data may be obtained either from primary source or the secondary source. If the data is to be obtained from the primary source, then it needs organization. The data are organized by editing, classifying and tabulating them.

3. Presentation: After the collection and organization of data, they are presented in systematic form such as table, diagram and graphical form.

4. Analysis: After the collection, organization and presentation of data, the next step is to analyze the data. To analyze the data we use average, correction, regression, time series etc. The statistical tools of analysis depend upon the nature of data.

5. Interpretation: The last step of a statistical method is the interpretation of the result obtained from the analysis. Interpretation means to draw the valid conclusion.

**1.2 Characteristics of Statistics**

1. Statistics are aggregate of facts: A single age of 20 or 30 years is not statistics, a series of ages are. Similarly, a single figure relating to production, sales, birth, death etc., would not be statistics although aggregates of such figures would be statistics because of their comparability and relationship.

2. Statistics are affected to a marked extent by a multiplicity of causes: A number of causes affect statistics in a particular field of enquiry, e.g., in production statistics are affected by climate, soil, fertility, availability of raw materials and methods of quick transport.

3. Statistics are numerically expressed, enumerated or estimated: The subject of statistics is concerned essentially with facts expressed in numerical form -with their quantitative details but not qualitative descriptions. Therefore, facts indicated by terms such as 'good', 'poor' are not statistics unless a numerical equivalent is assigned to each expression. Also this may either be enumerated or estimated, where actual enumeration is either not possible or is very difficult.

4. Statistics are numerated or estimated according to reasonable standard of accuracy: Personal bias and prejudices of the enumeration should not enter into the counting or estimation of figures, otherwise conclusions from the figures would not be accurate. The figures should be counted or estimated according to reasonable standards of accuracy. Absolute accuracy is neither necessary nor sometimes possible in social sciences. But whatever standard of accuracy is once adopted, should be used throughout the process of collection or estimation.

5. Statistics should be collected in a systematic manner for a predetermined purpose: The statistical methods to be applied on the purpose of enquiry since figures are always collected with some purpose. If there is no predetermined purpose, all the efforts in collecting the figures may prove to be wasteful. The purpose of a series of ages of husbands and wives may be to find whether young husbands have young wives and the old husbands have old wives.

6. Statistics should be capable of being placed in relation to each other: The collected figure should be comparable and well-connected in the same department of inquiry. Ages of husbands are to be compared only with the corresponding ages of wives, and not with, say, heights of trees.

## 1.3 Functions of statistics

### 1. Statistics enable realization magnitude

Bare statement of facts relating to a phenomenon not expressed in numbers enable us no doubt, to visualize the whole picture, but we cannot get an idea of the magnitude involved. It is only when such facts are expressed in numbers that we can get such an idea. This is what statistics

### 2. Statistics simplifies complexity

A huge mass of complicated data relating to a phenomenon is confusing to human mind-human mind is not able to assimilate complicated data. By the application of appropriate statistically methods complex data can be condensed into a few and simple numerical expressions and these the human mind can easily graph. Statistical measure like average .these measure describe a phenomenon in a simple way and bring to light the fundamental features of the data relating to the phenomenon.

### 3. Statistics enables comparison of simplified data

It is with the help of statistical methods such as the methods of presentation like diagram and graph and statistical measure like average, index numbers, measure of variation etc., that one quantity or estimate as compared with another and variation or difference noted between one aspects of a phenomenon and the other.

### 4. Statistics enables to study relationships between sets of related phenomenon

In all types of studies- national, social, economic, business, etc., the importance of observing between different phenomena is very great. For example, it may be necessary, for some practical purpose to study whether there is any relationship between say, heights and weights persons, between ages and sizes of shoes worn by them, etc the presence or absence and the kinds and the degree of relationships can be studies by the application of statistical methods.

### 5. Statistics enlarge human experience

The use and application of statistical methods enlarge human knowledge and experience by making it easier for him to understand and measure phenomena relating to practically all fields of human knowledge- naturals science, social science etc. many fields of knowledge which hither to closed to mankind, have been opened up by the application of the statistical techniques.

6. **Forecasting**

By the word forecasting, we mean to predict or to estimate beforehand. Given the data of the last ten years connected to rainfall of a particular district in Kerala, it is possible to predictor forecast the rainfall for the near future. In business also forecasting plays a dominant role in connection with production, sales, profits etc. The analysis of time series and regression analysis plays an important role in forecasting.

7. **Comparison**:

Classification and tabulation are the two methods that are used to condense the data. They help us to compare data collected from different sources. Grand totals, measures of central tendency measures of dispersion, graphs and diagrams, coefficient of correlation etc provide ample scope for comparison. If we have one group of data, we can compare within it. If the rice production (in Tonnes) in Palakkad district is known, then we can compare it with another district in the state. Or if the rice production (in Tonnes) of two different districts within Kerala is known, then also a comparative study can be made. As statistics is an aggregate of facts and figures, comparison is always possible and in fact comparison helps us to understand the data in a better way.

8. **Estimation**:

One of the main objectives of statistics is drawing  inference about a population from the analysis for the sample drawn from that population. The four major branches of statistical inference are1.Estimation theory 2.Tests of Hypothesis  3.Non Parametric tests     4.Sequential analysis. In estimation theory, we estimate the unknown value of the population parameter based on the sample observations. Suppose we are given a sample of heights of hundred students in a school, based upon the heights of these 100 students, it is possible to estimate the average height of all students in that school.

9. **Tests of Hypothesis**

A statistical hypothesis is some statement about the probability distribution, characterizing a population on the basis of the information available from the sample observations.   In the formulation and testing of hypothesis, statistical methods are extremely useful.  Whether crop yield has increased because of the use of new fertilizer or whether the new medicine is effective in eliminating a particular disease are some examples of statements of hypothesis and these are tested by proper statistical tools.

**1.4 Uses of Statistics**

Statistics is primarily used either to make predictions based on the data available or to make conclusions about a population of interest when only sample data is available. In both cases statistics tries to make sense of the uncertainty in the available data.

Statisticians apply statistical thinking and methods to a wide variety of scientific, social, and business endeavours in such areas as astronomy, biology, education, economics, engineering, genetics, marketing, medicine, psychology, public health, sports, among many. Many economic, social, political, and military decisions cannot be made without statistical techniques, such as the design of experiments to gain federal approval of a newly manufactured drug.

Statistics is of two types (a) Descriptive statistics involves methods of organizing, picturing and summarizing information from data. (b) Inferential statistics involves methods of using information from a sample to draw conclusions about the population.

These days statistical methods are applicable everywhere. There is no field of work in which statistical methods are not applied. According to A L. Bowley, 'knowledge of statistics is like knowledge of foreign languages or of Algebra; it may prove of use at any time under any circumstances". The importance of the statistical science is increasing in almost all spheres of knowledge, e g., astronomy, biology, meteorology, demography, economics and mathematics. Economic planning without statistics is bound to be baseless. Statistics serve in administration, and facilitate the work of formulation of new policies. Financial institutions and investors utilize statistical data to summaries the past experience. Statistics are also helpful to an auditor, when he uses sampling techniques or test checking to audit the accounts of his client.

(a) **Statistics and Economics**: In the year 1890 Alfred Marshall, the renowned economist observed that "statistics are the straw out of which I, like every other economist, have to make bricks". This proves the significance of statistics in economics. Economics is concerned with production and distribution of wealth as well as with the complex institutional set-up connected with the consumption, saving and investment of income. Statistical data and statistical methods are of immense help in the proper understanding of the economic problems and in the formulation of economic policies. In fact these are the tools and appliances of an economist's laboratory. In the field of economics it is almost impossible to find a problem which does not require an extensive uses of statistical data. As economic theory advances use of statistical methods also increase. The laws of economics like law of demand, law of supply etc can be considered true and established with the help of statistical methods. Statistics of consumption tells us about the relative strength of the desire of a section of people. Statistics of production describe the wealth of a nation. Exchange statistics through light on commercial development of a nation. Distribution statistics disclose the economic conditions of various classes of people. There for statistical methods are necessary for economics.

(b) **Statistics and business**: Statistics is an aid to business and commerce. When a person enters business, he enters into the profession of fore casting. Modern statistical devices have made business forecasting more precise and accurate. A business man needs statistics right from the time he proposes to start business. He should have relevant fact and figures to prepare the financial plan of the proposed business. Statistical methods are necessary for these purposes. In

industrial concern statistical devices are being used not only to determine and control the quality of products manufactured by also to reduce wastage to a minimum. The technique of statistical control is used to maintain quality of products.

(c) **Statistics and Research**: Statistics is an indispensable tool of research. Most of the advancement in knowledge has taken place because of experiments conducted with the help of statistical methods. For example, experiments about crop yield and different types of fertilizers and different types of soils of the growth of animals under different diets and environments are frequently designed and analysed according to statistical methods. Statistical methods are also useful for the research in medicine and public health. In fact there is hardly any research work today that one can find complete without statistical data and statistical methods.

Other uses of statistics are as follows.

(1) Statistics helps in providing a better understanding and exact description of a phenomenon of nature.

(2) Statistical helps in proper and efficient planning of a statistical inquiry in any field of study.

(3) Statistical helps in collecting an appropriate quantitative data.

(4) Statistics helps in presenting complex data in a suitable tabular, diagrammatic and graphic form for an easy and clear comprehension of the data.

(5) Statistics helps in understanding the nature and pattern of variability of a phenomenon through quantitative observations.

(6) Statistics helps in drawing valid inference, along with a measure of their reliability about the population parameters from the sample data.

**1.5 Scope of Statistics:**

"Today, there is hardly a phase of human activity which does not findStatisticaldevices at least occasionally useful. Economics, anthropology, psychology, agriculture, businessand education-all depend heavily upon statistics. Statistical methods are applied to the result of physical chemistry and biological experiments and observation as well to result to obtain in social and economics investigations". It is clear from the above that statistical analysis includes in its fold all quantitative analysis to whatever field of inquiry they might relate. The scope of statistical methods is stretched over all those branches of human knowledge in which a grasp of the significance of large numbers is looked for. The scope of statistical methods therefore, is wide the limiting factor being its applicability to studies of quantitative character only. The statistical methods can be used in studying a problem relating to any phenomenon, provided the problem or its aspects are susceptible to numerical measurement.

1. **Statistics and Industry**:

Statistics is widely used in many industries. In industries, control charts are widely used to maintain a certain quality level. In production engineering, to find whether the product is conforming to specifications or not, statistical tools, namely inspection plans, control charts, etc.,

are of extreme importance. In inspection plans we have to resort to some kind of sampling–a very important aspect of Statistics

2. **Statistics and Commerce:**

Statistics are lifeblood of successful commerce. Any businessman cannot afford to either by under stocking or having overstock of his goods. In the beginning he estimates the demand for his goods and then takes steps to adjust with his output or purchases. Thus statistics is indispensable in business and commerce.

3. **Statistics and Agriculture:**

Analysis of variance (ANOVA) is one of the statistical tools developed by R.A. Fisher, plays a prominent role in agriculture experiments. In tests of significance based on small samples, it can be shown that statistics is adequate to test the significant difference between two sample means. In analysis of variance, we are concerned with the testing of equality of several population means.

4. **Statistics and Economics:**

Statistical methods are useful in measuring numerical changes in complex groups and interpreting collective phenomenon. Nowadays the uses of statistics are abundantly made in any economic study. Both in economic theory and practice, statistical methods play an important role. Alfred Marshall said, "Statistics are the straw only which like every other economist have to make the bricks". It may also be noted that statistical data and techniques of statistical tools are immensely useful in solving many economic problems such as wages, prices, production, distribution of income and wealth and soon. Statistical tools like Index numbers, time series Analysis, Estimation theory, Testing Statistical Hypothesis are extensively used in economics.

5. **Statistics and Planning:**

Statistics is indispensable in planning. In the modern world, which can be termed as the "world of planning", almost all the organisations in the government are seeking the help of planning for efficient working, for the formulation of policy decisions and execution of the same. In order to achieve the above goals, the statistical data relating to production, consumption, demand, supply, prices, investments, income expenditure etc and various advanced statistical techniques for processing, analyzing and interpreting such complex data are of importance. In India statistics play an important role in planning, commissioning both at the central and state government levels.

**1.6 The Use of Statistics in Economics and Other Social Sciences**

Statistics play an important role in economics. Economics largely depends upon statistics. National income accounts are multipurpose indicators for the economists and administrators. Statistical methods are used for preparation of these accounts. In economics research statistical methods are used for collecting and analysis the data and testing hypothesis. The relationship between supply and demands is studies by statistical methods, the imports and

exports, the inflation rate, the per capita income are the problems which require good knowledge of statistics.

Businesses use statistical methodology and thinking to make decisions about which products to produce, how much to spend advertising them, how to evaluate their employees, how often to service their machinery and equipment, how large their inventories should be, and nearly every aspect of running their operations. The motivation for using statistics in the study of economics and other social sciences is somewhat different. The object of the social sciences and of economics in particular is to understand how the social and economic system functions. While our approach to statistics will concentrate on its uses in the study of economics, you will also learn  business uses of statistics because many of the exercises in your text  book, and some of the ones used here, will focus on business problems.

Views and understandings of how things work are called theories. Economic theories are descriptions and interpretations of how the economic system functions. They are composed of two parts—a logical structure which is tautological (that is, true by definition), and a set of parameters in that logical structure which gives the theory empirical content (that is, an ability to be consistent or inconsistent with facts or data). The logical structure, being true by definition, is uninteresting except insofar as it enables us to construct testable propositions about how the economic system works. If the facts turn out to be consistent with the testable implications of the theory, then we accept the theory as true until new evidence inconsistent with it is uncovered. A theory is valuable if it is logically consistent both within itself and with other theories established as "true" and is capable of being rejected by but nevertheless consistent with available evidence. Its logical structure is judged on two grounds—internal consistency and usefulness as a framework for generating empirically testable propositions.

To illustrate this, consider the statement: "People maximize utility."This statement is true by definition—behaviour is defined as what people do and utility is defined as what people maximize when they choose to do one thing rather than something else. These definitions and the associated utility maximizing approach form a useful logical structure for generating empirically testable propositions. One can choose the parameters in this tautological utility maximization structure so that the marginal utility of good declines relative to the marginal utility of other goods as the quantity of those good consumed increases relative to the quantities of other goods consumed. Downward sloping demand curves emerge, leading to the empirically testable statement: "Demand curves slope downward." This theory of demand (which consists of both the utility maximization structure and the proposition about how the individual's marginal utilities behave) can then be either supported or falsified by examining data on prices and quantities and incomes for groups of individuals and commodities. The set of tautologies derived using the concept of utility maximization are valuable because they are internally consistent and generate empirically testable propositions such as those represented by the theory of demand. If it didn't

yield testable propositions about the real world, the logical structure of utility maximization would be of little interest.

### 1.7 Limitations of Statistics

Statistics is indispensable to almost all sciences - social, physical and natural. It is very often used in most of the spheres of human activity. In spite of the wide scope of the subject it has certain limitations. Some important limitations of statistics are the following:

1. Statistics does not study qualitative phenomena: Statistics deals with facts and figures. So the quality aspect of a variable or the subjective phenomenon falls out of the scope of statistics. For example, qualities like beauty, honesty, intelligence etc. cannot be numerically expressed. So these characteristics cannot be examined statistically. This limits the scope of the subject.

2. Statistical laws are not exact: Statistical laws are not exact as in case of natural sciences. These laws are true only on average. They hold good under certain conditions. They cannot be universally applied. So statistics has less practical utility.

3. Statistics does not study individuals: Statistics deals with aggregate of facts. Single or isolated figures are not statistics. This is considered to be a major handicap of statistics.

4. Statistics can be misused: Statistics is mostly a tool of analysis. Statistical techniques are used to analyse and interpret the collected information in an enquiry. As it is, statistics does not prove or disprove anything. It is just a means to an end. Statements supported by statistics are more appealing and are commonly believed. For this, statistics is often misused. Statistical methods rightly used are beneficial but if misused these become harmful. Statistical methods used by less expert hands will lead to inaccurate results. Here the fault does not lie with the subject of statistics but with the person who makes wrong use of it.

Other limitations are as follows.

(1) Statistics laws are true on average. Statistics are aggregates of facts. So single observation is not a statistics, it deals with groups and aggregates only.

(2) Statistical methods are best applicable on quantitative data.

(3) Statistical cannot be applied to heterogeneous data.

(4) It sufficient care is not exercised in collecting, analysing and interpretation the data, statistical results might be misleading.

(5) Only a person who has an expert knowledge of statistics can handle statistical data efficiently.

(6) Some errors are possible in statistical decisions. Particularly the inferential statistics involves certain errors. We do not know whether an error has been committed or not.

### 2. FREQUENCY DISTRIBUTIONS

Frequency distribution is a specification of the way in which the frequencies of members of a population are distributed according to the values of the variants which they exhibit. For observed data the distribution is usually specified in tabular form, with some grouping for continuous variants.

The frequency distribution or frequency table is a tabular organization of statistical data, assigning to each piece of data its corresponding frequency.

**Types of Frequencies**

(a) Absolute Frequency

The absolute frequency is the number of times that a certain value appears in a statistical study.

It is denoted by $f_i$ .

The sum of the absolute frequencies is equal to the total number of data, which is denoted by N.

$$f_1 + f_2 + f_3 + \ldots + f_n = N$$

This sum is commonly denoted by the Greek letter Σ (capital sigma) which represents 'sum'.

$$\sum_1^n f_i = N$$

**(b) Relative Frequency**

The relative frequency is the quotient between the absolute frequency of a certain value and the total number of data. It can be expressed as a percentage and is denoted by $n_i$ .

$$n_i = \frac{f_i}{N}$$

The sum of the relative frequency is equal to 1.

**(c) Cumulative Frequency**

The cumulative frequency is the sum of the absolute frequencies of all values less than or equal to the value considered.

It is denoted by $F_i$.

**(d) Relative Cumulative Frequency**

The relative cumulative frequency is the quotient between the cumulative frequency of a particular value and the total number of data. It can be expressed as a percentage.

Example

A city has recorded the following daily maximum temperatures during a month:

32, 31, 28, 29, 33, 32, 31, 30, 31, 31, 27, 28, 29, 30, 32, 31, 31, 30, 30, 29, 29, 30, 30, 31, 30, 31, 34, 33, 33, 29, 29.

Let us form a table based on this information. In the first column of the table are the variables ordered from lowest to highest, in the second column is the count or the number or times this variable has occurred and in the third column is the score of the absolute frequency.

| $x_i$ | Count | $f_i$ | $F_i$ | $n_i$ | $N_i$ |
|---|---|---|---|---|---|
| 27 | I | 1 | 1 | 0.032 | 0.032 |
| 28 | II | 2 | 3 | 0.065 | 0.097 |
| 29 | ⊞⊞ I | 6 | 9 | 0.194 | 0.290 |
| 30 | ⊞⊞ II | 7 | 16 | 0.226 | 0.516 |
| 31 | ⊞⊞ III | 8 | 24 | 0.258 | 0.774 |
| 32 | III | 3 | 27 | 0.097 | 0.871 |
| 33 | III | 3 | 30 | 0.097 | 0.968 |
| 34 | I | 1 | 31 | 0.032 | 1 |
| | | 31 | | 1 | |

Discrete variables are used for this type of frequency table.

## 2.1 Graphs of frequency distribution

A frequency distribution can be represented graphically in any of the following ways. The most commonly used graphs and curves for representation a frequency distribution are

Bar Charts

Histogram

Frequency Polygon

Smoothened frequency curve

Ogives or cumulative frequency curves.

**(a)Bar Charts**

A bar chart is used to present categorical, quantitative or discrete data.

The information is presented on a coordinate axis. The values of the variable are represented on the horizontal axis and the absolute, relative or cumulative frequencies are represented on the vertical axis.

The data is represented by bars whose height is proportional to the frequency.

Example

A study has been conducted to determine the blood group of a class of 20 students. The results are as follows:

| Blood Group | $f_i$ |
|---|---|
| A | 6 |
| B | 4 |
| AB | 1 |
| O | 9 |
| **Total** | 20 |

Based on this we can draw a bar chart as follows.

Step 1: Number the Y-axis with the dependent variable. The dependent variable is the one being tested in an experiment. In this sample question, the study wanted to know how many students belonged to each blood group. So the number of students is the dependent variable. So it is marked on the Y-axis.

**Step 2:** Label the X-axis with what the bars represent. For this problem, label the x-axis "Blood Group" and then label the Y-axis with what the Y-axis represents: "number of students."

**Step 3:** Draw your bars. The height of the bar should be even with the correct number on the Y-axis. Don't forget to label each bar under the x-axis.

Finally, give your graph a name. For this problem, call the graph 'Blood group of students'.


## (b) Histogram:

A histogram is a set of vertical bars whose one as are proportional to the frequencies represented. While constructing histogram, the variable is always taken on the X axis and the frequencies on the Y axis. The width of the bars in the histogram will be proportional to the class interval. The bars are drawn without leaving space between them. A histogram generally represents a continuous curve. If the class intervals are uniform for a frequency distribution, then the width of all the bars will by equal.

Example:

| Marks | No. of students |
|-------|-----------------|
| 10-15 | 5 |
| 15-20 | 20 |
| 20-25 | 47 |
| 25-30 | 38 |
| 30-35 | 10 |

```
No.     Y
of    50  -
stud
ents  40  -

      30  -

      20  -

      10  -

      X
   0     5    10   15   20   25   30   35
  Marks
```

## 2.2 Frequency Polygon (or line graphs)

Frequency Polygon is a graph of frequency distribution. Frequency polygons are a graphical device for understanding the shapes of distributions. They serve the same purpose as histograms, but are especially helpful for comparing sets of data.

To create a frequency polygon, start just as for histograms, by choosing a class interval. Then draw an X-axis representing the values of the scores in your data. Mark the middle of each class interval with a tick mark, and label it with the middle value represented by the class. Draw the Y-axis to indicate the frequency of each class. Place a point in the middle of each class interval at the height corresponding to its frequency. Finally, connect the points. You should include one class interval below the lowest value in your data and one above the highest value. The graph will then touch the X-axis on both sides.

Another method of constructing frequency polygon is to take the mid points of the various class intervals and then plot frequency corresponding to each point and to join all these points by straight lines. Here need not construct a histogram:-

Frequency Distribution for internal marks

| Lower Bound | Upper Bound | Frequency | Cumulative Frequency |
|---|---|---|---|
| 49.5 | 59.5 | 5 | 5 |
| 59.5 | 69.5 | 10 | 15 |
| 69.5 | 79.5 | 30 | 45 |
| 79.5 | 89.5 | 40 | 85 |
| 89.5 | 99.5 | 15 | 100 |

Frequency polygon can also be drawn with the help of histogram by joining their mid points of rectangle.



### 2.3 Frequency Curves

Frequency curves are derived from frequency polygons. Frequency curve is obtained by joining the points of frequency polygon by a freehand smoothed curve. Unlike frequency polygon, where the points we joined by straight lines, we make use of free hand joining of those points in order to get a smoothed frequency curve. It is used to remove the ruggedness of polygon and to present it in a good form or shape. We smoothen the angularities of the polygon only without making any basic change in the shape of the curve. In this case also the curve begins and ends at base line, as is in case of polygon. Area under the curve must remain almost the same as in the case of polygon.

### Difference between frequency polygon and frequency curve

Frequency polygon is drawn to frequency distribution of discrete or continuous nature. Frequency curves are drawn to continuous frequency distribution. Frequency polygon is obtained by joining the plotted points by straight lines. Frequency curves are smooth. They are obtained by joining plotted points by smooth curve.

### 2.4 Ogives (Cumulative frequency curve)

A frequency distribution when cumulated, we get cumulative frequency distribution. A series can be cumulated in two ways. One method is frequencies of all the preceding classes one

added to the frequency of the classes. This series is called less than cumulative series. Another method is frequencies of succeeding classes are added to the frequency of a class. This is called more than cumulative series. Smoothed frequency curves drawn for these two cumulative series are called cumulative frequency curve or ogives. Thus corresponding to the two cumulative series we get two ogive curves, known as less than ogive and more than ogive.

Less than ogive curve is obtained by plotting frequencies (cumulated) against the upper limits of class intervals. More than ogive curve is obtained by plotting cumulated frequencies against the lower limits of class intervals. Less than ogive is an increasing curve, slopping – upwards from left to right. More than ogive is a decreasing curve and slopes from left to right.

**Example:**

From less than and more than cumulative frequency distribution given below draw a less than and more then ogive.

| Marks | No. of Students |
|-------|-----------------|
| 10-20 | 4 |
| 20-30 | 6 |
| 30-40 | 10 |
| 40-50 | 20 |
| 50-60 | 18 |
| 60-70 | 2 |

| Marks less than | No. of Students | Marks More than | No. of Students |
|-----------------|-----------------|-----------------|-----------------|
| 10 | 0 | 10 | 60 |
| 20 | 4 | 20 | 56 |
| 30 | 10 | 30 | 50 |
| 40 | 20 | 40 | 40 |
| 50 | 40 | 50 | 20 |
| 60 | 58 | 60 | 2 |
| 70 | 60 | 70 | 0 |

No. of Students

**2.5 Pie Diagrams**

One of the most common ways to represent data graphically is called a pie chart. It gets its name Less than ogive by how it looks, just like a circular pie that has been cut into several slices. This kind of graph is helpful when graphing qualitative data, where the information describes a trait or attribute and is not numerical. Each trait corresponds to a different slice of the pie. By looking at all of the pie pieces, you can compare how much of the data fits in each category.

Pie charts are a form of an area chart that are easy to understand with a quick look. They show the part of the total More than ogive (percentage) in an easy-to-understand way. Pie charts are useful tools that help you figure out and

Marks

understand polls, statistics, complex data, and income or spending. They are so wonderful because everybody can see what is going on.

Pie diagrams are used when the aggregate and their division are to be shown together. The aggregate is shown by means of a circle and the division by the sectors of the circle. For example: to show the total expenditure of a government distributed over different departments like agriculture, irrigation, industry, transport etc. can be shown in a pie diagram. In constructing a pie diagram the various components are first expressed as a percentage and then the percentage is multiplied by 3.6. so we get angle for each component. Then the circle is divided into sectors such that angles of the components and angles of the sectors are equal. Therefore one sector represents one component. Usually components are with the angles in descending order are shown.

**Example:**

You conducted a survey as part of a project work. You had taken a sample of 20 individuals and you want to represent their occupation using a pie chart .

First, put your data into a table, then add up all the values to get a total:

| Farmer | Business | Teacher | Bank | Driver | TOTAL |
|--------|----------|---------|------|--------|-------|
| 4 | 5 | 6 | 1 | 4 | 20 |

Calculate the angle of each sector, using the formula

$$\text{Angle of sector} = \frac{\text{Frequency of data}}{\text{Total frequency}} \times 360°$$

Divide each value by the total and multiply by 100 to get a percent:

| Farmer | Business | Teacher | Bank | Driver | TOTAL |
|--------|----------|---------|------|--------|-------|
| 4 | 5 | 6 | 1 | 4 | 20 |
| 4/20 =**20%** | 5/20 =**25%** | 6/20 =**30%** | 1/20 = **5%** | 4/20 =**20%** | 100% |

Now you need to figure out how many degrees for each 'pie slice' (correctly called a sector).

A Full Circle has 360 degrees, so we do this calculation:

| Farmer | Business | Teacher | Bank | Driver | TOTAL |
|--------|----------|---------|------|--------|-------|
| 4 | 5 | 6 | 1 | 4 | 20 |
| 4/20 =**20%** | 5/20 =**25%** | 6/20 =**30%** | 1/20 = **5%** | 4/20 =**20%** | 100% |
| 4/20 × 360° = **72°** | 5/20 × 360° = **90°** | 6/20 × 360° = **108°** | 1/20 × 360° = **18°** | 4/20 × 360° = **72°** | **360°** |

Draw a circle using a pair of compasses.
Use a protractor to draw the angle for each sector.
Label the circle graph and all its sectors.

Pie charts are to be used with qualitative data, however there are some limitations in using them. If there are too many categories, then there will be a multitude of pie pieces. Some of these are likely to be very skinny, and can be difficult to compare to one another.

If we want to compare different categories that are close in size, a pie chart does not always help us to do this. If one slice has central angle of 30 degrees, and another has a central angle of 29 degrees, then it would be very hard to tell at a glance which pie piece is larger than the other.

## 3. SUMMARY MEASURE OF DISTRIBUTIONS

We will discuss three sets of summary measures namely Measures of Central Tendency, Variability and Shape. These are called summary measures because they summarise the data. For example, one of summary measure very familiar to you is mean. (Mean comes under measure of central tendency.) If we take mean mark of students in a class for a subject, it gives you a rough idea of what the marks are like. Thus based on just one summary value, we get idea of the entire data.

### 3.1 Measures of Central Tendency

A measure of central tendency is a measure that tells us where the middle of a bunch of data lies. A measure of central tendency is a single value that attempts to describe a set of data by identifying the central position within that set of data. As such, measures of central tendency are sometimes called measures of central location. They are also classed as summary statistics. The mean (often called the average) is most likely the measure of central tendency that you are most familiar with, but there are others, such as the median and the mode.

**Arithmetic mean:** Mean is the most common measure of central tendency. It is simply the sum of the numbers divided by the number of numbers in a set of data. This is also known as average.

**Median:** Median is the number present in the middle when the numbers in a set of data are arranged in ascending or descending order. If the number of numbers in a data set is even, then the median is the mean of the two middle numbers.

**Mode**: Mode is the value that occurs most frequently in a set of data.

The mean, median and mode are all valid measures of central tendency, but under different conditions, some measures of central tendency become more appropriate to use than others. In the following sections, we will look at the mean, mode and median, and learn how to calculate them.

We will also discuss Geometric Mean and Harmonic Mean.

**Requisites of a good average**

Since an average is a single value representing a group of values, it is desired that such a value satisfies the following properties.

1.      Easy to understand:- Since statistical methods are designed to simplify the complexities.

2.      Simple to compute: A good average should be easy to compute so that it can be used widely.  However, though case of computation is desirable, it should not be sought at the expense

of other averages. ie, if in the interest of greater accuracy, use of more difficult average is desirable.

3.      Based on all items:- The average should depend upon each and every item of the series, so that if any of the items is dropped, the average itself is altered.

4.      Not unduly affected by Extreme observations:- Although each and every item should influence the value of the average, non of the items should influence it unduly. If one or two very small or very large items unduly affect the average, ie, either increase its value or reduce its value, the average can't be really typical of entire series. In other words, extremes may distort the average and reduce its usefulness.

5.      Rigidly defined: An average should be properly defined so that it has only one interpretation. It should preferably be defined by algebraic formula so that if different people compute the average from the same figures they all get the same answer. The average should not depend upon the personal prejudice and bias of the investigator, other wise results can be misleading.

6.      Capable of further algebraic treatment: We should prefer to have an average that could be used for further statistical computation so that its utility is enhanced. For example, if we are given the data about the average income and number of employees of two or more factories, we should able to compute the combined average.

7.      Sampling stability: Last, but not least we should prefer to get a value which has what the statisticians called "sampling stability". This means that if we pick 10 different group of college students, and compute the average of each group, we should expect to get approximately the same value. It does not mean, however that there can be no difference in the value of different samples. There may be some differences but those samples in which this difference is less that are considered better than those in which the difference is more.

**(a) Mean (Arithmetic mean / average)**

The mean (or average) is the most popular and well known measure of central tendency. It can be used with both discrete and continuous data, although its use is most often with continuous data (see our Types of Variable guide for data types). The mean is equal to the sum of all the values in the data set divided by the number of values in the data set. So, if we have n values in a data set and they have values $x_1$, $x_2$, ..., $x_n$, the sample mean, usually denoted by $\bar{x}$ (pronounced x bar), is:

$$\bar{x} = \frac{(x_1 + x_2 + \cdots + x_n)}{n}$$

This formula is usually written in a slightly different manner using the Greek capitol letter, $\Sigma$, pronounced "sigma", which means "sum of...":

$$\bar{x} = \frac{\sum x}{n}$$

Example

In a survey you collected information on monthly spending for mobile recharge by 20 students of which 10 are male and 10 female. We illustrate below how the data is used to find mean.

| | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | Total | **Mean** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Male | 250 | 150 | 100 | 175 | 150 | 250 | 200 | 200 | 150 | 170 | 1795 | **179.50** |
| Female | 100 | 150 | 150 | 100 | 200 | 150 | 125 | 150 | 130 | 180 | 1435 | **143.50** |
| Both | 350 | 300 | 250 | 275 | 350 | 400 | 325 | 350 | 280 | 350 | 3230 | **161.50** |

First we found the mean for male students. Here $\sum x$= 1795. n =10. So 1795/10 = 179.5.

Similarly, the mean for female students. Here $\sum x$= 1435. n =10. So 1435/10 = 143.5.

We also find the mean for male and female taken together.

Here $\sum x$= 3230. n =20. So 3230/20 = 161.50.

Based on the above we can make certain observations. Male students spend Rs. 179.50 on an average in a month for mobile recharge. Female students spend Rs. 143.50. We may conclude that male students spend more on monthly mobile recharges. As a researcher, you may now use this information to make further studies as to why this is so. What are the factors that make male students to spend more on mobile recharges. We have also calculated the average for all students taken together. It is Rs. 161.50. Thus we observe that the male students spend more than the average for 'all students' while female students spend less than the total for 'all students'.

Mean is also calculated using another method called the shortcut method asexplained below.

<u>Short cut method:</u> The arithmetic mean can also be calculated by short cut method. This method reduces the amount of calculation. It involves the following steps

    i.   Assume any one value as an assumed mean, which is also known as working mean or arbitrary average (A).

    ii.   Find out the difference of each value from the assumed mean (d = X-A).

    iii.   Add all the deviations ($\sum$d)

    iv.   Apply the formula

$$\acute{X} = A + \frac{\sum d}{N}$$

Where $\acute{X} \rightarrow$ Mean, $\frac{\sum d}{N} \rightarrow$ Sum of deviation from assumed mean,

A $\rightarrow$ Assumed mean

Example:

Calculate arithmetic mean

| Roll No : | 1 | 2 | 3 | | 4 | 5 | 6 |
|---|---|---|---|---|---|---|---|
| Marks : | 40 | 50 | 55 | | 78 | 58 | 60 |

| Roll Nos. | Marks | d = X - 55 |
|---|---|---|
| 1 | 40 | -15 |

| | | |
|---|---|---|
| 2 | 50 | -5 |
| 3 | 55 | 0 |
| 4 | 78 | 23 |
| 5 | 58 | 3 |
| 6 | 60 | 5 |
| | $\sum d = 11$ | |

$$\acute{X} = A + \frac{\sum d}{N}$$

$$= 55 + \frac{11}{6} = \underline{56.83}$$

**Calculation of arithmetic mean** - Discrete series

To find out the total items in discrete series, frequency of each value is multiplies with the respective size.  The value so obtained are totaled up.  This total is then divided by the total number of frequencies to obtain arithmetic mean.

Steps
1. Multiply each size of the item by its frequency fX
2. Add all fX – ($\sum f X$)
3. Divide $\sum fX$ by total frequency (N).

The formula is $\acute{X} = \frac{\sum fX}{N}$

Example

| X | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| f | 10 | 12 | 8 | 7 | 11 |

Solution

| X | f | fX |
|---|---|---|
| 1 | 10 | 10 |
| 2 | 12 | 24 |
| 3 | 8 | 24 |
| 4 | 7 | 28 |
| 5 | 11 | 55 |
| | N = $\sum$fX = 141 | |

$$\acute{X} = \frac{\sum fX}{N} = \frac{141}{4.8}$$

$$= \underline{2.93}$$

**Short cut Method**

Steps:

- Take the value of assumed mean (A)
- Find out deviations of each variable from A ie d.
- Multiply d with respective frequencies (fd)
- Add up the product ($\sum$fd)
- Apply formula

$$\acute{X} = A \pm \frac{\sum fd}{N}$$

**Continuous series**

In continuous frequency distribution, the value of each individual frequency distribution is unknown. Therefore an assumption is made to make them precise or on the assumption that the frequency of the class intervals is concentrated at the centre that the mid point of each class intervals has to be found out. In continuous frequency distribution, the mean can be calculated by any of the following methods.

a. Direct method
b. Short cut method
c. Step deviation method

a. **Direct Method**

Steps:

1. Find out the mid value of each group or class. The mid value is obtained by adding the lower and upper limit of the class and dividing the total by two. (symbol = m)
2. Multiply the mid value of each class by the frequency of the class. In other words m will be multiplied by f.
3. Add up all the products - $\sum$fm
4. $\sum$fm is divided by N

Example:

From the following find out the mean profit

| Profit/Shop: | 100-200 | 200-300 | 300-400 | 400-500 | 500-600 | 600-700 | 700-800 |
|---|---|---|---|---|---|---|---|
| No. of shops: | 10 | 18 | 20 | 26 | 30 | 28 | 18 |

Solution

| Profit (₹) | Mid point - m | No of Shops (f) | fm |
|---|---|---|---|
| 100-200 | 150 | 10 | 1500 |
| 200-300 | 250 | 18 | 4500 |
| 300-400 | 350 | 20 | 7000 |
| 400-500 | 450 | 26 | 11700 |
| 500-600 | 550 | 30 | 16500 |
| 600-700 | 650 | 28 | 18200 |
| 700-800 | 750 | 18 | 13500 |
|  |  | $\sum$f = 150 | $\sum$fm = 72900 |

$$\acute{X} = \frac{\sum fd}{N}$$

$$\frac{72900}{150} = \underline{486}$$

b) **Short cut method**

Steps:
1. Find the mid value of each class or group (m)
2. Assume any one of the mid value as an average (A)
3. Find out the deviations of the mid value of each from the assumed mean (d)
4. Multiply the deviations of each class by its frequency (fd).
5. Add up the product of step 4 - $\sum fd$
6. Apply formula

$$\acute{X} = A + \frac{\sum fd}{N}$$

Example: (solving the last example)
Solving: Calculation of Mean

| Profit (₹) | m | d = m - 450 | f | fd |
|------------|-----|-------------|-----|------|
| 100-200 | 150 | -300 | 10 | -3000 |
| 200-300 | 250 | -200 | 18 | -3600 |
| 300-400 | 350 | -100 | 20 | -2000 |
| 400-500 | 450 | 0 | 26 | 0 |
| 500-600 | 550 | 100 | 30 | 3000 |
| 600-700 | 650 | 200 | 28 | 5600 |
| 700-800 | 750 | 300 | 18 | 5400 |
| | | | $\sum f = 150$ | $\sum fd = 5400$ |

$$\acute{X} = A + \frac{\sum fd}{N}$$

$$= 450 + \frac{5400}{150} = \underline{486}$$

c) **Step deviation method**

The short cut method discussed above is further simplified or calculations are reduced to a great extent by adopting step deviation methods.

Steps:
1. Find out the mid value of each class or group (m)
2. Assume any one of the mid value as an average (A)
3. Find out the deviations of the mid value of each from the assumed mean (d)
4. Deviations are divided by a common factor (d')
5. Multiply the d' of each class by its frequency (f d')
6. Add up the products ($\sum fd'$)
7. Then apply the formula

$$\acute{X} = A + \frac{\sum fd'}{N} \times c \qquad \text{Where c = Common factor}$$

Example:

Calculate mean for the last problem

Solution

| Profit | m | f | d | d' | f d' |
|--------|-----|------------|------|----|-------------|
| 100-200 | 150 | 10 | -300 | -3 | -30 |
| 200-300 | 250 | 18 | -200 | -2 | -36 |
| 300-400 | 350 | 20 | -100 | -1 | -20 |
| 400-500 | 450 | 26 | 0 | 0 | 0 |
| 500-600 | 550 | 30 | 100 | 1 | 30 |
| 600-700 | 650 | 28 | 200 | 2 | 56 |
| 700-800 | 750 | 18 | 300 | 3 | 54 |
| | | $\sum f$ = 150 | | | $\sum f d'$ = 540 |

$$\acute{X} = A + \frac{\sum fd'}{N} \times c$$

$$450 + \frac{540}{150} \times 100$$

$$450 + (0.36 \times 100) = \underline{486}$$

The mean is essentially a model of your data set. It is the value that is most common. You will notice, however, that the mean is not often one of the actual values that you have observed in your data set. However, one of its important properties is that it minimises error in the prediction of any one value in your data set. That is, it is the value that produces the lowest amount of error from all other values in the data set.

An important property of the mean is that it includes every value in your data set as part of the calculation. In addition, the mean is the only measure of central tendency where the sum of the deviations of each value from the mean is always zero.

We complete our discussion on arithmetic mean by listing the merits and demerits of it.

**Merits:**

- It is rigidly defined.
- It is easy to calculate and simple to follow.
- It is based on all the observations.
- It is determined for almost every kind of data.
- It is finite and not indefinite.
- It is readily put to algebraic treatment.
- It is least affected by fluctuations of sampling.

**Demerits:**

- The arithmetic mean is highly affected by extreme values.
- It cannot average the ratios and percentages properly.
- It is not an appropriate average for highly skewed distributions.
- It cannot be computed accurately if any item is missing.

- The mean sometimes does not coincide with any of the observed value.

We elaborate on only one of the demerits for your better understanding. The first demerit says the arithmetic mean is highly affected by extreme values. What does this mean. See the following example.

Consider the following table which gives information on the marks obtained by students in a test.

| Student: | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|----------|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Mark : | 15 | 18 | 16 | 14 | 15 | 15 | 12 | 17 | 90 | 95 |

The mean mark for these ten students is 30.7. However, inspecting the raw data suggests that this mean value might not be the best way to accurately reflect the typical mark obtained by a student, as most students have marks in the 12 to 18 range. Here we see that the mean is being affected by the two large figures 90 and 95. This shows that arithmetic mean is highly affected by extreme values.

Therefore, in this situation, we would like to have a better measure of central tendency. As we will find out later, taking the median would be a better measure of central tendency in this situation.

**Weighted Mean (weighted averages)**

Simple arithmetic mean gives equal importance to all items. Sometimes the items in a series may not have equal importance. So the simple arithmetic mean is not suitable for those series and weighted average will be appropriate.

Weighted means are obtained by taking in to account these weights (or importance). Each value is multiplied by its weight and sum of these products is divided by the total weight to get weighted mean.

Weighted average often gives a fair measure of central tendency. In many cases it is better to have weighted average than a simple average. It is invariably used in the following circumstances.

1. When the importance of all items in a series are not equal. We associate weights to the items.
2. For comparing the average of one group with the average of another group, when the frequencies in the two groups are different, weighted averages are used.
3. When rations percentages and rates are to be averaged, weighted average is used.
4. It is also used in the calculations of birth and death rate index number etc.
5. When average of a number of series is to be found out together weighted average is used.

Formula: Let $x_1 + x_2 + x_3 \ldots + x_n$ be in values with corresponding weights $w_1 + w_2 + w_3 \ldots + w_n$. Then the weighted average is

$$= \frac{w_1 x_1 + w_2 x_2 + - - - - + w_n x_n}{w_1 + w_2 + - - - + w_n}$$

$$= \frac{\sum wx}{\sum w}$$

**(b) Median**

The median is also a frequently used measure of central tendency. The median is the midpoint of a distribution: the same number of data points is above the median as below it. The median is the middle score for a set of data that has been arranged in order of magnitude.

The median is determined by sorting the data set from lowest to highest values and taking the data point in the middle of the sequence. There is an equal number of points above and below the

median. For example, in the data 7,8,9,10,11, the median is 9; there are two data points greater than this value and two data points less than this value. Thus to find the median, we arrange the observations in order from smallest to largest value. If there is an odd number of observations, the median is the middle value.

If there is an even number of observations, the median is the average of the two middle values. Thus, the median of the numbers 2, 4, 7, 12 is (4+7)/2 = 5.5.

In certain situations the mean and median of the distribution will be the same, and in some situations it will be different. For example, in the data 1,2,3,4,5 the median is 3; there are two data points greater than this value and two data points less than this value. In this case, the median is equal to the mean. But consider the data 1,2,3,4,10. In this dataset, the median still is three, but the mean is equal to 4.

The median can be determined for ordinal data as well as interval and ratio data. Unlike the mean, the median is not influenced by outliers at the extremes of the data set. For this reason, the median often is used when there are a few extreme values that could greatly influence the mean and distort what might be considered typical. For data which is very skewed, the median often is used instead of the mean.

## Calculation of Median : Discrete series

Steps:
- Arrange the date in ascending or descending order
- Find cumulative frequencies
- Apply the formula Median

$$\text{Median} = \text{Size of } \left[\frac{N+1}{2}\right]^{th} \text{ item}$$

Example: Calculate median from the following

| Size of shoes: | 5 | 5.5 | 6 | 6.5 | 7 | 7.5 | 8 |
|---|---|---|---|---|---|---|---|
| Frequency : | | 10 | 16 | 28 | 15 | 30 | 40 | 34 |

Solution

| Size | f | Cumulative f  (f) |
|---|---|---|
| 5 | 10 | 10 |
| 5.5 | 16 | 26 |
| 6 | 28 | 54 |
| 6.5 | 15 | 69 |
| 7 | 30 | 99 |
| 7.5 | 40 | 139 |
| 8 | 34 | 173 |

$$\text{Median} = \text{Size of } \left[\frac{N+1}{2}\right]^{th} \text{ item}$$

N = 173

$$\text{Median} = \frac{173+1}{2} = 87^{th} \text{ item} = 7$$

Median = 7

## Calculation of median – Continuous frequency distribution

Steps:
- Find out the median by using N/2
- Find out the class which median lies
- Apply the formula

$$Median = L + \frac{h}{f}\left(\frac{N}{2} - C\right)$$

Where L = lower limit of the median class

h = class interval of the median class

f = frequency of the median class

N = $\sum f$, is the total frequency

c = cumulative frequency of the preceding median class

Example: Calculate median from the following data

| Age in years | Below 10 | Below 20 | Below 30 | Below 40 | Below 50 | Below 60 | Below 70 | 70 and over |
|---|---|---|---|---|---|---|---|---|
| No. of persons | 2 | 5 | 9 | 12 | 14 | 15 | 15.5 | 15.6 |

Solution:

First we have to convert the distribution to a continuous frequency distribution as in the following table and then compute median.

| Age in years | No. of persons (f) | Cumulative frequency (cf) – less than |
|---|---|---|
| 0-10 | 2 | 2 |
| 10-20 | 5-2=3 | 5 |
| **20-30** | 9-5=**4** | 9 |
| 30-40 | 12-9=3 | 12 |
| 40-50 | 14-12=2 | 14 |
| 50-60 | 15-14=1 | 15 |
| 60-70 | 15.5-15=0.5 | 15.5 |
| 70 and above | 15.6-15.5=0.1 | 15.6 |
| | $N = \sum f = 15.6$ | |

Median item = $\dfrac{N}{2} = \dfrac{15.6}{2} = 7.8$

Find the cumulative frequency (*c.f*) greater than 7.8 is 9. Thus the corresponding class 20-30 is the median class.

$$Here\ L = 20, h = 10, f = 4, N = 15.6, C = 5$$

Use the formula $$Median = L + \frac{h}{f}\left(\frac{N}{2} - C\right)$$

$$Median = 20 + \frac{10}{4}(7.8 - 5) = 20 + \frac{5}{2} \times 2.8$$

¿$20 + 5 \times 1.4 = 27$.

So the median age is 27.

**The Mean vs Median**

As measures of central tendency, the mean and the median each have advantages and disadvantages. Some pros and cons of each measure are summarized below.

The median may be a better indicator of the most typical value if a set of scores has an outlier. An outlier is an extreme value that differs greatly from other values.

However, when the sample size is large and does not include outliers, the mean score usually provides a better measure of central tendency.

**(b) Mode**

The mode of a data set is the value that occurs with the most frequency. This measurement is crude, yet is very easy to calculate. Suppose that a history class of eleven students scored the following (out of 100) on a test: 60, 64, 70, 70, 70, 75, 80, 90, 95, 95, 100. We see that 70 is in the list three times, 95 occurs twice, and each of the other scores are each listed only once. Since 70 appears in the list more than any other score, it is the mode. If there are two values that tie for the most frequency, then the data is said to be bimodal.

The mode can be very useful for dealing with categorical data. For example, if a pizza shop sells 10 different types of sandwiches, the mode would represent the most popular pizza. The mode also can be used with ordinal, interval, and ratio data. However, in interval and ratio scales, the data may be spread thinly with no data points having the same value. In such cases, the mode may not exist or may not be very meaningful.

To find mode in the case of a continuous frequency distribution, mode is found using the formula

$$Mode = l + \frac{h(f_1 - f_0)}{(f_1 - f_0) - (f_2 - f_1)}$$

Rearranging we get

$$Mode = l + \frac{h(f_1 - f_0)}{2f_1 - f_0 - f_2}$$

Where

$l$ is the lower limit of the model class

$f_1$ is the frequency of the model class

$f_0$ is the frequency of the class preceding the model class

$f_2$ is the frequency of the class succeeding the model class

$h$ is the class interval of the model class

See the following example where we compute mode using the above formula.(mean and median are also computed)

Example: Find the values of mean, mode and median from the following data.

| Weight (kg) | 93-97 | 98-102 | 103-107 | 108-112 | 113-117 | 118-122 | 123-127 | 128-132 |
|---|---|---|---|---|---|---|---|---|
| No. of students | 3 | 5 | 12 | 17 | 14 | 6 | 3 | 1 |

Solution: Since the formula for mode requires the distribution to be continuous with 'exclusive type' classes, we first convert the classes into class boundaries.

| Wight | Class boundaries | Mid value ($X$) | Number of students ($f$) | $d=\dfrac{X-110}{5}$ | $fd$ | Less than c.f |
|---|---|---|---|---|---|---|
| 93-97 | 92.5-97.5 | 95 | 3 | -3 | -9 | 3 |
| 98-102 | 97.5-102.5 | 100 | 5 | -2 | -10 | 8 |
| 103-107 | 102.5-107.5 | 105 | 12 | -1 | -12 | 20 |
| 108-112 | 107.5-112.5 | 110 | 17 | 0 | 0 | 37 |
| 113-117 | 112.5-117.5 | 115 | 14 | 1 | 14 | 51 |
| 118-122 | 117.5-122.5 | 120 | 6 | 2 | 12 | 57 |
| 123-127 | 122.5-127.5 | 125 | 3 | 3 | 9 | 60 |
| 128-132 | 127.5-132.5 | 130 | 1 | 4 | 4 | 61 |
| | | | $N=\sum f=$ | | $N=\sum fd$ | |

## Mean

$$Mean=A+\frac{h\sum fd}{N}$$

$$¿110+\frac{5\times 8}{61}=110.66.$$

Mean = 110.66kgs.

## Mode

Here maximum frequency is 17. The corresponding class 107.5-112.5 is the model class.

Using the formula of mode

$$Mode=l+\frac{h(f_1-f_0)}{2f_1-f_0-f_2}$$

We get

$$Mode=107.5+\frac{5(17-12)}{2(17)-12-14}$$

$$¿107.5+\frac{25}{8}=107.5+3.125=110.625$$

Hence mode is 110.63 kgs.

## Median

Use the formula

$$Median = L + \frac{h}{f}\left(\frac{N}{2} - C\right)$$

Here $\frac{N}{2} = \frac{61}{2} = 30.5$

The cumulative frequency (*c.f.*) just greater than 30.5 is 37. So the corresponding class 107.5-112.5 is the median class.

Substituting values in the median formula

$$Median = 107.5 + \frac{5}{17}\left(\frac{61}{2} - 20\right)$$

$$¿ \, 107.5 + \frac{5}{17}(30.5 - 20)$$

$$¿ \, 107.5 + \frac{5 \times 10.5}{17}$$

$$¿ \, 107.5 + 3.09 = 110.59$$

Median is 110.59 Kgs.

**When to use Mean, Median, and Mode**

The following table summarizes the appropriate methods of determining the middle or typical value of a data set based on the measurement scale of the data.

| Measurement Scale | Best Measure |
|---|---|
| Nominal (Categorical) | Mode |
| Ordinal | Median |
| Interval | Symmetrical data: Mean<br>Skewed data: Median |
| Ratio | Symmetrical data: Mean<br>Skewed data: Median |

**Merits and demerits of mean, median and mode**

Merits and demerits of arithmetic mean has already been discussed. Please refer to that. Here we discuss only median and mode.

**Median:**

The median is that value of the series which divides the group into two equal parts, one part comprising all values greater than the median value and the other part comprising all the values smaller than the median value.

**Merits of median**

(1) Simplicity:- It is very simple measure of the central tendency of the series. I the case of simple statistical series, just a glance at the data is enough to locate the median value.

(2) Free from the effect of extreme values: - Unlike arithmetic mean, median value is not destroyed by the extreme values of the series.

(3) Certainty: - Certainty is another merits is the median. Median values are always a certain specific value in the series.

(4) Real value: - Median value is real value and is a better representative value of the series compared to arithmetic mean average, the value of which may not exist in the series at all.

(5) Graphic presentation: - Besides algebraic approach, the median value can be estimated also through the graphic presentation of data.

(6) Possible even when data is incomplete: - Median can be estimated even in the case of certain incomplete series. It is enough if one knows the number of items and the middle item of the series.

**Demerits of median:**

Following are the various demerits of median:

(1) Lack of representative character: - Median fails to be a representative measure in case of such series the different values of which are wide apart from each other. Also, median is of limited representative character as it is not based on all the items in the series.

(2) Unrealistic:- When the median is located somewhere between the two middle values, it remains only an approximate measure, not a precise value.

(3) Lack of algebraic treatment: - Arithmetic mean is capable of further algebraic treatment, but median is not. For example, multiplying the median with the number of items in the series will not give us the sum total of the values of the series.

However, median is quite a simple method finding an average of a series. It is quite a commonly used measure in the case of such series which are related to qualitative observation as and health of the student.

**Mode**:
The value of the variable which occurs most frequently in a distribution is called the mode.
**Merits of mode:**
Following are the various merits of mode:
(1) Simple and popular: - Mode is very simple measure of central tendency. Sometimes, just at the series is enough to locate the model value. Because of its simplicity, it s a very popular measure of the central tendency.

(2) Less effect of marginal values: - Compared top mean, mode is less affected by marginal values in the series. Mode is determined only by the value with highest frequencies.

(3) Graphic presentation:- Mode can be located graphically, with the help of histogram.

(4) Best representative: - Mode is that value which occurs most frequently in the series. Accordingly, mode is the best representative value of the series.

(5) No need of knowing all the items or frequencies: - The calculation of mode does not require knowledge of all the items and frequencies of a distribution. In simple series, it is enough if one knows the items with highest frequencies in the distribution.

**Demerits of mode:**

Following are the various demerits of mode:
(1) Uncertain and vague: - Mode is an uncertain and vague measure of the central tendency.

(2) Not capable of algebraic treatment: - Unlike mean, mode is not capable of further algebraic treatment.

(3) Difficult: - With frequencies of all items are identical, it is difficult to identify the modal value.

(4) Complex procedure of grouping:- Calculation of mode involves cumbersome procedure of grouping the data. If the extent of grouping changes there will be a change in the model value.

(5) Ignores extreme marginal frequencies:- It ignores extreme marginal frequencies. To that extent model value is not a representative value of all the items in a series.

Besides, one can question the representative character of the model value as its calculation does not involve all items of the series.

**Exercises**

$$Mean = \frac{sum\ of\ all\ values}{total\ number\ of\ values}$$

$$Median = middle\ value\ (when\ the\ data\ are\ arranged\ in\ order)$$

$$Mode = most\ common\ value$$

1. Find the measures of central tendency for the data set 3, 7, 9, 4, 5, 4, 6, 7, and 9.
Mean = 6, median = 6 and modes are 4, 7 and 9.Note that here mode is bimodal.
2. Four friends take an IQ test. Their scores are 96, 100, 106, 114. Which of the following statements is true?
I. The mean is 103.
II. The mean is 104.
III. The median is 100.
IV. The median is 106.
(A) I only
(B) II only
(C) III only
(D) IV only
(E) None is true
The correct answer is (B). The mean score is computed from the equation:

Mean score = $\Sigma x / n$ = (96 + 100 + 106 + 114) / 4 = 104

Since there are an even number of scores (4 scores), the median is the average of the two middle scores. Thus, the median is (100 + 106) / 2 = 103.

3. The owner of a shoe shop recorded the sizes of the feet of all the customers who bought shoes in his shop in one morning. These sizes are listed below:

| 8 | 7 | 4 | 5 | 9 | 13 | 10 | 8 | 8 | 7 | 6 | 5 | 3 | 11 | 10 | 8 | 5 | 4 | 8 | 6 |
|---|---|---|---|---|----|----|---|---|---|---|---|---|----|----|---|---|---|---|---|

What is the mean of these values: 7.25

What is the median of these values: 7.5

What is the mode of these values: 8.

4. Eight people work in a shop. Their hourly wage rates of pay are:

| Worker | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|--------|---|----|---|---|---|---|---|---|
| Wage Rs. | 4 | 14 | 6 | 5 | 4 | 5 | 4 | 4 |

Work out the mean, median and mode for the values above.
Mean = 5.75, Median = 4.50, Mode = 4.00.
Using the above findings, if the owner of the shop wants to argue that the staff are paid well. Which measure would they use? He will use mean. Because mean shows the highest value.
Using the above findings, if the staff in the shop want to argue that they are badly paid. Which measure would they use? The staff will use mode as it is the lowest of the three measures of central tendencies.

5. The table below gives the number of accidents each year at a particular road junction:

| 1991 | 1992 | 1993 | 1994 | 1995 | 1996 | 1997 | 1998 |
|------|------|------|------|------|------|------|------|
| 4 | 5 | 4 | 2 | 10 | 5 | 3 | 5 |

Work out the mean, median and mode for the values above.
Mean =4.75,    Median =4.5,    Mode =5
Using the above measures, a road safety group want to get the council to make this junction safer.

Which measure will they use to argue for this? They will use mode as it is the figure which will help them to justify their argument that the junction has a large number of accidents.

Using the same data the council do not want to spend money on the road junction. Which measure will they use to argue that safety work is not necessary? The council will use median as this figure will help them to argue that the junction has less number of accidents.

6. Mr Sasi grows two different types of tomato plant in his greenhouse.

One week he keeps a record of the number of tomatoes he picks from each type of plant.

| Day | Mon | Tue | Wed | | Thu | Fri | Sat | Sun |
|--------|-----|-----|-----|---|-----|-----|-----|-----|
| Type A | 5 | 5 | 4 | | 1 | 0 | 1 | 5 |
| Type B | 3 | 4 | 3 | | 3 | 7 | 9 | 6 |

(a) Calculate the mean, median and mode for the Type A plants.
        Mean =3,  Median = 4,  Mode = 5.
(b) Calculate the mean, median and mode for the Type B plants.
        Mean =5,  Median = 4,  Mode = 3.
(c) Which measure would you use to argue that there is no difference between the types?
        We will use median as it is the same for both plants.
(d) Which measure would you use to argue that Type A is the best plant?

We will use mode as mode for type A is higher than B. Note that for type A mean is lower than type B and median is the same for both types.

(e) Which measure would you use to argue that Type B is the best plant?

We will use mean as mean for type A is higher than type B.

## c. Geometric Mean:

The geometric mean is a type of mean or average, which indicates the central tendency or typical value of a set of numbers. It is similar to the arithmetic mean, which is what most people think of with the word "average", except that the numbers are multiplied and then the $n^{th}$ root (where n is the count of numbers in the set) of the resulting product is taken.

Geometric mean is defined as the $n^{th}$ root of the product of N items of series. If there are two items, take the square root; if there are three items, we take the cube root; and so on. Symbolically;

$$GM = \sqrt[n]{(X_1)(X_2)\ldots\ldots(X)_n}$$

Where $X_1$, $X_2$ ….. $X_n$ refer to the various items of the series.

For instance, the geometric mean of two numbers, say 2 and 8, is just the square root of their product; that is $\sqrt[2]{2 \times 8}$ = 4. As another example, the geometric mean of three numbers 1, ½, ¼ is the cube root of their product (1/8), which is 1/2; that is $\sqrt[3]{1 \times \frac{1}{2} \times \frac{1}{4}} = \sqrt[3]{\frac{1}{8}} = \frac{1}{2}$ .

When the number of items is three or more, the task of multiplying the numbers and of extracting the root becomes excessively difficult. To simplify calculations, logarithms are used. GM then is calculated as follows.

$$\log G.M = \frac{\log X_1 + \log X_2 + \ldots\ldots \log X_N}{N}$$

$$G.M. = \frac{\sum \log X}{N}$$

$$G.M. = \text{Antilog} \left[ \frac{\sum \log X}{N} \right]$$

In discrete series GM = Antilog $\left[ \frac{\sum f \log X}{N} \right]$

In continuous series GM = Antilog $\left[ \frac{\sum f \log m}{N} \right]$

Where f = frequency
M = mid point

## Merits of G.M
1. It is based on each and every item of the series.
2. It is rigidly defined.
3. It is useful in averaging ratios and percentages and in determining rates of increase and decrease.
4. It is capable of algebraic manipulation.

## Limitations
1. It is difficult to understand
2. It is difficult to compute and to interpret

3. It can't be computed when there are negative and positive values in a series or one or more of values is zero.
4. G.M has very limited applications.

## d. Harmonic Mean:

Harmonic mean is a kind of average. It is the mean of a set of positive variables. It is calculated by dividing the number of observations by the reciprocal of each number in the series.

Harmonic Mean of a set of numbers is the number of items divided by the sum of the reciprocals of the numbers. Hence, the Harmonic Mean of a set of n numbers i.e. $a_1$, $a_2$, $a_3$, ... $a_n$, is given as

$$Harmonic\ mean = \frac{n}{a_1 + a_2 + a_3 + \ldots + a_n}$$

Example: Find the harmonic mean for the numbers 3 and 4.
Take the reciprocals of the given numbers and sum them.

$$\frac{1}{3} + \frac{1}{4} = \frac{4+3}{12} = \frac{7}{12}$$

Now apply the formula. Since the number of observations is two, here n = 2.

$$Harmonic\ mean = \frac{2}{\frac{7}{12}} = 2 \times \frac{12}{7} = \frac{24}{7} = 3.43$$

In discrete series, H.M = $\dfrac{N}{\sum\left[f.\dfrac{1}{x}\right]}$

In continuous series, H.M = $\dfrac{N}{\sum\left[f.\dfrac{1}{m}\right]}$ = $\dfrac{N}{\sum\left[\dfrac{f}{m}\right]}$

## Merits of Harmonic mean:

1. Its value is based on every item of the series.
2. It lends itself to algebraic manipulation.

## Limitations

1. It is not easily understood
2. It is difficult to compute
3. It gives larges weight to smallest item.

## 4. POSITIONAL VALUES

Statisticians often talk about the position of a value, relative to other values in a set of observations. The most common measures of position are Quartiles, deciles, and percentiles. Measures of position are techniques that divide a set of data into equal groups. Quartiles, deciles, and percentiles divide the data set into equal parts.

The data must be arranged in order to find these measures of position. To determine the measurement of position, the data must be sorted from lowest to highest.

We discuss them in detail in the next section.

## 4. MEASURES OF DISPERSION

The terms variability, spread, and dispersion are synonyms, and refer to how spread out a distribution is. Just as in the section on central tendency where we discussed measures of the centre of a distribution of scores, here we discuss measures of the variability of a distribution. Measures of variability provide information about the degree to which individual scores are clustered about or deviate from the average value in a distribution.

Quite often students find it difficult to understand what is meant by variability or dispersion and hence they find the measures of dispersion difficult. So we will discuss the meaning of the term in detail. First one should understand that dispersion or variability is a continuation of our discussion of measure of central tendency. So for any discussion on measure of dispersion we should use any of the measure of central tendency. We continue this discussion taking mean as an example. The mean or average measures the centre of the data. It is one aspect observations. Another feature of the observations is as to how the observations are spread about the centre. The observation may be close to the centre or they may be spread away from the centre. If the observations are close to the centre (usually the arithmetic mean or median), we say that dispersion or scatter or variation is small. If the observations are spread away from the centre, we say dispersion is large.

Let us make this clear with the help of an example. Suppose we have three groups of students who have obtained the following marks in a test. The arithmetic means of the three groups are also given below:

Group A: 46, 48, 50, 52, 54, for this the mean is 50.

Group B: 30, 40, 50, 60, 70, for this the mean is 50.

Group C: 40, 50, 60, 70, 80, for this the mean is 60.

In a group A and B arithmetic means are equal i.e. mean of Group A = Mean of Group B = 50. But in group A the observations are concentrated on the centre. All students of group A have almost the same level of performance. We say that there is consistence in the observations in group A. In group B the mean is 50 but the observations are not closed to the centre. One observation is as small as 30 and one observation is as large as 70. Thus there is greater dispersion in group B. In group C the mean is 60 but the spread of the observations with respect to the centre 60 is the same as the spread of the observations in group B with respect to their own centre which is 50. Thus in group B and C the means are different but their dispersion is the same. In group A and C the means are different and their dispersions are also different. Dispersion is an important feature of the observations and it is measured with the help of the measures of dispersion, scatter or variation. The word variability is also used for this idea of dispersion.

The study of dispersion is very important in statistical data. If in a certain factory there is consistence in the wages of workers, the workers will be satisfied. But if some workers have high wages and some have low wages, there will be unrest among the low paid workers and they might go on strikes and arrange demonstrations. If in a certain country some people are very poor and some are very high rich, we say there is economic disparity. It means that dispersion is large. The idea of dispersion is important in the study of wages of workers, prices of commodities, standard of living of different people, distribution of wealth, distribution of land among framers and various other fields of life. Some brief definitions of dispersion are:

The degree to which numerical data tend to spread about an average value is called the dispersion or variation of the data.

Dispersion or variation may be defined as a statistics signifying the extent of the scatter of items around a measure of central tendency.
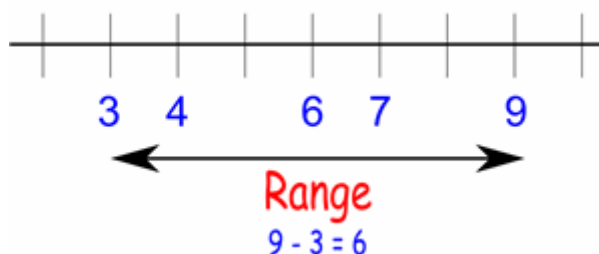
Dispersion or variation is the measurement of the scatter of the size of the items of a series about the average.

There are five frequently used measures of variability: the Range, Interquartile range or quartile deviation, Mean deviation or average deviation, Standard deviation and Lorenz curve.

**4.1 Range**

The range is the simplest measure of variability to calculate, and one you have probably encountered many times in your life. The range is simply the highest score minus the lowest score.

Range: R = maximum – minimum



Let's take a few examples. What is the range of the following group of numbers: 10, 2, 5, 6, 7, 3, 4. Well, the highest number is 10, and the lowest number is 2, so 10 - 2 = 8. The range is 8.

Let's take another example. Here's a dataset with 10 numbers: 99, 45, 23, 67, 45, 91, 82, 78, 62, 51. What is the range. The highest number is 99 and the lowest number is 23, so 99 - 23 equals 76; the range is 76.

Example 2: Ms. Kesavan listed 9 integers on the blackboard. What is the range of these integers? 14, -12, 7, 0, -5, -8, 17, -11, 19

Ordering the data from least to greatest, we get:   -12,  -11,  -8,  -5,  0,  7,  14,  17,  19

Range: R = highest - lowest = 19 - -12 = 19 + +12 = +31

The range of these integers is +31.

Example 3:  A marathon race was completed by 5 participants. What is the range of times given in hours below

2.7 hr,  8.3 hr,  3.5 hr,  5.1 hr,  4.9 hr

Ordering the data from least to greatest, we get:  2.7,  3.5,  4.9,  5.1,  8.3

Range: R = highest – lowest = 8.3 hr - 2.7 hr = 5.6 hr

The range of marathon race is 5.6 hr.

**<u>Merits and Limitations</u>**

**Merits**

➢ Amongst all the methods of studying dispersion, range is the simplest to understand easiest to compute.

➢ It takes minimum time to calculate the value of range Hence if one is interested in getting a quick rather than very accurate picture of variability one may compute range.
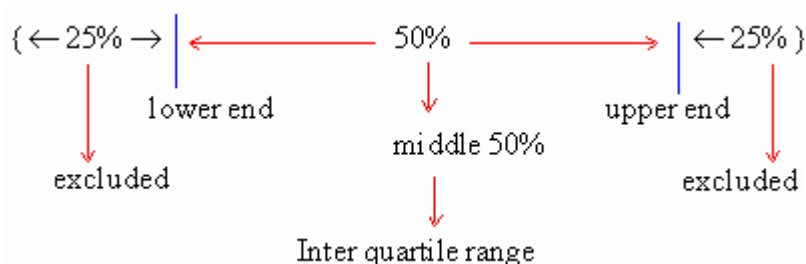
**Limitation**

➢ Range is not based on each and every item of the distribution.

➢ It is subject to fluctuation of considerable magnitude from sample to sample.

➢ Range can't tell us anything about the character of the distribution with the two.

➢ According to kind "Range is too indefinite to be used as a practical measure of dispersion

**Uses of Range**

➢ Range is useful in studying the variations in the prices of stocks, shares and other commodities that are sensitive to price changes from one period to another period.

➢ The meteorological department uses the range for weather forecasts since public is interested to know the limits within which the temperature is likely to vary on a particular day.

### 4.2 Inter – Quartile Range or Quartile Deviation

So we have seen Range which is a measure of variability which concentrates on two extreme values. If we concentrate on two extreme values as in the case of range, we do not get any idea about the scatter of the data within the range (i.e. what happens within the two extreme values). If we discard these two values the limited range thus available might be more informative. For this reason the concept of inter-quartile range is developed. It is the range which includes middle 50% of the distribution. Here 1/4 (one quarter of the lower end and 1/4 (one quarter) of the upper end of the observations are excluded.



Now the lower quartile (Q1) is the 25th percentile and the upper quartile (Q3) is the 75th percentile. It is interesting to note that the 50th percentile is the middle quartile (Q2) which is in fact what you have studied under the title' Median. Thus symbolically

Inter quartile range = Q3 - Q1
If we divide ( Q3 - Q1 ) by 2 we get what is known as Semi-inter quartile range.

$$\frac{Q_3 - Q_1}{2}$$

i.e. . It is known as Quartile deviation ( Q. D or SI QR ).

Another look at the same issue is given here to make the concept more clear for the student. In the same way that the median divides a dataset into two halves, it can be further divided into quarters by identifying the upper and lower quartiles. The lower quartile is found one quarter of the way along a dataset when the values have been arranged in order of magnitude; the upper quartile is found three quarters along the dataset. Therefore, the upper quartile lies half way between the median and the highest value in the dataset whilst the lower quartile lies halfway between the median and the lowest value in the dataset. The inter-quartile range is found by subtracting the lower quartile from the upper quartile.

For example, the examination marks for 20 students following a particular module are arranged in order of magnitude.

| Student | A | B | C | D | E | F | G | H | I | J | K | L | M | N | O | P | Q | R | S | T |
|---------|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mark | 43 | 48 | 50 | 50 | 52 | 53 | 56 | 58 | 59 | 60 | 62 | 65 | 66 | 68 | 70 | 71 | 74 | 76 | 78 | 80 |

median lies at the mid-point between the two central values (10th and 11th)

= half-way between 60 and 62 = 61

The lower quartile lies at the mid-point between the 5th and 6th values

= half-way between 52 and 53 = 52.5

The upper quartile lies at the mid-point between the 15th and 16th values

= half-way between 70 and 71 = 70.5

The inter-quartile range for this dataset is therefore 70.5 - 52.5 = 18 whereas the range is: 80 - 43 = 37.

The inter-quartile range provides a clearer picture of the overall dataset by removing/ignoring the outlying values.

Like the range however, the inter-quartile range is a measure of dispersion that is based upon only two values from the dataset. Statistically, the standard deviation is a more powerful measure of dispersion because it takes into account every value in the dataset. The standard deviation is explored in the next section.

**Example 1**

The wheat production (in Kg) of 20 acres is given as: 1120, 1240, 1320, 1040, 1080, 1200, 1440, 1360, 1680, 1730, 1785, 1342, 1960, 1880, 1755, 1720, 1600, 1470, 1750, and 1885. Find the quartile deviation and coefficient of quartile deviation.

After arranging the observations in ascending order, we get

1040, 1080, 1120, 1200, 1240, 1320, 1342, 1360, 1440, 1470, 1600, 1680, 1720, 1730, 1750, 1755, 1785, 1880, 1885, 1960.

$$Q_1 = Value\ of\left(\frac{n+1}{4}\right)th\ item$$

$$¿\ Value\ of\left(\frac{20+1}{4}\right)th\ item$$

$$¿\ Value\ of\ (5.25)\ th\ item$$

$$¿\ 5th\ item + 0.25\ (6th\ item - 5th\ item)$$

$$¿\ 1240 + 0.25\ (1320 - 1240)$$

$$Q_1 = 1240 + 20 = 1260$$

$$Q_3 = Value\ of\ \frac{3(n+1)}{4}th\ item$$

$$¿\ Value\ of\ \frac{3(20+1)}{4}th\ item$$

$$¿\ Value\ of\ (15.75)\ th\ item$$

$$¿\ 15th\ item + 0.75\ (16th\ item - 15th\ item)$$

$$¿\ 1750 + 0.75\ (1755 - 1750)$$

$$Q_3 = 1750 + 3.75 = 1753.75$$

$$Quartile\ deviation\ (Q.D.) = \frac{Q_3 - Q_1}{2} = \frac{1753.75 - 1260}{2} = \frac{492.75}{2} = 246.88$$

$$Coefficient\ of\ Quartile\ Deviation = \frac{Q_3 - Q_1}{Q_3 + Q_1} = \frac{1753.75 - 1260}{1753.75 + 1260} = 0.164$$

**Example 2**
Calculate the range and Quartile deviation of wages.

| Wages | No. of Labourers |
|-------|------------------|
| 30-32 | 12 |
| 32-34 | 18 |
| 34-36 | 16 |
| 36-38 | 14 |
| 38-40 | 12 |
| 40-42 | 8 |
| 42-44 | 6 |

| Wages (₹) | Labourers |
|-----------|-----------|
| | |

| 30 – 32 | 12 |
|---|---|
| 32 – 34 | 18 |
| 34 – 36 | 16 |
| 36 – 38 | 14 |
| 38 – 40 | 12 |
| 40 – 42 | 8 |
| 42 - 44 | 6 |

**Solution**

Range : $= L - S$

Calculation of Quartiles :

| X | f | c.f |
|---|---|---|
| 30 – 32 | 12 | 12 |
| 32 – 34 | 18 | 30 |
| 34 – 36 | 16 | 46 |
| 36 – 38 | 14 | 60 |
| 38 – 40 | 12 | 72 |
| 40 – 42 | 8 | 80 |
| 42 - 44 | 6 | 86 |

$$Q_1 = \text{Size of } \left(\frac{N}{4}\right)^{th} \text{ item}$$

$$= \frac{86}{4} = 21.5$$

ie. Q. lies in the group 32 – 34

$$Q_1 = L + \frac{\frac{N}{4} - c.f}{f} \times I \qquad = 32 + \frac{21.5 - 12}{18} \times 2$$

$$= 32 + \frac{19}{18} \qquad = 32 + 1.06$$

$$= 33.06$$
$$====$$

$$Q_3 = \text{Size of } \left(\frac{3N}{4}\right)^{th} \text{ item} \qquad = 3 \times \frac{86}{4} = 64.5^{th} \text{ item}$$

$$Q_3 \text{ lies in the group 38 – 40}$$

$$Q_3 = L + \frac{\frac{3N}{4} - c.f}{f} \times I \qquad = 38 + \frac{64.5 - 60}{12} \times 2$$

$$= 38 + 0.75 \qquad = 38.75$$

$$Q.D = \frac{Q_3 - Q_1}{2} \qquad = \frac{38.75 - 33.06}{2}$$

$$= \frac{5.69}{2} \quad = 2.85$$

$$\text{Coefficient of Q.D.} \quad = \frac{Q_3 - Q_1}{Q_3 + Q_1} \qquad = \frac{38.75 - 33.0}{38.75 + 33.06}$$

$$= \frac{5.69}{71.81} \qquad = 0.08$$

## Merits of Quartile Deviation

1. It is simple to understand and easy to calculate.
2. It is not influenced by extreme values.
3. It can be found out with open end distribution.
4. It is not affected by the presence of extreme values.

## Demerits

1. It ignores the first 25% of the items and the last 25% of the items.
2. It is a positional average: hence not amenable to further mathematical treatment.
3. The value is affected by sampling fluctuations.

## 4.3 Mean Deviation or Average Deviation

An average deviation (mean deviation) is the average amount of variations (scatter) of the items in a distribution from either the mean or the median or the mode, ignoring the signs of these deviations. In other words, the mean deviation or average deviation is the arithmetic mean of the absolute deviations.
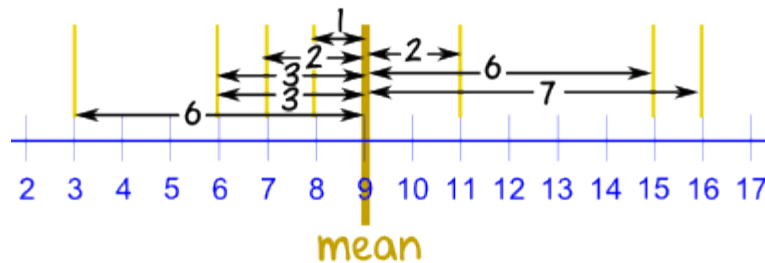
Example 1: Find the Mean Deviation of 3, 6, 6, 7, 8, 11, 15, 16

Step 1: Find the mean: $Mean = \frac{3 + 6 + 6 + 7 + 8 + 11 + 15 + 16}{8} = \frac{72}{8} = 9$

Step 2: Find the distance of each value from that mean:

| Value | Distance from 9 |
|-------|-----------------|
| 3 | 6 |
| 6 | 3 |
| 6 | 3 |
| 7 | 2 |
| 8 | 1 |
| 11 | 2 |
| 15 | 6 |
| 16 | 7 |

Which looks like this diagrammatically:

Step 3. Find the mean of those distances:

$$Mean\,Deviation = \frac{6+3+3+2+1+2+6+7}{8} = \frac{30}{8} = 3.75$$

So, the mean = 9, and the mean deviation = 3.75

It tells us how far, on average, all values are from the middle.

In that example the values are, on average, 3.75 away from the middle.

The formula is:

$$Mean\,Deviation = \frac{\sum |X - \mu|}{N}$$

Where

μ is the mean (in our example μ = 9)

x is each value (such as 3 or 16)

N is the number of values (in our example N = 8)

Each distance we calculated is called an Absolute Deviation, because it is the Absolute Value of the deviation (how far from the mean).To show "Absolute Value" we put "|" marks either side like this: |-3| = 3. Thus absolute value is one where we ignore sign. That is, if it is – or +, we consider it as +. Eg. -3 or +3 will be taken as just 3.

Let us redo example 1 using the formula: Find the Mean Deviation of 3, 6, 6, 7, 8, 11, 15, 16

Step 1: Find the mean:

$$\mu = \frac{3+6+6+7+8+11+15+16}{8} = \frac{72}{8} = 9$$

Step 2: Find the Absolute Deviations:

| x | x - μ | \|x - μ\| |
|---|---|---|
| 3 | -6 | 6 |
| 6 | -3 | 3 |
| 6 | -3 | 3 |
| 7 | -2 | 2 |
| 8 | -1 | 1 |
| 11 | 2 | 2 |

| 15 | 6 | 6 |
|----|---|---|
| 16 | 7 | 7 |
| | $\sum |x-\mu|=30$ | |

Step 3. Find the Mean Deviation:

$$MeanDeviation=\frac{\sum |X-\mu|}{N}=\frac{30}{8}=3.75$$

Example 2

Calculate the mean deviation using mean for the following data

| 2-4 | 4-6 | 6-8 | 8-10 |
|-----|-----|-----|------|
| 3 | 4 | 2 | 1 |

Solution

| Class | Mid Value (X) | Frequency (f) | d = X-5 | fd | $|X-\acute{X}|=|X-5.2|$ | $f|X-\acute{X}|$ |
|-------|---------------|---------------|---------|-----|------------------------|------------------|
| 2-4 | 3 | 3 | -2 | -6 | 2.2 | 6.6 |
| 4-6 | 5 | 4 | 0 | 0 | 0.2 | 0.8 |
| 6-8 | 7 | 2 | 2 | 4 | 1.8 | 3.6 |
| 8-10 | 9 | 1 | 4 | 4 | 3.8 | 3.8 |
| | | $\sum f=10$ | | $\sum fd=2$ | | $\sum f|X-\acute{X}|=1$ |

$$\acute{X}=A+\frac{\sum fd}{N}=5+\frac{2}{10}=5.2$$

$$Mean\, deviation=\frac{1}{N}\sum f|X-\acute{X}|=\frac{14.8}{10}=1.48$$

Example 3

Calculate mean deviation based on (a) Mean and (b) median

| Class Interval | 0-10 | 10-20 | 20-30 | 30-40 | 40-50 | 50-60 | 60-70 |
|----------------|------|-------|-------|-------|-------|-------|-------|
| Frequenc y f | 8 | 12 | 10 | 8 | 3 | 2 | 7 |

Solution

Let us first make the necessary computations.

| Class interval | Mid value (X) | Freq-uency (f) | Less than c.f. | fX | $\|X-\acute{X}\|$ | $f\|X-\acute{X}\|$ | $\|X-Md\|$ | $f\|X-M\|$ |
|---|---|---|---|---|---|---|---|---|
| 0-10 | 5 | 8 | 8 | 40 | 24 | 192 | 17 | 136 |
| 10-20 | 15 | 12 | 20 | 180 | 14 | 168 | 7 | 84 |
| 20-30 | 25 | 10 | 30 | 250 | 4 | 40 | 3 | 30 |
| 30-40 | 35 | 8 | 38 | 280 | 6 | 48 | 13 | 104 |
| 40-50 | 45 | 3 | 41 | 135 | 16 | 48 | 23 | 69 |
| 50-60 | 55 | 2 | 43 | 110 | 26 | 52 | 33 | 66 |
| 60-70 | 65 | 7 | 50 | 455 | 36 | 252 | 43 | 301 |
| | | N=50 | | $\sum fx=$ | | $\sum f\|X$ | | $\sum f\|X$ |

(a) M.D. from Mean

$$Mean(\acute{X})=\frac{1}{N}\sum fX=\frac{1450}{50}=29$$

So mean =29. Let us now find men deviation about mean

$$M.D.=\frac{1}{N}\sum f|X-\acute{X}|=\frac{800}{50}=16$$

We see that mean deviation based on mean is 16.

Now let us compute M.D. about median

(b) M.D. from median

(N/2) =(50/2) = 25. The c.f. just greater than 25 is 30 in the table above. So the corresponding class 20-30 is the median class.

So $l$ = lower limit of the median class = 20, $f$ = *frequency of the median class* = 25, $h$ = class interval of the median class =10, $c$ = cumulative frequency of the preceding median class =20.

Use the formula of median to substitute values.

$$Median=l+\frac{h}{f}\left(\frac{N}{2}-C\right)$$

$$¿20+\frac{10}{25}(25-20)=20+2=22$$

Median = 22. Let us now find Mean Deviation about median.

$$M.D.=\frac{1}{N}\sum f|X-Md|=\frac{790}{50}=15.8$$

Thus we have computed Mean Deviation from Mean and Median. Let us compare the two results. MD from Mean is 16 and MD from median is 15.8.

So, M.D. from Median < M.D. from Mean. This implies that M.D. is least when taken about median.

**Merits of M.D.**

i. It is simple to understand and easy to compute.

ii. It is not much affected by the fluctuations of sampling.

iii. It is based on all items of the series and gives weight according to their size.

iv. It is less affected by extreme items.

v. It is rigidly defined.

vi. It is a better measure for comparison.

**Demerits of M.D.**
i. It is a non-algebraic treatment
ii. Algebraic positive and negative signs are ignored. It is mathematically unsound and illogical.
iii. It is not as popular as standard deviation.

Uses :

It will help to understand the standard deviation. It is useful in marketing problems. It is used in statistical analysis of economic, business and social phenomena. It is useful in calculating the distribution of wealth in a community or nation.

## 4.4 Measures of Position (Positional values / Partition Values): Quartiles, Deciles and Percentiles

Statisticians often talk about the position of a value, relative to other values in a set of observations. A measure of position is a method by which the position that a particular data value has within a given data set can be identified. The most common measures of position are quartiles, deciles and percentiles.

**Quartiles**

The mean and median both describe the 'centre' of a distribution. This is usually what you want to summarize about a set of marks, but occasionally a different part of the distribution is of more interest.

The median of a distribution splits the data into two equally-sized groups. In the same way, the quartiles are the three values that split a data set into four equal parts. Note that the 'middle' quartile is the median.

The upper quartile describes a 'typical' mark for the top half of a class and the lower quartile is a 'typical' mark for the bottom half of the class.

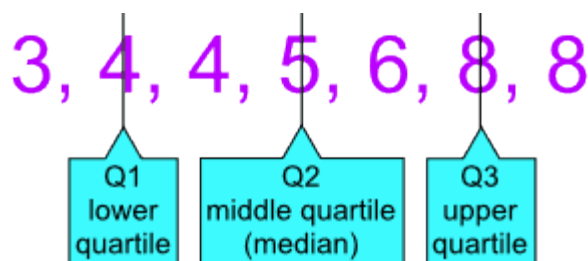Thus Quartiles are the values that divide a list of numbers into quarters.

If you are given a set of data this is how you find the quartile

- First put the list of numbers in order
- Then cut the list into four equal parts
- The Quartiles are at the "cuts"

**Example: 5, 8, 4, 4, 6, 3, 8**
Put them in order: 3, 4, 4, 5, 6, 8, 8
Cut the list into quarters:

3, 4, 4, 5, 6, 8, 8

Q1 lower quartile | Q2 middle quartile (median) | Q3 upper quartile

And the result is:

Quartile 1 (Q1) = 4

Quartile 2 (Q2), which is also the Median, = 5

Quartile 3 (Q3) = 8

**Deciles**

In a similar way, the deciles of a distribution are the ninevalues that split the data set into tenequal parts. It is called a decile. A decile is any of the nine values that divide the sorted data into ten equal parts, so that each part represents 1/10 of the sample or population.

Deciles are similar to quartiles. But while quartiles sort data into four quarters, deciles sort data into ten equal parts: The 10th, 20th, 30th, 40th, 50th, 60th, 70th, 80th, 90th and 100th percentiles.

**Deciles** (sounds like decimal and percentile together), which splits the data into 10% groups:

- The **1st decile** is the **10th percentile** (the value that divides the data so that 10% is below it)
- The **2nd decile** is the **20th percentile** (the value that divides the data so that 20% is below it)
- etc!

**Example: You are the fourth tallest person in a group of 20**

80% of people are shorter than you:



You are at the **8th decile** (the 80th percentile).

**Percentiles**

The percentiles divide the data into 100 equal regions.

A percentile (or a centile) is a measure used in statistics indicating the value below which a given percentage of observations in a group of observations fall. For example, the 20th percentile is the value (or score) below which 20 percent of the observations may be found.

Percentiles report the relative standing of a particular value within a statistical data set. For example, in the case of exam scores, assume in a tough exam you scored 40 points out of 100. In this case, your score itself is meaningless, but your percentile tells you everything. Suppose your exam score is better than 90% of the rest of the class. That means your exam score is at the 90th percentile (so $k = 90$). On the other hand, if your score is at the 10th percentile, then $k = 10$; that means only 10% of the other scores are below yours, and 90% of them are above yours. So percentile tells you where you stand in relation to other students in the class.

Example: You are the fourth tallest person in a group of 20

80% of people are shorter than you:

That means you are at the 80th percentile (8th decile).

If your height is 1.85m then "1.85m" is the 80th percentile height in that group.

A useful property of percentiles is they have a universal interpretation: Being at the 95th percentile means the same thing no matter if you are looking at exam scores or weights of students in a class etc; the 95th percentile always means 95% of the other values lie below yours, and 5% lie above it.

Please note that a percentile different from a percent. As we have seen a percentile is a value in the data set that marks a certain percentage of the way through the data. Suppose your score on the CAT exam was reported to be the 70th percentile. This does not mean that you scored 70% of the questions correctly. It means that 70% of the students' scores were lower than yours and 30% of the students' scores were higher than yours.
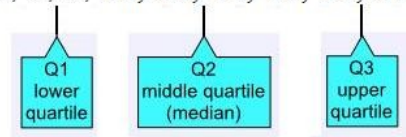
**Exercises**

(A) Individual series

**1.** What are the quartiles for the following set of numbers?

   8, 11, 20, 10, 2, 17, 15, 5, 16, 15, 25, 6

   First arrange the numbers in order: 2, 5, 6, 8, 10, 11, 15, 15, 16, 17, 20, 25

This list can be split up into four equal groups of three:



Therefore:

Q1 is the mean of 6 and 8 = (6 + 8) ÷ 2 = 7

Q2 is the mean of 11 and 15 = (11 + 15) ÷ 2 = 13

Q3 is the mean of 16 and 17 = (16 + 17) ÷ 2 = 16.5

**Formula for finding decile**

There are nine deciles namely $D_1, D_2, D_3, \ldots, D_9$.

<u>Deciles for **Individual** Observations</u> (Ungrouped Data):

First arrange the items in ascending order of magnitude, then apply the formula. Here 'n' stands for number of values.

$$k^{th} decile = D_k = Value\ of\ k\left(\frac{n+1}{10}\right)^{th} item$$

$$1st\ Decile = D_1 = Value\ of\left(\frac{n+1}{10}\right)^{th} item$$

$$2nd\ Decile = D_2 = Value\ of\ 2\left(\frac{n+1}{10}\right)^{th} item$$

$$3rd\ Decile = D_3 = Value\ of\ 3\left(\frac{n+1}{10}\right)^{th} item$$

Like this for $D_9$ we have

$$9th\,Decile=D_9=Value\,of\,9\left(\frac{n+1}{10}\right)^{th}item$$

Deciles for Grouped Frequency Distribution:

$$1st\,Decile=D_1=l+\frac{h}{f}\left(\frac{1n}{10}-c\right)$$

Where
$l$ is the lower class boundary of the class containing the 1st decile
$h$ is the class interval of the class containing $D_1$
$f$ is the frequency of the class containing $D_1$
c is the cumulative frequency of the class immediately preceding to the class containing $D_1$
$n=\sum f$, that is, $n$ is the is the total number of frequencies

Applying the above formula to find 2nd Decile $D_2$

$$2nd\,Decile=D_2=l+\frac{h}{f}\left(\frac{2n}{10}-c\right)$$

$$3rd\,Decile=D_3=l+\frac{h}{f}\left(\frac{3n}{10}-c\right)$$

$$9th\,Decile=D_9=l+\frac{h}{f}\left(\frac{9n}{10}-c\right)$$

Formula for finding Percentiles:
There are ninety nine percentiles namely $P_1, P_2, P_3, \ldots, P_{99}$.
Percentile for **Individual** Observations (Ungrouped Data):
First arrange the items in ascending order of magnitude, then apply the formula. Here 'n' stands for number of values.

$$k^{th}\,percentile=P_k=Value\,of\,k\left(\frac{n+1}{100}\right)^{th}item$$

$$1st\,Percentile=P_1=Value\,of\left(\frac{n+1}{100}\right)^{th}item$$

$$2nd\,Percentile=P_2=Value\,of\,2\left(\frac{n+1}{100}\right)^{th}item$$

$$99th\,Percentile=P_2=Value\,of\,99\left(\frac{n+1}{100}\right)^{th}item$$

Percentile for Grouped Frequency Distribution:
The percentiles are usually calculated for grouped data.

$$1st\,Percentile=P_1=l+\frac{h}{f}\left(\frac{1n}{100}-c\right)$$

Where
$l$ is the lower class boundary of the class containing the 1st percentile
$h$ is the class interval of the class containing $P_1$
$f$ is the frequency of the class containing $P_1$
c is the cumulative frequency of the class immediately preceding to the class containing $P_1$
$n=\sum f$, that is, $n$ is the is the total number of frequencies

Note that 50th percentile is the median by definition as half of the values in the data are smaller than the median and half of the values are larger than the median.

Applying the above formula to find $2^{nd}$ Percentile $P_2$

$$2nd\,Percentile = P_2 = l + \frac{h}{f}\left(\frac{2n}{100} - c\right)$$

$$3rd\,Percentile = P_3 = l + \frac{h}{f}\left(\frac{3n}{100} - c\right)$$

$$99th\,Percentile = P_{99} = l + \frac{h}{f}\left(\frac{99n}{100} - c\right)$$

Example 1: Find $4^{th}$ Decile, $3^{rd}$ Decile and $30^{th}$ percentile for the following observations.
65, 23, 95, 101, 89, 52, 43, 15, 55
First arrange the items in ascending order of magnitude.
15, 23, 43, 52, 55, 65, 89, 95, 101
Here 'n' (number of values) here is 9.
Now apply the formula.

$$4th\,Decile = D_4 = Value\,of\,4\left(\frac{n+1}{10}\right)^{th}item$$

$$¿\,Value\,of\,4\left(\frac{9+1}{10}\right)^{th}item$$

$$¿¿4th\,item = 52$$

$$3rd\,Decile = D_3 = Value\,of\,3\left(\frac{n+1}{10}\right)^{th}item$$

$$¿\,Value\,of\,3\left(\frac{9+1}{10}\right)^{th}item$$

$$¿¿3rd\,item = 43$$

$$30th\,Percentile = P_{30} = Value\,of\,30\left(\frac{n+1}{100}\right)^{th}item$$

$$¿\,Value\,of\,30\left(\frac{9+1}{100}\right)^{th}item$$

$$¿\,Value\,of\,30\left(\frac{10}{100}\right)^{th}item$$

$$¿¿3^{rd}\,item = 43$$

Note that $3^{rd}$ Decile is the same as $30^{th}$ Percentile.

Example 2: Find $3^{rd}$ Decile, and $80^{th}$ Percentile for the following observations.
18, 50, 15, 30, 13, 28, 18, 90, 51, 47
First arrange the items in ascending order of magnitude.
13, 15, 18, 18, 28, 30, 47, 50, 51, 90
Here 'n' (number of values) here is 10.
Now apply the formula.

$$3rd\,Decile = D_3 = Value\,of\,3\left(\frac{n+1}{10}\right)^{th}item$$

$¿ Value of \left(\dfrac{10+1}{10} \times 3\right)^{th} item$

$¿ Value of 1.1 \times 3 rd item = 3.3 rd items$

$¿ Value of 3 rd item + 0.3(4 th - 3 rd) = 18 + 0.3(18 - 18)$

$¿$

$¿ 18 + 0.3 ¿$ 0) = 18 + 0 = 18

$80 th Percentile = P_{80} = Value of 80\left(\dfrac{n+1}{100}\right)^{th} item$

$¿ Value of 80\left(\dfrac{10+1}{100}\right)^{th} item$

$¿ Value of 80\left(\dfrac{11}{100}\right)^{th} item$

$¿ Value of 80(0.11)^{th} item$

$¿ Value of 8.8^{th} item$

$$= 8^{th} item + .8(9^{th} item - 8^{th} item)$$
$$= 50 + .8(51 - 50)$$
$$= 50 + .8(1) = 50.8$$

Example 3: For the following grouped data compute $P_{10}$, $P_{25}$, $P_{50}$, and $P_{95}$. Also find $D_1$ and $D_7$.

| Class Boundaries | $X_i$ | $f_i$ | CF |
|---|---|---|---|
| 85.5-90.5 | 87 | 6 | 6 |
| 90.5-95.5 | 93 | 4 | 10 |
| 95.5-100.5 | 98 | 10 | 20 |
| 100.5-105.5 | 103 | 6 | 26 |
| 105.5-110.5 | 108 | 3 | 29 |
| 110.5-115.5 | 113 | 1 | 30 |
| | | 30 | |

Finding $P_{10}$

Locate the 10th Percentile by $\dfrac{10 \times n}{100} = \dfrac{10 \times 30}{100} = 3$ . So the third observation is $P_{10}$.

So, $P_{10}$ group is the one containing the 3rd observation in the CF column. Here it is 85.5–90.5. Now apply the formula

$10 th Percentile = P_{10} = l + \dfrac{h}{f}\left(\dfrac{10 n}{100} - c\right)$

$¿ 85.5 + \dfrac{5}{6}\left(\dfrac{10 \times 30}{100} - 0\right)$

$¿ 85.5 + \dfrac{5}{6}\left(\dfrac{300}{100} - 0\right)$

$¿ 85.5 + \dfrac{5}{6}(3 - 0)$

$¿85.5+2.5=88$

### Finding P$_{25}$

Locate the 25th Percentile by $\dfrac{25 \times n}{100}=\dfrac{25 \times 30}{100}=7.5$.

So the 7.5$^{th}$ observation is P$_{25}$.

So, P$_{25}$ group is the one containing the 7.5$^{th}$ observation in the CF column. Here it is 90.5–95.5.

Now apply the formula

$$25th\,Percentile=P_{25}=l+\frac{h}{f}\left(\frac{25n}{100}-c\right)$$

$¿90.5+\dfrac{5}{4}\left(\dfrac{25 \times 30}{100}-6\right)$

$¿90.5+\dfrac{5}{4}\left(\dfrac{750}{100}-6\right)$

$¿90.5+\dfrac{5}{4}(7.5-6)$

$¿90.5+\dfrac{5}{4}1.5$

$¿90.5+(1.25 \times 1.5)=90.5+1.875=92.375$

### Finding P$_{50}$

Locate the 50th Percentile by $\dfrac{50 \times n}{100}=\dfrac{50 \times 30}{100}=15$.

So the 15$^{th}$ observation is P$_{50}$.

So, P$_{50}$ group is the one containing the 15$^{th}$ observation in the CF column. Here it is 95.5–100.5.

Now apply the formula

$$50th\,Percentile=P_{50}=l+\frac{h}{f}\left(\frac{50n}{100}-c\right)$$

$¿95.5+\dfrac{5}{10}\left(\dfrac{50 \times 30}{100}-10\right)$

$¿95.5+\dfrac{5}{10}\left(\dfrac{150}{100}-10\right)$

$¿95.5+\dfrac{5}{10}(15-10)$

$¿95.5+\dfrac{5}{10}5$

$¿95.5+(0.5 \times 5)=95.5+2.5=98$

### Finding P$_{95}$

Locate the 95th Percentile by $\dfrac{95 \times n}{100}=\dfrac{95 \times 30}{100}=\dfrac{2850}{100}=28.5$.

So the 28.5$^{th}$ observation is P$_{95}$.

So, $P_{95}$ group is the one containing the $28.5^{th}$ observation in the CF column. Here it is 105.5–110.5.

Now apply the formula

$$95th\ Percentile = P_{95} = l + \frac{h}{f}\left(\frac{95n}{100} - c\right)$$

$$¿105.5 + \frac{5}{3}\left(\frac{95 \times 30}{100} - 26\right)$$

$$¿105.5 + \frac{5}{3}(28.5 - 26)$$

$$¿105.5 + \frac{5}{3}2.5$$

$$¿105.5 + 4.1667$$

$$¿109.6667$$

### Finding $D_1$

Locate the $1^{st}$ Decile by $\quad \frac{m \times n}{10} = \frac{1 \times 30}{10} = \frac{30}{10} = 3.$

So the $3^{rd}$ observation is $D_1$.

So, $D_1$ group is the one containing the 3rd observation in the CF column. Here 3rd observation lie in first class (first group), that is 85.5–90.5.

Now apply the formula

$$1st\ Decile = D_1 = l + \frac{h}{f}\left(\frac{1n}{10} - c\right)$$

$$¿85.5 + \frac{5}{6}\left(\frac{1 \times 30}{10} - 0\right)$$

$$¿85.5 + \frac{5}{6}\left(\frac{30}{10} - 0\right)$$

$$¿85.5 + 0.833(3)$$

$$¿85.5 + 2.5$$

$$¿88$$

### Finding $D_7$

Locate the $7^{th}$ decile by $\quad \frac{m \times n}{10} = \frac{7 \times 30}{10} = \frac{210}{10} = 21.$

So the $21^{st}$ observation is $D_7$.

So, $D_7$ group is the one containing the $21^{st}$ observation in the CF column. Here 3rd observation lie in first class (first group), that is 100.5–105.5.

Now apply the formula

$$7th\ Decile = D_7 = l + \frac{h}{f}\left(\frac{1n}{10} - c\right)$$

¿ $100.5 + \dfrac{5}{6}\left(\dfrac{7 \times 30}{10} - 20\right)$

¿ $100.5 + \dfrac{5}{6}\left(\dfrac{210}{10} - 20\right)$

¿ $100.5 + 0.833(1)$

¿ $100.5 + 0.833$

¿ $101.333$

The Percentiles may be read directly from the graphs of cumulative frequency function. We can estimate Percentiles from a line graph as shown below.
Example:
A total of 10,000 people visited the shopping mall over 12 hours:

| Time (hours) | People |
|---|---|
| 0 | 0 |
| 2 | 350 |
| 4 | 1100 |
| 6 | 2400 |
| 8 | 6500 |
| 10 | 8850 |
| 12 | 10,000 |

a) Estimate the 30th Percentile (when 30% of the visitors had arrived).
b) Estimate what Percentile of visitors had arrived after 11 hours.
Solution
First draw a line graph of the data: plot the points and join them with a smooth curve:



a) The 30th Percentile occurs when the visits reach 3,000.
Draw a line horizontally across from 3,000 until you hit the curve, then draw a line vertically downwards to read off the time on the horizontal axis:

So the 30th Percentile occurs after about 6.5 hours.

b) To estimate the Percentile of visits after 11 hours: draw a line vertically up from 11 until you hit the curve, then draw a line horizontally across to read off the population on the horizontal axis:



So the visits at 11 hours were about 9,500, which is the 95th percentile.
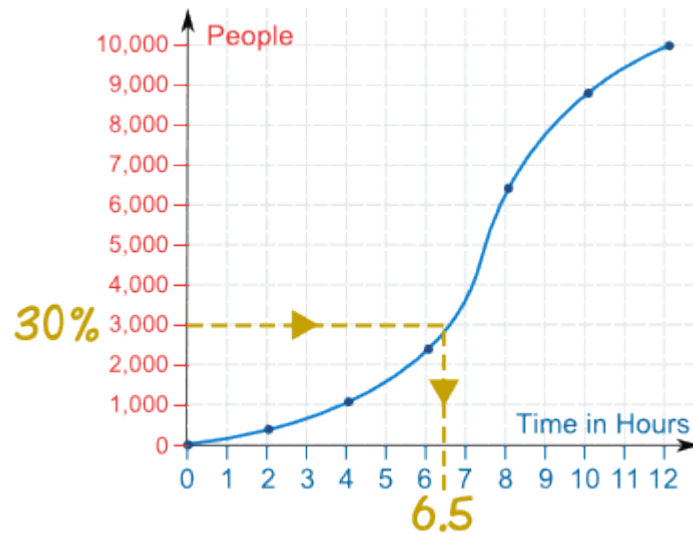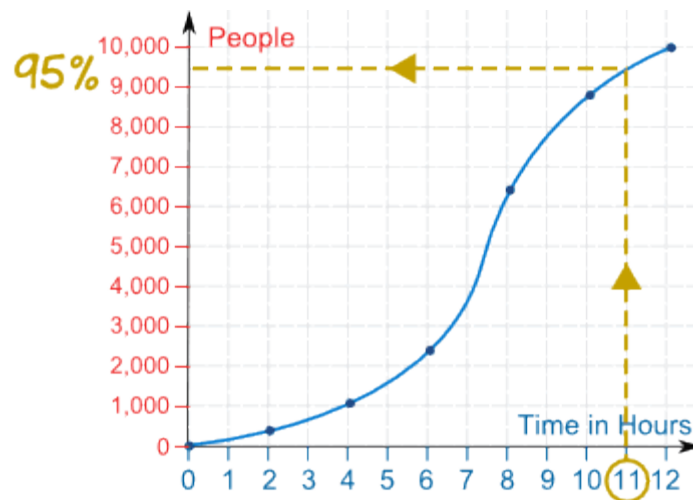
## 4.5    Standard Deviation

The concept, standard deviation was introduced by Karl Pearson in 1893.  It is the most important measure of dispersion and is widely used. It is a measure of the dispersion of a set of data from its mean. The standard deviation is kind of the "mean of the mean," and often can help you find the story behind the data.

The standard deviation is a measure that summarizes the amount by which every value within a dataset varies from the mean. Effectively it indicates how tightly the values in the dataset are bunched around the mean value. It is the most robust and widely used measure of dispersion since, unlike the range and inter-quartile range; it takes into account every variable in the dataset. When the values in a dataset are pretty tightly bunched together the standard deviation is small. When the values are spread apart the standard deviation will be relatively large.

Standard deviation is defined as a statistical measure of dispersion in the value of an asset around mean. The standard deviation calculation tells you how spread out the numbers are in your sample. Standard Deviation is represented using the symbol $\sigma\,(the\,greek\,letter\,sigma)$ .

For example if you want to measure the performance a mutual fund, SD can be used.  It gives an idea of how volatile a fund's performance is likely to be. It is an important measure of a fund's performance. It gives an idea of how much the return on the asset at a given time differs or deviates from the average return. Generally, it gives an idea of a fund's volatility i.e. a higher dispersion (indicated by a higher standard deviation) shows that the value of the asset has fluctuated over a wide range.

The formula for finding SD in a sentence form is : it is the square root of the Variance. So now you ask, 'What is the Variance'. Let us see what is variance.

The Variance is defined as:The average of the squared differences from the Mean.

We can calculate the variance follow these steps:

a. Work out the Mean (the simple average of the numbers)

b. Then for each number: subtract the Mean and square the result (the squared difference).

c. Then work out the average of those squared differences.

You may ask Why square the differences. If we just added up the differences from the mean ... the negatives would cancel the positives as shown below. So we take the square.



Example

You have figures of the marks obtained by your five bench mates which is as follows: 600, 470, 170, 430 and 300. Find out the Mean, the Variance, and the Standard Deviation.

Your first step is to find the Mean:

$$Mean=\frac{600+470+170+430+300}{5}=\frac{1970}{5}=394$$

So the mean (average) mark is 394. Let us plot this on the chart:

| x | $X-\acute{X}$ | $(X-\acute{X})^2$ |
|---|---|---|
| 600 | 206 | 42436 |
| 470 | 76 | 5776 |
| 170 | -224 | 50176 |
| 430 | 36 | 1296 |
| 300 | -94 | 8836 |
| | | $\sum(X-\acute{X})^2=108520$ |

To calculate the Variance, take each difference, square it, find the sum (108520) and find average:

$$Variance = \frac{108520}{5} = 21704$$

So, the Variance is 21,704.

The Standard Deviation is just the square root of Variance, so:

$$SD = \sigma = \sqrt{21704} = 147.32 \approx 147$$

Now we can see which heights are within one Standard Deviation (147) of the Mean. Please note that there is a slight difference when we find variance from a population and mean. In the above example we found out variance for data collected from all your bench mates. So it may be considered as population. Suppose now you collect data only from some of your bench mates. Now it may be considered as a sample. If you are finding variance for a sample data, in the formula to find variance, divide by N-1 instead of N.

For example, if we say that in our problem the marks are of some students in a class, it should be treated as a sample. In that case

Variance (or to be precise Sample Variance) = 108,520 / 4 = 27,130. Note that instead of N (i.e.5) we divided by N-1 (5-1=4).

Standard Deviation (Sample Standard Deviation) = $\sigma = \sqrt{27130} = 164.31 \approx 164$

Based on the above information, let us build the formula for finding SD. Since we use two different formulae for data which is population and data which is sample, we will have two different formula for SD also.

The "**Population** Standard Deviation":  $\sigma = \sqrt{\dfrac{1}{N}\sum_{i=1}^{N}(x_i - \mu)^2}$

$$s = \sqrt{\dfrac{1}{N-1}\sum_{i=1}^{N}(x_i - \overline{x})^2}$$

The "**Sample** Standard Deviation":

Computation of Standard Deviation: There are different methods to computeSD. They are illustrated through examples below.

Example 1

Calculate SD for the following observations using different methods.

160, 160, 161, 162, 163, 163, 163, 164, 164, 170

(a) Direct method No.1

Formula $\sigma = \sqrt{\dfrac{\sum d^2}{N}}$ where $d = x - \acute{x}$

| X | $d = x - \acute{x}$ | $d^2$ |
|---|---|---|
| 160 | -3 | 9 |
| 160 | -3 | 9 |
| 161 | -2 | 4 |
| 162 | -1 | 1 |
| 163 | 0 | 0 |
| 163 | 0 | 0 |
| 163 | 0 | 0 |
| 164 | 1 | 1 |
| 164 | 1 | 1 |
| 170 | 7 | 49 |
| $\sum X = 1630$ | | $\sum d^2 = 74$ |

$Where\ Mean = \acute{X} = \dfrac{\sum X}{N} = 163$

Now compute SD $\quad \sigma = \sqrt{\dfrac{\sum d^2}{N}}$

$\sigma = \sqrt{\dfrac{74}{10}} = \sqrt{7.4}\quad = 2.72$

(b) Direct method No.2

Here the formula is

$\sigma = \sqrt{\dfrac{\sum X^2 - \sum X^2 / N}{N}}$

| X | $X^2$ |
|---|---|
| 160 | 25600 |
| 160 | 25600 |
| 161 | 25921 |
| 162 | 26244 |
| 163 | 26569 |
| 163 | 26569 |
| 163 | 26569 |
| 164 | 26896 |
| 164 | 26896 |
| 170 | 28900 |
| $\sum X = 1630$ | $\sum X^2 = 2657640$ |

$$\sigma = \sqrt{\frac{265764 - 1630^2/10}{10}}$$

$\sigma = \sqrt{\frac{74}{10}} = \sqrt{7.4} = 2.72$    (c)Method 3 (Short Cut Method) – in this method instead of finding the

mean we assume a figure as mean. Here we have assumed 162 as mean arbitrarily.

We use the formula

$$\sigma = \sqrt{\frac{\sum dx^2}{N} - \left(\frac{\sum dx}{N}\right)^2}$$

| X | Deviation from assumed mean (here we assume mean as162) dx | $(dx)^2$ |
|---|---|---|
| 160 | -2 | 4 |
| 160 | -2 | 4 |
| 161 | -1 | 1 |
| 162 | 0 | 0 |
| 163 | 1 | 1 |
| 163 | 1 | 1 |
| 163 | 1 | 1 |
| 164 | 2 | 4 |
| 164 | 2 | 4 |
| 170 | 8 | 64 |
| 1630 | +10 | $\sum dx^2 = 84$ |

$$\sigma = \sqrt{\frac{84}{10} - \left(\frac{10}{10}\right)^2}$$

$¿\sqrt{8.4 - 1}$

$¿\sqrt{7.4} = 2.72$

Another example where we find many of the concepts together.

Example:

Given the series: 3, 5, 2, 7, 6, 4, 9.

Calculate:

The (a)mode, (b)median and (c)mean.

(d) variance (e)standard deviation and (f)The average deviation.

(a)*Mode*: **Does not exist because all the scores have the same frequency.**

(b) *Median*

2, 3, 4, 5, 6, 7, 9.

Median = 5

**(c)Mean**

$$\acute{x}=\frac{2+3+4+5+6+7+9}{7}=5.143$$

**(d)Variance**

$$Variance=\sigma^2=\frac{2^2+3^2+4^2+5^2+6^2+7^2+9^2}{7}-5.143^2=4.978$$

**(e)Standard Deviation**

$$\sigma=\sqrt{4.978}=2.231$$

(f) Average Deviation

| x | $\lvert x-\acute{x}\rvert=\lvert x-5.143\rvert$ |
|---|---|
| 2 | 3.143 |
| 3 | 2.143 |
| 4 | 1.143 |
| 5 | 0.143 |
| 6 | 0.857 |
| 7 | 1.857 |
| 9 | 3.857 |
| | $\sum\lvert x-\acute{x}\rvert=13.14$ |

$$Average\,Deviation=\frac{\sum\lvert x-\acute{x}\rvert}{N}=\frac{13.143}{7}=1.878$$

Calculation of SD for continuous series

The step deviation method is easy to use to find SD for continuous series.

$$\sigma=\sqrt{\frac{\sum f d^2}{N}-\left(\frac{\sum fd}{N}\right)^2}\times i$$

$$where\,d=\frac{(m-A)}{i}\,where\,m\,is\,midpoint \wedge i=class\,interval$$

Calculate Mean and SD for the following data

| 0-10 | 10-20 | 20-30 | 30-40 | 40-50 | 50-60 | 60-70 |
|---|---|---|---|---|---|---|
| 5 | 12 | 30 | 45 | 50 | 37 | 21 |

Make the necessary computations

| x | Midpoint (m) | $f$ | $d=\frac{(m-35)}{10}$ | $fd$ | f×d$^2$ |
|---|---|---|---|---|---|
| 0-10 | 5 | 5 | -3 | -15 | 45 |

| 10-20 | 15 | 12 | -2 | -24 | 48 |
|---|---|---|---|---|---|
| 20-30 | 25 | 30 | -1 | -30 | 30 |
| 30-40 | 35 | 45 | 0 | 0 | 0 |
| 40-50 | 45 | 50 | 1 | 50 | 50 |
| 50-60 | 55 | 37 | 2 | 74 | 148 |
| 60-70 | 65 | 21 | 3 | 63 | 189 |
| | | N = 200 | | $\sum fd = 118$ | $\sum f d^2 = 510$ |

$$Mean = \acute{X} = A + \frac{\sum fd}{N} \times i = 35 + \frac{118}{200} \times 10 = 35 + 5.9 = 40.9$$

$$\sigma = \sqrt{\frac{\sum f d^2}{N} - \left(\frac{\sum fd}{N}\right)^2} \times i = \sqrt{\frac{510}{200} - \left(\frac{118}{200}\right)^2} \times 10$$

$$¿\sqrt{2.55 - 3481} \times 10$$

$$= 1.4839 \times 10 = 14.839.$$

## Merits of Standard Deviation

1. It is rigidly defined and its value is always definite and based on all observation.

2. As it is based on arithmetic mean, it has all the merits of arithmetic mean.

3. It is possible for further algebraic treatment.

4. It is less affected by sampling fluctuations.

## Demerits

1. It is not easy to calculate.

It gives more weight to extreme values, because the values are squared up.

## 4.5 Coefficient of Variation

Standard deviation is the absolute measure of dispersion. It is expressed in terms of the units in which the original figures are collected and stated. The relative measure of standard deviation is known as coefficient of variation.

Variance : Square of Standard deviation

Symbolically;

$$Variance = \sigma^2$$

$$\sigma = \sqrt{Variance}$$

Coefficient of standard deviation = $\frac{\sigma}{\acute{X}}$

## 5. MEASURES OF VARIABILITY IN SHAPE
### - Graphic Method of Dispersion

Dispersion or variance can be represented using graphs also. We discuss here some of the graphical methods which rely on the shape of the curve to represent the deviations. We will see Lorenz Curve, Gini's Coefficient, Skewness and Kurtosis
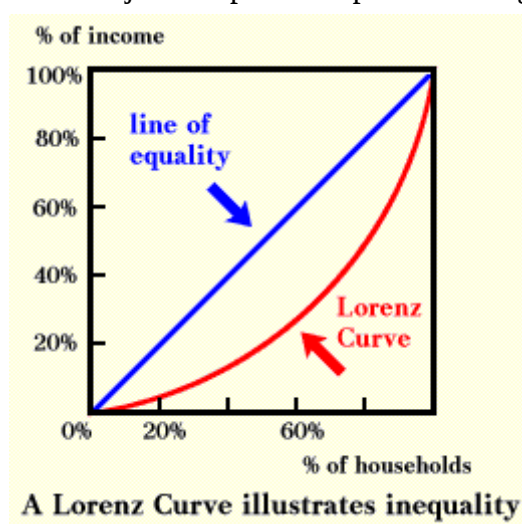
**5.1 - LORENZ CURVE**

Lorenz Curve is a graphical representation of wealth distribution developed by American economist Dr. Max O. Lorenz a popular Economic- Statistician in 1905. He studied distribution of Wealth and Income with its help.. On the graph, a straight diagonal line represents perfect equality of wealth distribution; the Lorenz curve lies beneath it, showing the reality of wealth distribution. The difference between the straight line and the curved line is the amount of inequality of wealth distribution, a figure described by the Gini coefficient. One practical use of The Lorenz curve is that it can be used to show what percentage of a nation's residents possess what percentage of that nation's wealth. For example, it might show that the country's poorest 10% possess 2% of the country's wealth.

It is graphic method to study dispersion. It helps in studying the variability in different components of distribution especially economic. The base of Lorenz Curve is that we take cumulative percentages along X and Y axis. Joining these points we get the Lorenz Curve. Lorenz Curve is of much importance in the comparison of two series graphically. It gives us a clear cut visual view of the series to be compared.

Steps to plot 'Lorenz Curve'

- Cumulate both values and their corresponding frequencies.
- Find the percentage of each of the cumulated figures taking the grand total of each corresponding column as 100.
- Represent the percentage of the cumulated frequencies on X axis and those of the values on the Y axis.
- Draw a diagonal line designated as the line of equal distribution.
- Plot the percentages of cumulated values against the percentages of the cumulated frequencies of a given distribution and join the points so plotted through a free hand curve.



**A Lorenz Curve illustrates inequality**

The greater the distance between the curve and the line of equal distribution, the greater the dispersion.  If the Lorenz curve is nearer to the line of equal distribution, the dispersion or variation is smaller.

Based on data of annual income of 8 individuals we have drawn a Lorenz curve below using MS Excel.

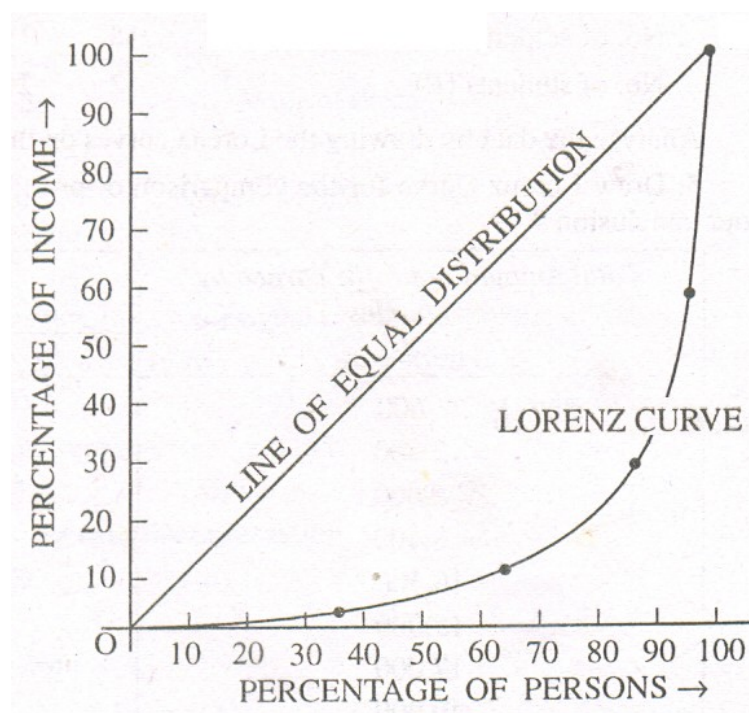| Individual | Income | % population | % income | Cumulative Income % |
|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 |
| 1 | 5000 | 12.5 | 1.204819 | 1.204819 |
| 2 | 12000 | 25 | 2.891566 | 4.096385 |
| 3 | 18000 | 37.5 | 4.337349 | 8.433735 |
| 4 | 30000 | 50 | 7.228916 | 15.66265 |
| 5 | 40000 | 62.5 | 9.638554 | 25.3012 |
| 6 | 60000 | 75 | 14.45783 | 39.75904 |
| 7 | 100000 | 87.5 | 24.09639 | 63.85542 |
| 8 | 150000 | 100 | 36.14458 | 100 |
|  | 415000 |  |  |  |



Example

From the following table giving data regarding income of workers in a factory, draw Lorenz Curve to study inequality of income

The following method for constructing Lorenz Curve.

1.      The size of the item and their frequencies are to be cumulated.

2.      Percentage must be calculated for each cumulation value of the size and frequency of items.

3.      Plot the percentage of the cumulated values of the variable against the percentage of the corresponding cumulated frequencies.  Join these points with as smooth free hand curve.  This curve is called Lorenz curve.

4.      Zero percentage on the X axis must be joined with 100% on Y axis.  This line is called the line of equal distribution.

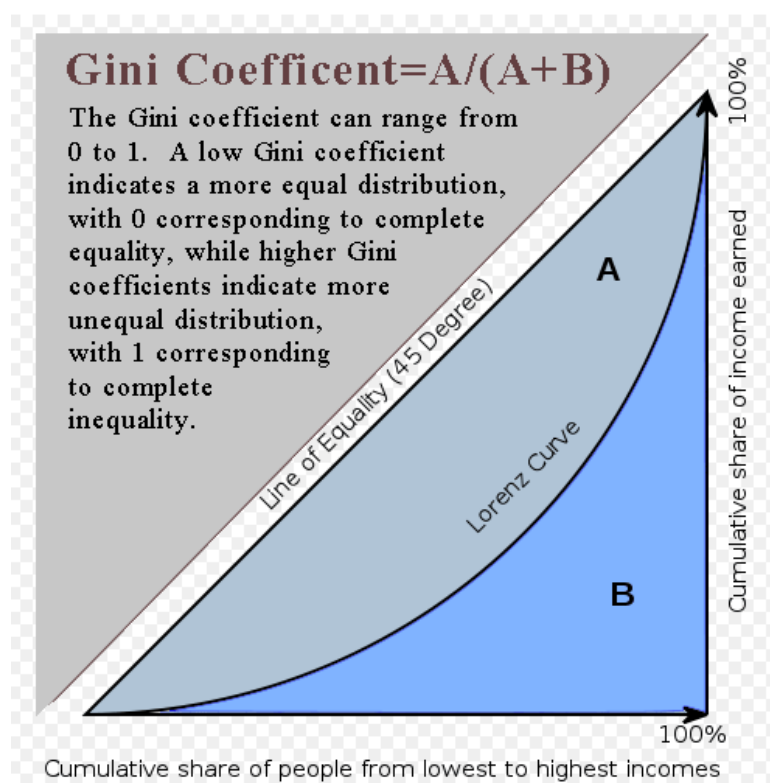| Income | Mid value | Cumulative income | % of cumulative income | No. of workers (f) | Cumulative no. of workers | % of Cumulative no. Of workers |
|---|---|---|---|---|---|---|
| 0-500 | 250 | 250 | 2.94 | 6000 | 6000 | 37.50 |
| 500-1000 | 750 | 1000 | 11.76 | 4250 | 10250 | 64.06 |
| 1000-2000 | 1500 | 2500 | 29.41 | 3600 | 13850 | 86.56 |
| 2000-3000 | 2500 | 5000 | 58.82 | 1500 | 15350 | 95.94 |
| 3000-4000 | 3500 | 8500 | 100.00 | 650 | 16000 | 100.00 |
|  | 8500 |  |  | 16000 |  |  |



**Uses of Lorenz Curve**

1.  To study the variability in a distribution.
2.  To compare the variability relating to a phenomenon for two regions.
3.  To study the changes in variability over a period.

**5.2 - Gini index / Gini coefficient**

A Lorenz curve plots the cumulative percentages of total income received against the cumulative number of recipients, starting with the poorest individual or household. The Gini index measures the area between the Lorenz curve and a hypothetical line of absolute equality, expressed as a percentage of the maximum area under the line. This is the most commonly used measure of inequality. The coefficient varies between 0, which reflects complete equality and 1(100), which indicates complete inequality (one person has all the income or consumption, all others have none). Gini coefficient is found by measuring the areas A and B as marked in the following diagram and using the formula A/(A+B). If the Gini coefficient is to be presented as a ratio or percentage, A/(A+B)×100.



The Gini coefficient (also known as the Gini index or Gini ratio) is a measure of statistical dispersion intended to represent the income distribution of a nation's residents. This is the most commonly used measure of inequality. The coefficient varies between 0, which reflects complete equality and 1, which indicates complete inequality (one person has all the income or consumption, all others have none). It was developed by the Italian statistician and sociologist Corrado Gini in 1912.

**5.3 - Skewness**

We have discussed earlier techniques to calculate the deviations of a distribution from its measure of central tendency (mean / median, mode ). Here we see another measure for that named Skewness. Skewness characterizes the degree of asymmetry of a distribution around its mean. If there is only one mode (peak) in our data (unimodal) , and if the other data are distributed evenly to the left and right of this value, if we plot it in a graph, we get a curve like

this, which is called a normal curve (See figure below). Here we say that there is no skewness or skewness = 0. If there is zero skewness (i.e., the distribution is symmetric) then the mean = median for this distribution.



However data need not always be like this. Sometimes the bulk of the data is at the left and the right tail is longer, we say that the distribution is skewed right or positively skewed. Positive skewness indicates a distribution with an asymmetric tail extending towards more positive values.On the other hand, sometimes the bulk of the data is at is at the right and the left tail is longer, we say that the distribution is skewed left or negatively skewed. Negative skewness indicates a distribution with an asymmetric tail extending towards more negative values"

| Skewed Left | Symmetric | Skewed Right |
|---|---|---|
|  |  |  |

**Tests of Skewness**

There are certain tests to know whether skewness does or does not exist in a frequency distribution.

They are :

1. In a skewed distribution, values of mean, median and mode would not coincide. The values of mean and mode are pulled away and the value of median will be at the centre. In this distribution, mean-Mode = 2/3 (Median - Mode).

2. Quartiles will not be equidistant from median.

3. When the asymmetrical distribution is drawn on the graph paper, it will not give a bell shapedcurve.

4. Sum of the positive deviations from the median is not equal to sum of negative deviations.

5. Frequencies are not equal at points of equal deviations from the mode.

**Nature of Skewness**

Skewness can be positive or negative or zero.

1. When the values of mean, median and mode are equal, there is no skewness.

2. When mean > median > mode, skewness will be positive.

3. When mean < median < mode, skewness will be negative.

**Characteristic of a good measure of skewness**

1. It should be a pure number in the sense that its value should be independent of the unit of the series and also degree of variation in the series.

2. It should have zero-value, when the distribution is symmetrical.

3. It should have a meaningful scale of measurement so that we could easily interpret the measured value.

**Measures of Skewness**

Skewness can be studied graphically and mathematically. When we study Skewness graphically, we can find out whether Skewness is positive or negative or zero. This is what we have shown above.
Mathematically Skewness can be studied as :
(a) Absolute Skewness
(b) Relative or coefficient of skewness
When the skewness is presented in absolute term i.e, in units, it is absolute skewness. If the value of skewness is obtained in ratios or percentages, it is called relative or coefficient of skewness. When skewness is measured in absolute terms, we can compare one distribution with the other if the units of measurement are same. When it is presented in ratios or percentages, comparison become easy. Relative measures of skewness is also called coefficient of skewness.

**(a) Absolute measure of Skewness:**

Skewness can be measured in absolute terms by taking the difference between mean and mode.

Absolute Skewness = $\acute{X}$ – mode

If the value of the mean is greater than mode, the Skewness is positive
If the value of mode is greater than mean, the Skewness is negative

Greater the amount of Skewness (negative or positive) the more tendency towards asymmetry. The absolute measure of Skewness will be proper measure for comparison, and hence, in each series a relative measure or coefficient of Skeweness have to be computed.

**(b) Relative measure of skewness**

There are three important measures of relative skewness.

1. Karl Pearson's coefficient of skewness.

2. Bowley's coefficient of skewness.

3. Kelly's coefficient of skewness.

**(b 1) Karl Pearson's coefficient of Skewness**

The mean, median and mode are not equal in a skewed distribution. The Karl Pearson's measure of skewness is based upon the divergence of mean from mode in a skewed distribution.

Karl Pearson's measure of skewness is sometimes referred to $S_{kp}$

$$S_{kp} = \frac{mean - mode}{standard\,deviation}$$

<u>Properties of Karl Pearson coefficient of Skewness</u>

(1) $-1 \leq Skp \leq 1$.

(2) $Skp = 0 \Rightarrow$ distribution is symmetrical about mean.

(3) $Skp > 0 \Rightarrow$ distribution is skewed to the right.

(4) $Skp < 0 \Rightarrow$ distribution is skewed to the left.

Advantage of Karl Pearson coefficient of Skewness

Skp is independent of the scale. Because (mean-mode) and standard deviation have same scale and it will be canceled out when taking the ratio.

Disadvantage of Karl Pearson coefficient of Skewness

Skp depends on the extreme values.

Example: 1
Calculate the coefficient of skewness of the following data by using Karl Pearson's method for the data 2 3 3 4 4 6 6

$$\mu = \frac{\sum X}{N} = \frac{1 + 2 + 3 + 3 + 4 + 4 + 4}{7} = 3$$

Step 1. Find the mean:

Step 2. Find the standard deviation:

$$\sigma^2 = \frac{(1-3)^2 + (2-3)^2 + (3-3)^2 + (3-3)^2 + (4-3)^2 + (4-3)^2 + (4-3)^2}{7}$$

$$= \frac{4 + 1 + 1 + 1 + 1}{7} = 1.14$$

$$\sigma = \sqrt{1.14} \approx 1.07$$

Then

$$\alpha_3 = \frac{3-4}{1.07} = -0.93$$

Step 3. Find the coefficient of skeness:
Here skewness is negative.

**(b 2) Bowley's coefficient of skewness**

Bowley's formula for measuring skewness is based on quartiles. For a symmetrical distribution, it is seen that $Q_1$, and $Q_3$ areequidistant from median ($Q_2$).

Thus $(Q3 - Q2) - (Q2 - Q1)$ can be taken as an absolute measure of skewness.

$$Skq = \frac{(Q_3 - Q_2) - (Q_2 - Q_1)}{(Q_3 - Q_2) + (Q_2 - Q_1)}$$

$$Skq = \frac{Q_3 - Q_2 - Q_2 + Q_1}{Q_3 - Q_2 + Q_2 - Q_1}$$

$$Skq = \frac{Q_3 + Q_1 - 2Q_2}{Q_3 - Q_1}$$

*Note:*

*In the above equation, where the Qs denote the interquartile ranges. Divide a set of data into two groups (high and low) of equal size at the statistical median if there is an even number of data points, **or** two groups consisting of points on either side of the statistical median itself plus the statistical median if there is an odd number of data points. Find the statistical medians of the low and high groups, denoting these first and third quartiles by Q1 and Q3. The interquartile range is then defined by IQR = Q₃ - Q₁.*

Properties of Bowley's coefficient of skewness

1 $-1 \leq Skq \leq 1$.

2 $Skq = 0 \Rightarrow$ distribution is symmetrical about mean.

3 $Skq > 0 \Rightarrow$ distribution is skewed to the right.

4 $Skq < 0 \Rightarrow$ distribution is skewed to the left.

Advantageof Bowley's coefficient of skewness

Skq does not depend on extreme values.

Disadvantage of Bowley's coefficient of skewness

Skq does not utilize the data fully.

Example

The following table shows the distribution of 128 families according to the number of children.

| No of children | No of families |
|---|---|
| 0 | 20 |
| 1 | 15 |
| 2 | 25 |
| 3 | 30 |
| 4 | 18 |
| 5 | 10 |
| 6 | 6 |
| 7 | 3 |
| 8 or more | 1 |

Compute Bowley's coefficient of skewness

We use formula for measuring Bowley's coefficient of skewness

$$Skq = \frac{Q_3 + Q_1 - 2Q_2}{Q_3 - Q_1}$$

Let us find the necessary values

| No of children | No of families | Cumulative frequency |
|---|---|---|
| 0 | 20 | 20 |
| 1 | 15 | 35 |
| 2 | 25 | 60 |
| 3 | 30 | 90 |
| 4 | 18 | 108 |

| 5 | 10 | 118 |
|---|----|-----|
| 6 | 6 | 124 |
| 7 | 3 | 127 |
| 8 or more | 1 | 128 |

$$Q_1 = \left(\frac{128+1}{4}\right)^{th} Observation$$

$$= (32.25)^{th} \text{ observation}$$

$$= 1$$

$$Q_2 = \left(\frac{128+1}{2}\right)^{th} Observation$$

$$= (64.5)^{th} \text{ observation}$$

$$= 3$$

$$Q_3 = 3\left(\frac{128+1}{4}\right)^{th} Observation$$

$$= (96.75)^{th} \text{ observation}$$

$$= 4$$

$$Skq = \frac{4+1-2(3)}{4-1}$$

$$¿ -\frac{1}{3} = -0.333$$

Since Skq < 0 distribution is skewed left

**(b 3) Kelly's coefficient of skewness**

Bowley's measure of skewness is based on the middle 50% of the observations because it leaves 25% of the observations on each extreme of the distribution. As an improvement over Bowley's measure, Kelly has suggested a measure based on $P_{10}$ and, $P_{90}$ so that only 10% of the observations on each extreme are ignored.

$$Sp = \frac{\left(P_{90} - P_{50}\right) - \left(P_{50} - P_{10}\right)}{\left(P_{90} - P_{50}\right) + \left(P_{50} - P_{10}\right)}$$

$$Sp = \frac{P_{90} - P_{50} - P_{50} + P_{10}}{P_{90} - P_{50} + P_{50} - P_{10}}$$
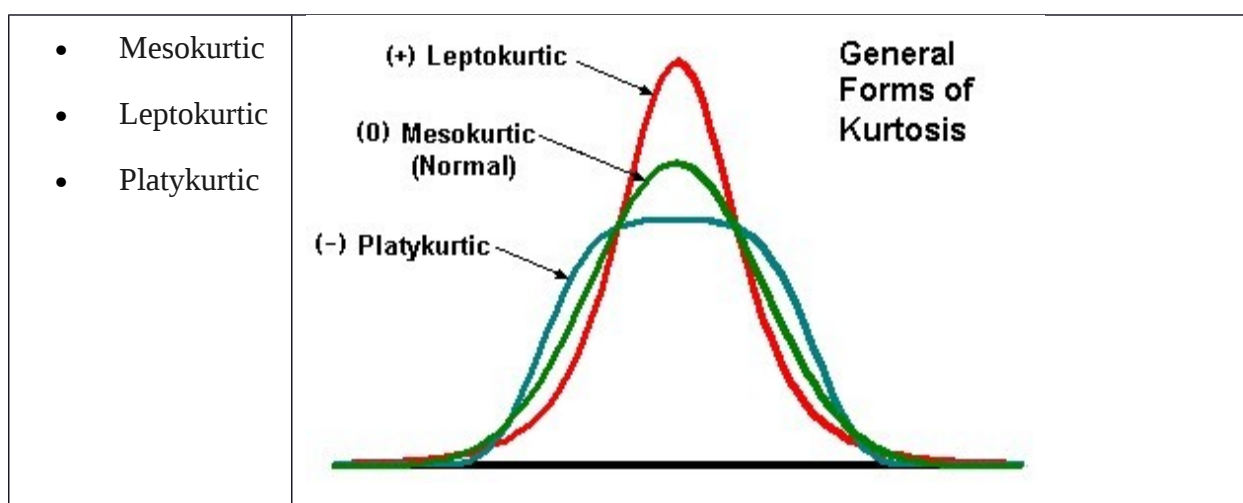
$$Sp = \frac{P_{90} + P_{10} - 2P_{50}}{P_{90} - P_{10}}$$

### 5.4 - KURTOSIS

As we saw above, Skewness is a measure of symmetry, or more precisely, the lack of symmetry. A distribution, or data set, is symmetric if it looks the same to the left and right of the center point.

Kurtosis is a measure of whether the data are peaked or flat relative to a normal distribution. That is, data sets with high kurtosis tend to have a distinct peak near the mean, decline rather rapidly, and have heavy tails. Data sets with low kurtosis tend to have a flat top near the mean rather than a sharp peak. A uniform distribution would be the extreme case. Kurtosis has its origin in the Greek word 'Bulginess.'

Distributions of data and probability distributions are not all the same shape. Some are asymmetric and skewed to the left or to the right. Other distributions are bimodal and have two peaks. In other words there are two values that dominate the distribution of values. Another feature to consider when talking about a distribution is not just the number of peaks but the shape of them. Kurtosis is the measure of the peak of a distribution, and indicates how high the distribution is around the mean. The kurtosis of a distributions is in one of three categories of classification:

| | |
|---|---|
| • Mesokurtic<br><br>• Leptokurtic<br><br>• Platykurtic |  |

We will consider each of these classifications in turn.

### Mesokurtic

Kurtosis is typically measured with respect to the normal distribution. A distribution that is peaked in the same way as any normal distribution, not just the standard normal distribution, is said to be mesokurtic. The peak of a mesokurtic distribution is neither high nor low, rather it is considered to be a baseline for the two other classifications. Besides normal distributions, binomial distributions for which p is close to 1/2 are considered to be mesokurtic.

### Leptokurtic

A leptokurtic distribution is one that has kurtosis greater than a mesokurtic distribution. Leptokurtic distributions are identified by peaks that are thin and tall. The tails of these distributions, to both the right and the left, are thick and heavy. Leptokurtic distributions are named by the prefix "lepto" meaning "skinny."

There are many examples of leptokurtic distributions. One of the most well known leptokiurtic distributions is Student's t distribution.

**Platykurtic**

The third classification for kurtosis is platykurtic. Platykurtic distributions are those that have a peak lower than a mesokurtic distribution. Platykurtic distributions are characterized by a certain flatness to the peak, and have slender tails. The name of these types of distributions come from the meaning of the prefix "platy" meaning "broad."

All uniform distributions are platykurtic. In addition to this the discrete probability distribution from a single flip of a coin is platykurtic.



**Measures of Kurtosis**

Moment ratio and Percentile Coefficient of kurtosis are used to measure the kurtosis

Moment Coefficient of Kurtosis= $\beta_2$ = $\dfrac{M_4}{M_2^2}$

Where M4 = 4th moment and M2 = 2nd moment

If $\beta_2$ = 3, the distribution is said to be normal. (ie mesokurtic)

If $\beta_2 > ¿$ 3, the distribution is more peaked to curve is lepto kurtic.

If $\beta_2 < ¿$ 3, the distribution is said to be flat topped and the curve is platy kurtic.

Percentile Coefficient of Kurtosis = $k = \dfrac{Q.D.}{P_{90} - P_{10}}$

where $Q.D. = \dfrac{1}{2}(Q_3 - Q_1)$ is the semi-interquartile range. For normal distribution this has the value 0.263.

A normal random variable has a kurtosis of 3 irrespective of its mean or standard deviation. If a random variable's kurtosis is greater than 3, it is said to be Leptokurtic. If its kurtosis is less than 3, it is said to be Platykurtic.

Thus we conclude our discussion by saying that kurtosis is any measure of the 'peakedness' of a distribution. The height and sharpness of the peak relative to the rest of the data are measured by a number called kurtosis. Higher values indicate a higher, sharper peak; lower values indicate a lower, less distinct peak. This occurs because, higher kurtosis means more of the variability is due to a few extreme differences from the mean, rather than a lot of modest differences from the mean. A normal distribution has kurtosis exactly 3. Any distribution with kurtosis =3 is called mesokurtic. A distribution with kurtosis <3 is called platykurtic. Compared to a normal distribution, its central peak is lower and broader, and its tails are shorter and thinner. A distribution with kurtosis >3 is called leptokurtic. Compared to a normal distribution, its central peak is higher and sharper, and its tails are longer and fatter.

**Comparison among dispersion, skewness and kurtosis**

Dispersion, Skewness and Kurtosis are different characteristics of frequency distribution. Dispersion studies the scatter of the items round a central value or among themselves. It does not show the extent to which deviations cluster below an average or above it. Skewness tells us about the cluster of the deviations above and below a measure of central tendency. Kurtosis studies the concentration of the items at the central part of a series. If items concentrate too much at the centre, the curve becomes 'leptokurtic' and if the concentration at the centre is comparatively less, the curve becomes 'platykurtic'.

## Module V - CORRELATION AND REGRESSION ANALYSIS

**Meaning of Correlation**

Variables which are related in some way are commonly used in economics and other fields of study. For example, relation between the price of gold and the demand for it.

- between economic growth and life expectancy
- between fertiliser use and crop yield
- between hours of work and wage

Correlation examines the relationships between these pairs of variables. Correlation measures the association between two variables. Correlation is a statistical technique which tells us if two variables are related. For example, consider the variables family income and family expenditure. It is well known that income and expenditure increase or decrease together. Thus they are related in the sense that change in any one variable is accompanied by change in the other variable. Again price and demand of a commodity are related variables; when price increases demand will tend to decreases and vice versa. If the change in one variable is accompanied by a change in the other, then the variables are said to be correlated. We can therefore say that family income and family expenditure, price and demand are correlated.

Correlation can tell us something about the relationship between variables. It is used to understand:  a) whether the relationship is positive or negative

b) the strength of relationship.

Correlation is a powerful tool that provides these vital pieces of information. In the case of family income and family expenditure, it is easy to see that they both rise or fall together in the same direction. This is called positive correlation. In case of price and demand, change occurs in the opposite direction so that increase in one is accompanied by decrease in the other. This is called negative correlation.

According to the number of variables, correlation is said to be of the following three types viz;

**(i)** Simple Correlation: In simple correlation, we study the relationship between two variables. Of these two variables one is principal and the other is secondary? For instance, income and expenditure, price and demand etc. Here income and price are principal variables while expenditure and demand are secondary variables.

**(ii)** Partial Correlation: If in a given problem, more than two variables are involved and of these variables we study the relationship between only two variables keeping the other variables constant, correlation is said to be partial. It is so because the effect of other variables is assumed to be constant

**(iii)** Multiple Correlations: Under multiple correlations, the relationship between two and more variables is studied jointly. For instance, relationship between rainfall, use of fertilizer, manure on per hectare productivity of wheat crop.

**Coefficient of Correlation**

Correlation is measured by what is called coefficient of correlation (r). A correlation coefficient is a statistical measure of the degree to which changes to the value of one variable predict change to the value of another. Correlation coefficients are expressed as values between +1 and -1. Its numerical value gives us an indication of the strength of relationship. In general, r > 0 indicates positive relationship, r < 0 indicates negative relationship while r = 0 indicates no relationship (or that the variables are independent and not related). Here r = +1.0 describes a perfect positive correlation and r = −1.0 describes a perfect negative correlation. Closer the coefficients are to +1.0 and −1.0, greater is the strength of the relationship between the variables. As a rule of thumb, the following guidelines on strength of relationship are often useful (though many experts would somewhat disagree on the choice of boundaries).

Correlation is only appropriate for examining the relationship between meaningful quantifiable data (e.g. air pressure, temperature) rather than categorical data such as gender, favourite colour etc. A key thing to remember when working with correlations is never to assume a correlation means that a change in one variable causes a change in another. Sales of personal computers and athletic shoes have both risen strongly in the last several years and there is a high correlation between them, but you cannot assume that buying computers causes people to buy athletic shoes (or vice versa).

The second caution is that the Pearson correlation technique (which we are about to see) works best with linear relationships: as one variable gets larger (or smaller), the other gets larger (or smaller) in direct proportion. It does not work well with curvilinear relationships (in which the relationship does not follow a straight line). An example of a curvilinear relationship is age and health care. They are related, but the relationship doesn't follow a straight line. Young children and older people both tend to use much more health care than teenagers or young adults. (In such cases, the technique of 'multiple regression' can be used to examine curvilinear relationships)

### METHODS OF MEASURING CORRELATION

I.  Graphical Method

   (a) Scatter Diagram

   (b) Correlation Graph

II.  Algebraic Method (Coefficient of Correlation)

   (a) Karl Pearson's Coefficient of Correlation

   (b) Spearman's Rank Correlation Coefficient

   (c)


I.  **(a) Scatter Diagram**

Scatter Diagram (also called scatter plot, X–Y graph) is a graph that shows the relationship between two quantitative variables measured on the same individual. Each individual in the data set is represented by a point in the scatter diagram. The predictor variable is plotted on the

horizontal axis and the response variable is plotted on the vertical axis. Do not connect the points when drawing a scatter diagram. The scatter diagram graphs pairs of numerical data, with one variable on each axis, to look for a relationship between them. If the variables are correlated, the points will fall along a line or curve. The better the correlation, the tighter the points will hug the line. Scatter Diagram is a graphical measure of correlation.

Examples of Scatter Diagram. Given below each diagram is the value of correlation.



Note that the value shows how good the correlation is (not how steep the line is), and if it is positive or negative.

Scatter Diagram Procedure

1. Collect pairs of data where a relationship is suspected.

2. Draw a graph with the independent variable on the horizontal axis and the dependent variable on the vertical axis. For each pair of data, put a dot or a symbol where the x-axis value intersects the y-axis value. (If two dots fall together, put them side by side, touching, so that you can see both.)

3. Look at the pattern of points to see if a relationship is obvious. If the data clearly form a line or a curve, you may stop. The variables are correlated.

The data set below represents a random sample of 5 workers in a particular industry. The productivity of each worker was measured at one point in time, and the worker was asked the number of years of job experience. The dependent variable is productivity, measured in number of units produced per day, and the independent variable is experience, measured in years.

| Worker | y=Productivity(output/day) | x=Experience(in years) |
|--------|----------------------------|------------------------|
| 1 | 33 | 10 |
| 2 | 19 | 6 |
| 3 | 32 | 12 |
| 4 | 26 | 8 |
| 5 | 15 | 4 |

This scatter diagram tell us that the two variables, productivity and experience, are positively correlated.

**Merits of Scatter Diagram Method:**

1. It is an easy way of finding the nature of correlation between two variables.

2. By drawing a line of best fit by free hand method through the plotted dots, the method can be used for estimating the missing value of the dependent variable for a given value of independent variable.

3. Scatter diagram can be used to find out the nature of linear as well as non-linear correlation.

4. The values of extreme observations do not affect the method.

**<u>Demerits of Scatter Diagram Method:</u>**

It gives only rough idea of how the two variables are related. It gives an idea about the direction of correlation and also whether it is high or low. But this method does not give any quantitative measure of the degree or extent of correlation.

**I (b) Correlation Graph**

Correlation graph is also used as a measure of correlation. When this method is used the correlation graph is drawn and the direction of curve is examined to understand the nature of correlation. Under this method, separate curves are drawn for the X variable and Y variable on the same graph paper. The values of the variable are taken as ordinates of the points plotted. From the direction and closeness of the two curves we can infer whether the variables are related. If both the curves move in the same direction (upward or downward), correlation is said to be positive. If the curves are moving in the opposite direction, correlation is said to be negative.

But correlation graphs are not capable of doing anything more than suggesting the fact of a possible relationship between two variables. We can neither establish any casual relationship between two variables nor obtain the exact degree of correlation through them. They only tell us whether the two variables are positively or negatively correlated. Example of a graph is given below.

All Test Scores



## II.    Algebraic Method (Coefficient of Correlation)

II. (a)  **Karl Pearson's Coefficient of Correlation (Pearson product-moment correlation coefficient)**

Karl Pearson's Product-Moment Correlation Coefficient or simply Pearson's Correlation Coefficient for short, is one of the important methods used in Statistics to measure Correlation between two variables. Karl Pearson was a British mathematician, statistician, lawyer and a eugenicist. He established the discipline of mathematical statistics. He founded the world's first statistics department In the University of London in the year 1911. He along with his colleagues Weldon and Galton founded the journal 'Biometrika' whose object was the development of statistical theory.

The Pearson product-moment correlation coefficient (r) is a common measure of the correlation between two variables X and Y. When measured in a population the Pearson Product Moment correlation is designated by the Greek letter rho (?). When computed in a sample, it is designated by the letter "r" and is sometimes called "Pearson's r." Pearson's correlation reflects the degree of linear relationship between two variables.

Mathematical Formula:--

The quantity r, called the linear correlation coefficient, measures the strength and the direction of a linear relationship between two variables. (The linear correlation coefficient is a measure of the strength of linear relation between two quantitative variables. We use the Greek

letter ρ (rho) to represent the population correlation coefficient and r to represent the sample correlation coefficient.)

Correlation coefficient for ungrouped data

$$r = \frac{\sum_{i=1}^{n} (X_i - \acute{X})(Y_i - \acute{Y})}{n \sigma_X \sigma_Y}$$

Where

$X_i$ is the i$^{th}$ observation of the variable X

$Y_i$ is the i$^{th}$ observation of the variable Y

$\acute{X}$ is the mean of the observations of the variable X

$\acute{Y}$ is the mean of the observations of the variable Y

n is the number of pairs of observations of X and Y

$\sigma_X$ is the standard deviation of the variable X

$\sigma_Y$ is the standard deviation of the variable Y

The above formula may be presented in the following form

$$r = \frac{\sum_{i=1}^{n} (X_i - \acute{X})(Y_i - \acute{Y})}{\sqrt{\sum_{i=1}^{n} (X_i - \acute{X})^2} \sqrt{\sum_{i=1}^{n} (Y_i - \acute{Y})^2}}$$

The same may be computed using Pearson product-moment correlation coefficient formula as shown below.

$$r = \frac{n \sum_{i=1}^{n} X_i Y_i - \sum_{i=1}^{n} X_i \sum_{i=1}^{n} Y_i}{\sqrt{n \sum_{i=1}^{n} X_i^2 - \left( \sum_{i=1}^{n} X_i \right)^2} \sqrt{n \sum_{i=1}^{n} Y_i^2 - \left( \sum_{i=1}^{n} Y_1 \right)^2}}$$

| Year (i) | Annual advertising expenditure Xi | Annual Sales |
|---|---|---|
| 1 | 10 | 20 |
| 2 | 12 | 30 |
| 3 | 14 | 37 |
| 4 | 16 | 50 |
| 5 | 18 | 56 |
| 6 | 20 | 78 |
| 7 | 22 | 89 |
| 8 | 24 | 100 |
| 9 | 26 | 120 |
| 10 | 28 | 110 |

Compute the necessary values and substitute in the formula, we will solve using both formula.

We get $\acute{X}=\left(\sum X_i/n\right)=\dfrac{190}{10}=19.\ \acute{Y}=\left(\sum Y_i/n\right)=\dfrac{690}{10}=69.$

| Year (i) | Xi | Annual Sales (Yi) | $\left(X_i-\acute{X}\right)$ | $\left(Y_i-Y\right)$ | $\left(X_i-\acute{X}\right)^2$ | $\left(Y_i-Y\right)^2$ | $\left(X_i-\acute{X}\right)\left(Y_i-\acute{Y}\right)$ |
|---|---|---|---|---|---|---|---|
| 1 | 10 | 20 | -9 | -49 | 81 | 2401 | 441 |
| 2 | 12 | 30 | -7 | -39 | 49 | 1521 | 273 |
| 3 | 14 | 37 | -5 | -32 | 25 | 1024 | 160 |
| 4 | 16 | 50 | -3 | -19 | 9 | 361 | 57 |
| 5 | 18 | 56 | -1 | -13 | 1 | 169 | 13 |
| 6 | 20 | 78 | 1 | 9 | 1 | 81 | 9 |
| 7 | 22 | 89 | 3 | 20 | 9 | 400 | 60 |
| 8 | 24 | 100 | 5 | 31 | 25 | 961 | 155 |
| 9 | 26 | 120 | 7 | 51 | 49 | 2601 | 357 |
| 10 | 28 | 110 | 9 | 41 | 81 | 1681 | 369 |
|  | 190 | 690 | 0 | 0 | 330 | 11200 | 1894 |

We make the additional computations for the Pearson product-moment correlation coefficient formula.

| $X_iY_i$ | $X_i^2$ | $Y_i^2$ |
|---|---|---|
| 200 | 100 | 400 |
| 360 | 144 | 900 |
| 518 | 196 | 1369 |
| 800 | 256 | 2500 |
| 1008 | 324 | 3136 |
| 1560 | 400 | 6084 |
| 1958 | 484 | 7921 |
| 2400 | 576 | 10000 |
| 3120 | 676 | 14400 |
| 3080 | 784 | 12100 |
| 15004 | 3940 | 58810 |

Substitute the values in the respective formula.

Using the basic formula $r=\dfrac{\displaystyle\sum_{i=1}^{n}\left(X_i-\acute{X}\right)\left(Y_i-\acute{Y}\right)}{\sqrt{\displaystyle\sum_{i=1}^{n}\left(X_i-\acute{X}\right)^2}\sqrt{\displaystyle\sum_{i=1}^{n}\left(Y_i-\acute{Y}\right)^2}}$

$$r = \frac{1894}{\sqrt{330}\sqrt{11200}} = 0.985$$

Now let us re do the problem using Pearson product-moment correlation coefficient formula

$$r = \frac{n\sum_{i=1}^{n} X_i Y_i - \sum_{i=1}^{n} X_i \sum_{i=1}^{n} Y_i}{\sqrt{n\sum_{i=1}^{n} X_i^2 - \left(\sum_{i=1}^{n} X_i\right)^2}\sqrt{n\sum_{i=1}^{n} Y_i^2 - \left(\sum_{i=1}^{n} Y_1\right)^2}}$$

$$r = \frac{10 \times 15004 - 190 \times 690}{\sqrt{10 \times 3940 - 190^2}\sqrt{10 \times 58810 - 690^2}} = 0.985$$

The correlation coefficient between annual advertising expenditure and annual sales revenue is 0.985. This is a positive value and is very close to 1. So it implies there is very strong corelation between annual advertising expenditure and annual sales revenue.

**Properties of Correlation coefficient**

1. The correlation coefficient lies between -1 & +1 symbolically ( $-1 \leq r \geq 1$ )

2. The correlation coefficient is independent of the change of origin & scale.

3. The coefficient of correlation is the geometric mean of two regression coefficient.

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

The one regression coefficient is (+ve) other regression coefficient is also (+ve) correlation coefficient is (+ve)

**Assumptions of Pearson's Correlation Coefficient**

1. 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.

2. Cause and effect relation exists between different forces operating on the item of the two variable series.

**Advantages of Pearson's Coefficient**

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

**Disadvantages**

While 'r' (correlation coefficient) is a powerful tool, it has to be handled with care.

1. The most used correlation coefficients only measure linear relationship. It is therefore perfectly possible that while there is strong non-linear relationship between the variables, r is close to 0 or even 0. In such a case, a scatter diagram can roughly indicate the existence or otherwise of a non-linear relationship.

2. One has to be careful in interpreting the value of 'r'. For example, one could compute 'r' between the size of shoe and intelligence of individuals, heights and income. Irrespective of the value of 'r', it makes no sense and is hence termed chance or non-sense correlation.

3. 'r' should not be used to say anything about cause and effect relationship. Put differently, by examining the value of 'r', we could conclude that variables X and Y are related. However the same value of 'r' does not tell us if X influences Y or the other way round. Statistical correlation should not be the primary tool used to study causation, because of the problem with third variables.

## 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 r2 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: An example**

Suppose: r = 0.60 in one case and r = 0.30 in another case. It does not mean that the first correlation is twice as strong as the second the 'r' can be understood by computing the value of $r^2$.

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

When 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.

II. (b) **Spearman's Rank Correlation Coefficient**

The Spearman's rank-order correlation is the nonparametric version of the Pearson product-moment correlation. Spearman's correlation coefficient, ( $\rho\ Greek\ lap\ habet\ Rho, \vee r_s ¿$ measures the strength of association between two ranked variables. Data which are arranged in numerical order, usually from largest to smallest and numbered 1,2,3 ---- are said to be in ranks or ranked data.. These ranks prove useful at certain times when two or more values of one variable are the same. The coefficient of correlation for such type of data is given by Spearman rank difference correlation coefficient.

Spearman Rank Correlation Coefficient uses ranks to calculate correlation. The Spearman Rank Correlation Coefficient is its analogue when the data is in terms of ranks. One can therefore also call it correlation coefficient between the ranks. The Spearman's rank-order correlation is

used when there is a monotonic relationship between our variables. A monotonic relationship is a relationship that does one of the following: (1) as the value of one variable increases, so does the value of the other variable; or (2) as the value of one variable increases, the other variable value decreases. A monotonic relationship is an important underlying assumption of the Spearman rank-order correlation. It is also important to recognize the assumption of a monotonic relationship is less restrictive than a linear relationship (an assumption that has to be met by the Pearson product-moment correlation). The middle image above illustrates this point well: A non-linear relationship exists, but the relationship is monotonic and is suitable for analysis by Spearman's correlation, but not by Pearson's correlation.

Let us make the relevance of use of Spearman Rank Correlation Coefficient with the aid of an example. As an example, let us consider a musical talent contest where 10 competitors are evaluated by two judges, A and B. Usually judges award numerical scores for each contestant after his/her performance.

A product moment correlation coefficient of scores by the two judges hardly makes sense here as we are not interested in examining the existence or otherwise of a linear relationship between the scores. What makes more sense is correlation between ranks of contestants as judged by the two judges. Spearman Rank Correlation Coefficient can indicate if judges agree to each other's views as far as talent of the contestants are concerned (though they might award different numerical scores) - in other words if the judges are unanimous.

The numerical value of the correlation coefficient, $r_s$, ranges between -1 and +1. The correlation coefficient is the number indicating the how the scores are relating.

In general,

- $r_s > 0$ implies positive agreement among ranks
- $r_s < 0$ implies negative agreement (or agreement in the reverse direction)
- $r_s = 0$ implies no agreement

Closer $r_s$ is to 1, better is the agreement while rs closer to -1 indicates strong agreement in the reverse direction.

The formula for finding Spearman Rank Correlation Coefficient is

$$r_s = 1 - \frac{6 \sum_{i=1}^{n} (X_i + Y_i)^2}{n(n^2 - 1)}$$

Where
$X_i$ is the rank of the $i^{th}$ observation of the variable X
$Y_i$ is the rank of the $i^{th}$ observation of the variable Y
n is the number of payers of observations

Let us calculate Spearman Rank Correlation Coefficient for our example of the musical talent contest where 10 competitors are evaluated by two judges, A and B. The scores are given below.

| Contestant | Rating by judge 1 | Rating by judge 2 |
|---|---|---|
| 1 | 1 | 2 |
| 2 | 2 | 4 |

| | | |
|---|---|---|
| 3 | 3 | 5 |
| 4 | 4 | 1 |
| 5 | 5 | 3 |
| 6 | 6 | 6 |
| 7 | 7 | 7 |
| 8 | 8 | 9 |
| 9 | 9 | 10 |
| 10 | 10 | 8 |

Let us first make the necessary calculations

| Contestant | Rating by judge 1 ($X_i$) | Rating by judge 2($Y_i$) | $X_i - Y_i$ | $(X_1 - Y_1)^2$ |
|---|---|---|---|---|
| 1 | 1 | 2 | -1 | 1 |
| 2 | 2 | 4 | -2 | 4 |
| 3 | 3 | 5 | -2 | 4 |
| 4 | 4 | 1 | 3 | 9 |
| 5 | 5 | 3 | 2 | 4 |
| 6 | 6 | 6 | 0 | 0 |
| 7 | 7 | 7 | 0 | 0 |
| 8 | 8 | 9 | -1 | 1 |
| 9 | 9 | 10 | -1 | 1 |
| 10 | 10 | 8 | 2 | 4 |
| | | | | 28 |

$$r_s = 1 - \frac{6 \sum_{i=1}^{n} (X_i + Y_i)^2}{n(n^2 - 1)} = r_s = 1 - \frac{6 \times 28}{\frac{10 \times \dot{\iota}}{\dot{\iota}} \quad 10^2 - \dot{\iota}}$$

Spearman Rank Correlation Coefficient tries to assess the relationship between ranks without making any assumptions about the nature of their relationship. Hence it is a non-parametric measure - a feature which has contributed to its popularity and wide spread use.

**Interpretation of Rank Correlation Coefficient (R)**

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

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

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

4. If R = 0, then there is no correlation

**Advantages Spearman's Rank Correlation**

1. This method is simpler to understand and easier to apply compared to karlearson's correlation method.

2. This method is useful where we can give the ranks and not the actual data. (qualitative term)

3. This method is to use where the initial data in the form of ranks.

## Disadvantages Spearman's Rank Correlation

1. It cannot be used for finding out correlation in a grouped frequency distribution.

2. This method should be applied where N exceeds 30.

3. As Spearman's rank only uses rank, it is not affected by significant variations in readings. As long as the order remains the same, the coefficient will stay the same. As with any comparison, the possibility of chance will have to be evaluated to ensure that the two quantities are actually connected.

4. A significant correlation does not necessarily mean cause and effect.

## Advantages of Correlation studies

1. Show the amount (strength) of relationship present.

2. Can be used to make predictions about the variables under study.

3. Can be used in many places, including natural settings, libraries, etc.

4. Easier to collect co relational data

## Importance of Correlation
1. Most of the variables show some kind of relationship. For instance, there is relationship between price and supply, income and expenditure etc. With the help of correlation analysis we can measure in one figure the degree of relationship.
2. Once we know that two variables are closely related, we can estimate the value of one variable given the value of another. This is known with the help of regression.
3. Correlation analysis contributes to the understanding of economic behaviour, aids in locating the critically important variables on which others depend.
4. Progressive development in the methods of science and philosophy has been characterized by increase in the knowledge of relationship. In nature also one finds multiplicity of interrelated forces.
5. The effect of correlation is to reduce the range of uncertainty. The prediction based on correlation analysis is likely to be more variable and near to reality.
## Uses of Correlation in Economics

Correlation is used in Economics for decision making. Correlation concepts used in Economics are more explanatory in nature and less prescriptive. Many observations made by economists are parlayed into normative policy proposals. Economics is data-driven, statistically presented fields of study. As data has been collected over time, analysts have looked to identify meaningful statistical correlations to help explain phenomena, identify trends, make predictions and better understand exchanges between actors. The empirical observations of economic agents are used to drive and test economic assumption.

Economics is a social science, so any correlations would be a means of explaining human action. Not all economists agree about the usefulness of statistical correlation, but almost all macroeconomic analysis is done through correlation analysis. This reaches its apex with

econometrics, which uses regression analysis to distinguish between correlation and causation in the hopes of making accurate forecasts.

**Properties of Correlation Coefficient**

Property 1:   Correlation Coefficientis independent of the change of origin and scale. That is, the correlation coefficient does not change the measurement scale.

Property 2:    The sign of the linear correlation coefficient is shared by the covariance.

Property 3: Linear correlation coefficient cannot exceed 1 numerically. In other words it lies between 1and -1.

Property 4:   If the linear correlation coefficient takes values closer to 1, the correlation is strong and negative, and will become stronger the closer it approaches 1.

Property 5:     If the linear correlation coefficient takes values closer to −1, the correlation is strong and negative, and will become stronger the closer r approaches −1.

Property 6: If the linear correlation coefficient takes values close to 0, the correlation is weak.

Property 7: If the linear correlation coefficient takes values close to 1 the correlation is strong and positive, and will become stronger the closer r approaches 1.

Property 8: Two independent variables are uncorrelated but the converse is not true.

Property 9: If r = 1 or r = −1, there is perfect correlation and the line on the scatter plot is increasing or decreasing respectively.

Property 10: If r = 0, there is no linear correlation.

**Interpretation of Correlation Coefficient**

Correlation refers to a technique used to measure the relationship between two or more variables. When two things are correlated, it means that they vary together. Positive correlation means that high scores on one are associated with high scores on the other, and that low scores on one are associated with low scores on the other. Negative correlation, on the other hand, means that high scores on the first thing are associated with low scores on the second. Negative correlation also means that low scores on the first are associated with high scores on the second. An example is the correlation between body weight and the time spent on a weight-loss program. If the program is effective, the higher the amount of time spent on the program, the lower the body weight. Also, the lower the amount of time spent on the program, the higher the body weight.

As we have already stated the correlation coefficient is a number between -1 and 1 that indicates the strength of the linear relationship between two variables. The interpretation of correlation is mainly based on the value of correlation.

To   interpret   correlations,   four   pieces   of   information   are   necessary.
1.   The   numerical   value   of   the   correlation   coefficient.   Correlation   coefficients   can   vary numerically between 0.0 and 1.0. The closer the correlation is to 1.0, the stronger the relationship

between the two variables. A correlation of 0.0 indicates the absence of a relationship. If the correlation coefficient is –0.80, which indicates the presence of a strong relationship.

2. The sign of the correlation coefficient. A positive correlation coefficient means that as variable 1 increases, variable 2 increases, and conversely, as variable 1 decreases, variable 2 decreases. In other words, the variables move in the same direction when there is a positive correlation. A negative correlation means that as variable 1 increases, variable 2 decreases and vice versa. In other words, the variables move in opposite directions when there is a negative correlation. The negative sign indicates that as class size increases, mean reading scores decrease.

3. The statistical significance of the correlation. A statistically significant correlation is indicated by a probability value of less than 0.05. This means that the probability of obtaining such a correlation coefficient by chance is less than five times out of 100, so the result indicates the presence of a relationship. For -0.80 there is a statistically significant negative relationship between class size and reading score ($p < .001$), such that the probability of this correlation occurring by chance is less than one time out of 1000.

4. The effect size of the correlation. For correlations, the effect size is called the coefficient of determination and is defined as $r^2$. The coefficient of determination can vary from 0 to 1.00 and indicates that the proportion of variation in the scores can be predicted from the relationship between the two variables. For r = -0.80 the coefficient of determination is 0.65, which means that 65% of the variation in mean reading scores among the different classes can be predicted from the relationship between class size and reading scores. (Conversely, 35% of the variation in mean reading scores cannot be explained.)

For a quick interpretation you may use the following.

| Value of 'r' | Interpretation |
|---|---|
| 1.0 | Perfect correlation |
| 0 to 1 | The two variables tend to increase or decrease together. |
| 0.0 | The two variables do not vary together at all. |
| -1 to 0 | One variable increases as the other decreases. |
| -1.0 | Perfect negative or inverse correlation. |

If 'r' is far from zero, there are four possible explanations:

: Changes in the X variable causes a change the value of the Y variable.

: Changes in the Y variable causes a change the value of the X variable.

: Changes in another variable influence both X and Y.

: X and Y don't really correlate at all, and you just happened to observe such a strong correlation by chance.

Another quick reference table

      If r = +.70 or higher Very strong positive relationship

+.40 to +.69 Strong positive relationship

+.30 to +.39 Moderate positive relationship

+.20 to +.29 weak positive relationship

+.01 to +.19 No or negligible relationship

-.01 to -.19 No or negligible relationship

-.20 to -.29 weak negative relationship

-.30 to -.39 Moderate negative relationship

-.40 to -.69 Strong negative relationship

-.70 or higher Very strong negative relationship

Note that a correlation can only indicate the presence or absence of a relationship, not the nature of the relationship. Correlation is not causation. There is always the possibility that a third variable influenced the results. For example, in a college the students in the small classes scored higher in maths exam than the students in the large classes, but it could also be because they were from better schools or they had higher quality teachers.

## REGRESSION ANALYSIS

If two variables are significantly correlated, and if there is some theoretical basis for doing so, it is possible to predict values of one variable from the other. This observation leads to a very important concept known as 'Regression Analysis'.

Regression analysis, in general sense, means the estimation or prediction of the unknown value of one variable from the known value of the other variable. It is one of the most important statistical tools which is extensively used in almost all sciences – Natural, Social and Physical. It is specially used in business and economics to study the relationship between two or more variables that are related causally and for the estimation of demand and supply graphs, cost functions, production and consumption functions and so on.

Prediction or estimation is one of the major problems in almost all the spheres of human activity. The estimation or prediction of future production, consumption, prices, investments, sales, profits, income etc. are of very great importance to business professionals. Similarly, population estimates and population projections, GNP, Revenue and Expenditure etc. are indispensable for economists and efficient planning of an economy.

Regression analysis was explained by M. M. Blair as follows:
"Regression analysis is a mathematical measure of the average relationship between two or more variables in terms of the original units of the data."
Regression Analysis is a very powerful tool in the field of statistical analysis in predicting the value of one variable, given the value of another variable, when those variables are related to each other. Regression Analysis is mathematical measure of average relationship between two or more variables. Regression analysis is a statistical tool used in prediction of value of unknown variable from known variable.

Advantages of Regression Analysis

1. Regression analysis provides estimates of values of the dependent variables from the values of independent variables.

2. Regression analysis also helps to obtain a measure of the error involved in using the regression line as a basis for estimations.

3. Regression analysis helps in obtaining a measure of the degree of association or correlation that exists between the two variable.

Assumptions in Regression Analysis

1. Existence of actual linear relationship.

2. The regression analysis is used to estimate the values within the range for which it is valid.

3. The relationship between the dependent and independent variables remains the same till the regression equation is calculated.

4. The dependent variable takes any random value but the values of the independent variables are fixed.

5. In regression, we have only one dependant variable in our estimating equation. However, we can use more than one independent variable.

**Regression line**

A regression line summarizes the relationship between two variables in the setting when one of the variables helps explain or predict the other.

A regression line is a straight line that describes how a response variable y changes as an explanatory variable x changes. A regression line is used to predict the value of y for a given value of x. Regression, unlike correlation, requires that we have an explanatory variable and a response variable.

Regression line is the line which gives the best estimate of one variable from the value of any other given variable. The regression line gives the average relationship between the two variables in mathematical form.

For two variables X and Y, there are always two lines of regression –

Regression line of X on Y : gives the best estimate for the value of X for any specific given values of Y :

$$X = a + b\ Y$$

Where

      a = X – intercept
      b = Slope of the line
      X = Dependent variable
      Y = Independent variable

Regression line of Y on X : gives the best estimate for the value of Y for any specific given values of X

$$Y = a + bx$$

Where

a = Y – intercept
b = Slope of the line
Y = Dependent variable
x= Independent variable

## Simple Linear Regression

Regression analysis is most often used for prediction. The goal in regression analysis is to create a mathematical model that can be used to predict the values of a dependent variable based upon the values of an independent variable. In other words, we use the model to predict the value of Y when we know the value of X. (The dependent variable is the one to be predicted). Correlation analysis is often used with regression analysis because correlation analysis is used to measure the strength of association between the two variables X and Y.

In regression analysis involving one independent variable and one dependent variable the values are frequently plotted in two dimensions as a scatter plot. The scatter plot allows us to visually inspect the data prior to running a regression analysis. Often this step allows us to see if the relationship between the two variables is increasing or decreasing and gives only a rough idea of the relationship. The simplest relationship between two variables is a straight-line or linear relationship. Of course the data may well be curvilinear and in that case we would have to use a different model to describe the relationship. Simple linear regression analysis finds the straight line that best fits the data.

## Fitting a Line to Data

Fitting a Line to data means drawing a line that comes as close as possible to the points. (Note that, no straight line passes exactly through all of the points). The overall pattern can be described by drawing a straight line through the points.

Example:

The data in the table below were obtained by measuring the heights of 161 children from a village each month from 18 to 29 months of age.

Table: Mean height of children

| Age in months (x) | Height in centimeters (y) |
|---|---|
| 18 | 76.1 |
| 19 | 77 |
| 20 | 78.1 |
| 21 | 78.2 |
| 22 | 78.8 |
| 23 | 79.7 |
| 24 | 79.9 |
| 25 | 81.1 |
| 26 | 81.2 |
| 27 | 81.8 |
| 28 | 82.8 |
| 29 | 83.5 |

Figure below is a scatter plot of the data in the above table.

Age is the explanatory variable, which is plotted on the x axis. Mean height (in cm) is the response variable.

We can see on the plot a strong positive linear association with no outliers. The correlation is r=0.994, close to the r = 1 of points that lie exactly on a line.

If we draw a line through the points, it will describe these data very well. This line is called the regression line and the process of doing so is called 'Fitting a line'. This is done in figure below.

Let y is a response variable and x is an explanatory variable.

A straight line relating y to x has an equation of the form  y = a + bx.

In this equation, b is the slope, the amount by which y changes when x increases by one unit.

The number a is the intercept, the value of y when x = 0

The straight line describing the data has the form

height = a + (b × age).

In Figure below  the regression line has been drawn with the following equation

height = 64.93 + (0.635 × age).

⇒The figure above shows that this line fits the data well.

The slope b = 0.635 tells us that the height of children increases by about 0.6 cm for each month of age.

The slope b of a line y = a + bx is the rate of change in the response y as the explanatory variable x changes.

The slope of a regression line is an important numerical description of the relationship between the two variables.

**Regression for prediction**

We use the regression equation for prediction of the value of a variable,

Suppose we have a sample of size 'n' and it has two sets of measures, denoted by x and y. We can predict the values of 'y' given the values of 'x' by using the equation, called the regression equation given below.

y* = a + bx

where the coefficients a and b are given by

$$a = \frac{\sum y - b \sum x}{n}$$

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

In the regression equation the symbol y* refers to the predicted value of y from a given value of x from the regression equation.

Let us see with the aid of an example how regressions used for prediction.

Example:

Scores made by students in a statistics class in the mid - term and final examination are given here. Develop a regression equation which may be used to predict final examination scores from the mid – term score.

| STUDENT | MID TERM | FINAL |
|---------|----------|-------|
| 1 | 98 | 90 |
| 2 | 66 | 74 |
| 3 | 100 | 98 |
| 4 | 96 | 88 |
| 5 | 88 | 80 |
| 6 | 45 | 62 |
| 7 | 76 | 78 |
| 8 | 60 | 74 |
| 9 | 74 | 86 |
| 10 | 82 | 80 |

Solution:

We want to predict the final exam scores from the mid term scores. So let us designate 'y' for the final exam scores and 'x' for the mid term exam scores. We open the following table for the calculations.

| STUDENT | X | Y | $X^2$ | XY |
|---------|-----|-----|-------|------|
| 1 | 98 | 90 | 9604 | 8820 |
| 2 | 66 | 74 | 4356 | 4884 |
| 3 | 100 | 98 | 10000 | 9800 |
| 4 | 96 | 88 | 9216 | 8448 |
| 5 | 88 | 80 | 7744 | 7040 |
| 6 | 45 | 62 | 2025 | 2790 |
| 7 | 76 | 78 | 5776 | 5928 |
| 8 | 60 | 74 | 3600 | 4440 |
| 9 | 74 | 86 | 5476 | 6364 |
| 10 | 82 | 80 | 6724 | 6560 |
| | 785 | 810 | 64521 | 65074 |

First find $b$ and then find $a$ and substitute in the equation.

$$b = \frac{n \sum xy - \left( \sum x \right)\left( \sum y \right)}{n \left( \sum x^2 \right) - \left( \sum x \right)^2} = \frac{10(65074) - (785)(810)}{10(64521) - (785)^2}$$

$$¿ \frac{650740 - 635850}{645210 - 616225} = \frac{14890}{28985} = 0.514$$

$$a = \frac{\sum y - b \sum x}{n} = \frac{810 - (0.514)(785)}{10} = \frac{810 - 403.49}{10} = \frac{406.51}{10} = 40.651$$

So a = 40.651 and b =0.514

Substitute in the equation for regression line y* = a + bx

y* = 40.651 + (0.514)x

Now we can use this for making predictions.

We can use this to find the projected or estimated final scores of the students.

For example, for the midterm score of 50 the projected final score is

y* = 40.651 + (0.514) 50 = 40.651 + 25.70 = 66.351, which is a quite a good estimation.

To give another example, consider the midterm score of 70. Then the projected final score is

y* = 40.651 + (0.514) 70 = 40.651 + 35.98= 76.631, which is again a very good estimation.

**Applications (uses) of regression analysis**

1. **Predicting the Future** :The most common use of regression in business is to predict events that have yet to occur. Demand analysis, for example, predicts how many units consumers will purchase. Many other key parameters other than demand are dependent variables in regression models, however. Predicting the number of shoppers who will pass in front of a particular billboard or the number of viewers who will watch the Champions Trophy Cricket may help management assess what to pay for an advertisement.

2. Insurance companies heavily rely on regression analysis to estimate, for example, how many policy holders will be involved in accidents or be victims of theft,.

**3. Optimization**: Another key use of regression models is the optimization of business processes. A factory manager might, for example, build a model to understand the relationship between oven temperature and the shelf life of the cookies baked in those ovens. A company operating a call center may wish to know the relationship between wait times of callers and number of complaints.

4. A fundamental driver of enhanced productivity in business and rapid economic advancement around the globe during the 20th century was the frequent use of statistical tools in manufacturing as well as service industries. Today, managers considers regression an indispensable tool.

**Limitations of Regression Analysis:**          There are three main limitations:

1. Parameter Instability - This is the tendency for relationships between variables to change over time due to changes in the economy or the markets, among other uncertainties. If a mutual fund

produced a return history in a market where technology was a leadership sector, the model may not work when foreign and small-cap markets are leaders.

2. Public Dissemination of the Relationship - In an efficient market, this can limit the effectiveness of that relationship in future periods. For example, the discovery that low price-to-book value stocks outperform high price-to-book value means that these stocks can be bid higher, and value-based investment approaches will not retain the same relationship as in the past.

3. Violation of Regression Relationships - Earlier we summarized the six classic assumptions of a linear regression. In the real world these assumptions are often unrealistic - e.g. assuming the independent variable X is not random.

**Correlation or Regression**

Correlation and regression analysis are related in the sense that both deal with relationships among variables. Whether to use Correlation or Regression in an analysis is often confusing for researchers.

In regression the emphasis is on predicting one variable from the other, in correlation the emphasis is on the degree to which a linear model may describe the relationship between two variables. In regression the interest is directional, one variable is predicted and the other is the predictor; in correlation the interest is non-directional, the relationship is the critical aspect.

Correlation makes no a priori assumption as to whether one variable is dependent on the other(s) and is not concerned with the relationship between variables; instead it gives an estimate as to the degree of association between the variables. In fact, correlation analysis tests for interdependence of the variables.

As regression attempts to describe the dependence of a variable on one (or more) explanatory variables; it implicitly assumes that there is a one-way causal effect from the explanatory variable(s) to the response variable, regardless of whether the path of effect is direct or indirect. There are advanced regression methods that allow a non-dependence based relationship to be described (eg. Principal Components Analysis or PCA) and these will be touched on later.

The best way to appreciate this difference is by example.

Take for instance samples of the leg length and skull size from a population of elephants. It would be reasonable to suggest that these two variables are associated in some way, as elephants with short legs tend to have small heads and elephants with long legs tend to have big heads. We may, therefore, formally demonstrate an association exists by performing a correlation analysis. However, would regression be an appropriate tool to describe a relationship between head size and leg length? Does an increase in skull size cause an increase in leg length? Does a decrease in leg length cause the skull to shrink? As you can see, it is meaningless to apply a

causal regression analysis to these variables as they are interdependent and one is not wholly dependent on the other, but more likely some other factor that affects them both (eg. food supply, genetic makeup).

Consider two variables: crop yield and temperature. These are measured independently, one by the weather station thermometer and the other by Farmer Giles' scales. While correlation anaylsis would show a high degree of association between these two variables, regression anaylsis would be able to demonstrate the dependence of crop yield on temperature. However, careless use of regression analysis could also demonstrate that temperature is dependent on crop yield: this would suggest that if you grow really big crops you will be guaranteed a hot summer.

Thus, neither regression nor correlation analyses can be interpreted as establishing cause-and-effect relationships. They can indicate only how or to what extent variables are associated with each other. The correlation coefficient measures only the degree of linear association between two variables. Any conclusions about a cause-and-effect relationship must be based on the judgment of the analyst.

## **Uses of Correlation and Regression**

There are three main uses for correlation and regression.

1. One is to test hypotheses about cause-and-effect relationships. In this case, the experimenter determines the values of the X-variable and sees whether variation in X causes variation in Y. For example, giving people different amounts of a drug and measuring their blood pressure.

2. The second main use for correlation and regression is to see whether two variables are associated, without necessarily inferring a cause-and-effect relationship. In this case, neither variable is determined by the experimenter; both are naturally variable. If an association is found, the inference is that variation in X may cause variation in Y, or variation in Y may cause variation in X, or variation in some other factor may affect both X and Y.

3.The third common use of linear regression is estimating the value of one variable corresponding to a particular value of the other variable.

***********************