



Param Pujya Dr. Babasaheb Ambedkar Smarak Samiti's  
**Dr. Ambedkar Institute of Management Studies & Research**

Deeksha Bhoomi, Nagpur - 440010 (Maharashtra State) INDIA

NAAC Accredited with 'A' Grade

Tel: +91 712 6521204, 6521203, 6501379

Email: info@daimsr.in

# **UNIT II**

# **MEASURES OF DISPERSION**

# Programme Educational Objectives

*Our program will create graduates who:*

- 1. Will be recognized as a creative and an enterprising team leader.*
- 2. Will be a flexible, adaptable and an ethical individual.*
- 3. Will have a holistic approach to problem solving in the dynamic business environment.*

# Research Methodology & Quantitative Techniques

## Course Outcomes

- CO1-Given a managerial problem and associated frequency distribution data, the student manager will be able to apply descriptive and inferential statistics to facilitate quick and rationale managerial decision making.
- CO2-Given the data for two or more variables, the student manager will be able to estimate the strength of the relationship between two variables using 'Karl Pearson' and 'Spearman's Rank' correlation coefficient.
- CO3-Given the data for two or more variables, the student manager will be able to predict / forecast using as moving averages, regression and time series analysis.

CO4-Given a managerial problem, the student manager will be able to formulate it as 'research problem' and also will be able to suggest suitable research methodology to identify workable solutions.

CO5-Given a business Problem/situation, the student manager will be able to develop methods and instruments (questionnaire/ interview schedule) for collection and measurement of qualitative as well as quantitative data using primary and secondary sources from a given sampling framework.

CO6-Given the sample statistics, the student manager will be able to apply Z, t and Chi-square tests to accept or reject the stated hypotheses for making sound decisions.

# Learning Objectives

To understand the limitation of averages

To measure the extent to which the items vary from central value

To compare the data set in terms of variability, consistency by using absolute and relative measures

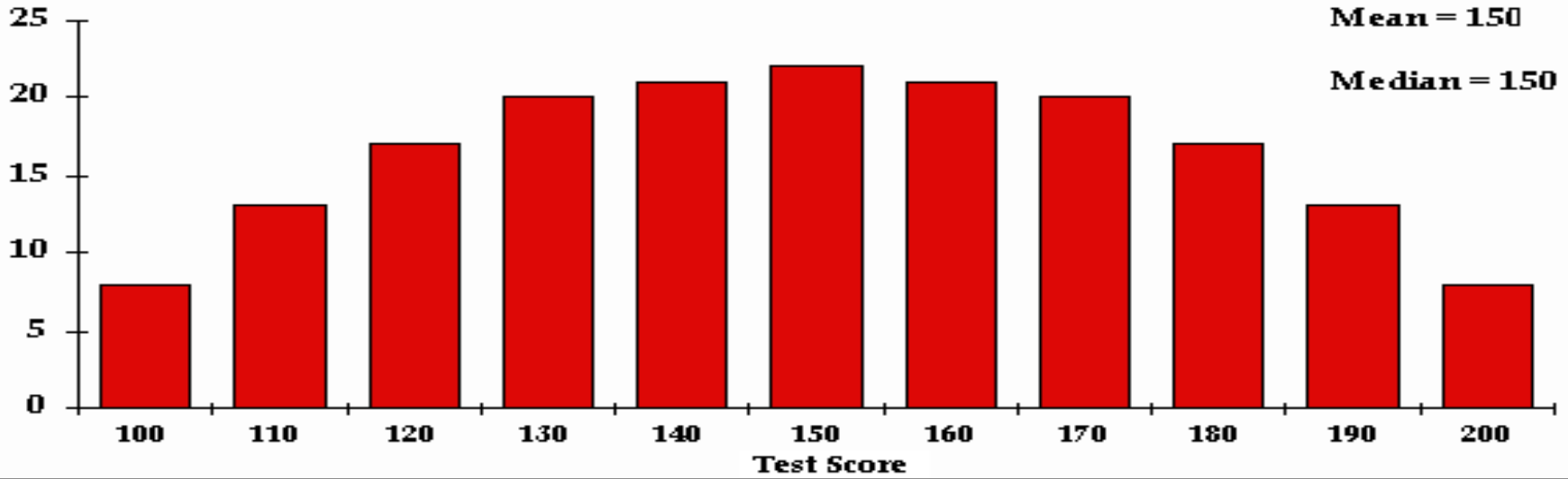
**Flat Distribution**

No. of People

**N = 180**

**Mean = 150**

**Median = 150**



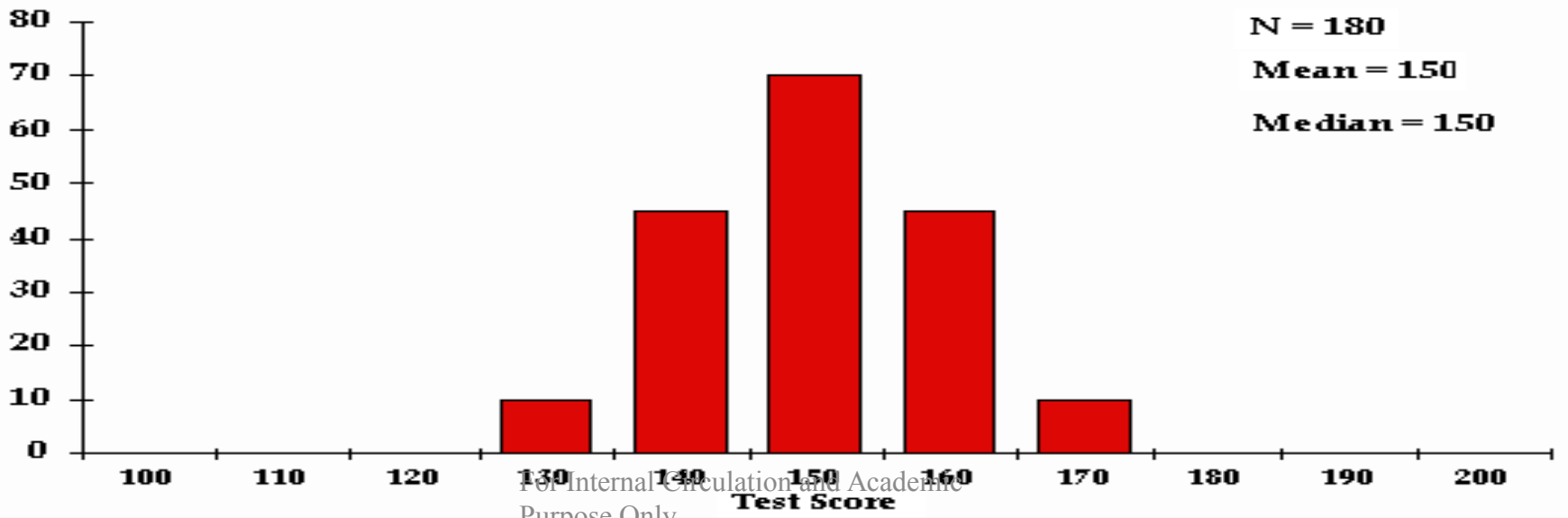
**Narrow Distribution**

No. of People

**N = 180**

**Mean = 150**

**Median = 150**



# WHAT IS DISPERSION?

*Dispersion or spread is the degree of the scatter or variation of the variable about a central value.*

*Dispersion is the measure of the variations of the item*

*The degree to which numerical data tend to spread about an average value is called the variation or dispersion of the data*

*Measures of variability are usually used to indicate how tightly bunched the sample values are around the mean*

# PURPOSE OF MEASURING DISPERSION.

- To judge the reliability of measures of central tendency*
- To make a comparative study of the variability of two series*
- To identify the causes of variability with a view to control it*
- To serve as a basis for further statistical analysis*



# RANGE.

*Absolute Measures of Dispersion: The measures of dispersion which are expressed in terms of **original units** of a data are termed as Absolute Measures.*

*Relative Measures of Dispersion: Relative measures of dispersion, are also known as coefficients of dispersion, are obtained as **ratios or percentages**. These are pure numbers **independent of the units** of measurement and used to compare two or more sets of data values.*

*Absolute Measures • Range • Quartile Deviation • Mean Deviation • Standard Deviation*

*Relative Measure • Co-efficient of Range • Co-efficient of Quartile Deviation • Co-efficient of mean Deviation • co-efficient of Variation.*

# RANGE.

*Range is the simplest possible measure of dispersion. It is the difference between the values of the extreme items of a series.*

$$\text{Range} = L - S$$

$$\text{Co-efficient of Range} = (L-S)/(L+S)$$

1. The profits of a company for the last 8 years are given below. Calculate the Range and its Co-efficient

	<b>197</b>	<b>197</b>	<b>197</b>	<b>197</b>	<b>197</b>	<b>198</b>	<b>198</b>	<b>198</b>
<b>Year</b>	<b>5</b>	<b>6</b>	<b>7</b>	<b>8</b>	<b>9</b>	<b>0</b>	<b>1</b>	<b>2</b>
<b>Profits (in '000)</b>								<b>230</b>

**(Answer: Range = 200, and its Co-efficient = 0.77)**

2. Calculate Co-efficient of Range from the following data

<b>Weekly Wages (Rs.)</b>	<b>50-60</b>	<b>60-70</b>	<b>70-80</b>	<b>80-90</b>	<b>90-100</b>	<b>100-110</b>	<b>110-120</b>
<b>No. of Laborers</b>	<b>50</b>	<b>45</b>	<b>45</b>	<b>40</b>	<b>35</b>	<b>30</b>	<b>30</b>

*First Method (By Taking Lower limits of first and upper limit of last interval)*

$$\begin{aligned}\text{Coefficient of Range} &= (L-S)/(L+S) \\ &= (120-50)/(120+50) \\ &= 70/170 = 0.41\end{aligned}$$

*Second Method (By Taking mid value of first and last interval)*

$$\begin{aligned}\text{Coefficient of Range} &= (L-S)/(L+S) \\ &= (115-55)/(115+55) \\ &= 60/170 = 0.35\end{aligned}$$

*It should be noted that in the calculation of Range only the values of the variable are taken into account and the frequencies are completely ignored.*

# INTER QUARTILE RANGE.

$Q1 = \text{the value of } \left(\frac{N + 1}{4}\right)^{\text{th}}$

$Q3 = \text{the value of } 3 \left(\frac{N + 1}{4}\right)^{\text{th}}$

*Quartile Deviation*

*Coefficient of Quartile Deviation*  $= \frac{Q3 - Q1}{2}$

$$\frac{Q3 - Q1}{Q3 + Q1}$$

*Find the Quartile Deviation and its Co-efficient from the following data, relating to the weekly of seven laborers*

<b>Weekly Wages (Rs.)</b>	<b>50</b>	<b>70</b>	<b>80</b>	<b>60</b>	<b>65</b>	<b>40</b>	<b>90</b>
---------------------------	-----------	-----------	-----------	-----------	-----------	-----------	-----------

**Answer: Quartile Deviation = Rs. 15, and its Co-efficient = 0.23**

*Calculate Quartile Deviation and its Co-efficient from the following data*

<b>Weight in Pounds</b>	<b>120</b>	<b>122</b>	<b>124</b>	<b>126</b>	<b>130</b>	<b>140</b>	<b>150</b>	<b>160</b>
<b>No. of Students</b>	<b>1</b>	<b>3</b>	<b>5</b>	<b>7</b>	<b>10</b>	<b>3</b>	<b>1</b>	<b>1</b>

**(Answer: Q1=124, Q3=130, Co-efficient of Q.D. = 0.0236)**

# QD in Continuous series.

Step 1: Calculate the class in which  $Q_1$  lies using formula  $N/4$ .

Step 2: Calculate the class in which  $Q_3$  lies using formula  $3(N/4)$

Step 3: 
$$Q_1 = l_1 + \frac{l_2 - l_1}{f_1} (q_1 - c)$$

Step 4: 
$$Q_3 = l_1 + \frac{l_2 - l_1}{f_1} (q_3 - c)$$

Step 5: Quartile Deviation = 
$$\frac{Q_3 - Q_1}{2}$$

Step 6: Coefficient Of QD 
$$\frac{Q_3 - Q_1}{Q_3 + Q_1}$$

*Calculate semi-inter quartile range and it's co-efficient from the following data*

<b>Marks</b>	<b>0-10</b>	<b>10-20</b>	<b>20-30</b>	<b>30-40</b>	<b>40-50</b>	<b>50-60</b>	<b>60-70</b>	<b>70-80</b>	<b>80-90</b>
<b>No. of Students</b>	<b>11</b>	<b>18</b>	<b>25</b>	<b>28</b>	<b>30</b>	<b>33</b>	<b>22</b>	<b>15</b>	<b>22</b>

**(Answer: Q.D. = 17.42 marks, co-efficient of Q.D. = 0.37)**

*Calculate Quartile Deviation and its relative measure*

<b>Variable</b>	<b>Frequency</b>	<b>Variable</b>	<b>Frequency</b>
<b>20-29</b>	<b>306</b>	<b>50-59</b>	<b>96</b>
<b>30-39</b>	<b>182</b>	<b>60-69</b>	<b>42</b>
<b>40-49</b>	<b>144</b>	<b>70-79</b>	<b>34</b>

**(Answer: Q.D. = 17.42 marks, co-efficient of Q.D. = 0.37)**



*Estimate an appropriate measure of dispersion of the following data*

<b>Income (Rs.)</b>	<b>No. of Persons</b>	<b>Income (Rs.)</b>	<b>No. of Persons</b>
<b>Less than 50</b>	<b>54</b>	<b>110-130</b>	<b>230</b>
<b>50-70</b>	<b>100</b>	<b>130-150</b>	<b>125</b>
<b>70-90</b>	<b>140</b>	<b>Above 150</b>	<b>51</b>
<b>90-110</b>	<b>300</b>		

**(Answer: Q.D. = 19.9 Rs.)**

# MEAN DEVIATION.

*Mean deviation of a series is the arithmetic averages of the deviations of various items from a measure of central tendency (mean, median or mode).*

$$\delta\bar{X} = \frac{\sum |d\bar{X}|}{N}$$

$$\delta M = \frac{\sum |dM|}{N}$$

$$\delta Z = \frac{\sum |dZ|}{N}$$

*The following are the marks obtained by a batch of 9 students in a certain test. Calculate the mean deviation from mean and median.*

<b>SR. No.</b>	<b>Marks (out of 100)</b>	<b>SR. No.</b>	<b>Marks (out of 100)</b>
<b>1</b>	<b>68</b>	<b>6</b>	<b>38</b>
<b>2</b>	<b>49</b>	<b>7</b>	<b>59</b>
<b>3</b>	<b>32</b>	<b>8</b>	<b>66</b>
<b>4</b>	<b>21</b>	<b>9</b>	<b>41</b>
<b>5</b>	<b>54</b>		

$$\text{Mean} = 428 / 9 = 47.55$$

$$\text{Median} = 49$$

<b>X</b>	<b>X - MEAN</b>	<b>X - MEAN</b>
<b>68</b>	<b>20.45</b>	<b>20.45</b>
<b>49</b>	<b>1.45</b>	<b>1.45</b>
<b>32</b>	<b>-15.55</b>	<b>15.55</b>
<b>21</b>	<b>-26.55</b>	<b>26.55</b>
<b>54</b>	<b>6.45</b>	<b>6.45</b>
<b>38</b>	<b>-9.55</b>	<b>9.55</b>
<b>59</b>	<b>11.45</b>	<b>11.45</b>
<b>66</b>	<b>18.45</b>	<b>18.45</b>
<b>41</b>	<b>-6.55</b>	<b>6.55</b>
		<b>116.45</b>

$$\delta\bar{X} = \frac{\sum |d\bar{X}|}{N}$$

$$\delta\bar{X} = \frac{116.45}{9}$$

$$\delta\bar{X} = 12.93$$

<b>X</b>	<b>X - Median</b>	<b>X - Median</b>
<b>68</b>		
<b>49</b>	<b>19</b>	<b>19</b>
<b>32</b>	<b>0</b>	<b>0</b>
<b>21</b>	<b>-17</b>	<b>17</b>
<b>54</b>	<b>-28</b>	<b>28</b>
<b>38</b>	<b>5</b>	<b>5</b>
<b>59</b>	<b>-11</b>	<b>11</b>
<b>66</b>	<b>10</b>	<b>10</b>
<b>41</b>	<b>17</b>	<b>17</b>
	<b>-8</b>	<b>8</b>
		<b>115</b>

$$\delta M = \frac{\sum |dM|}{N}$$

$$\delta M = \frac{115}{9}$$

$$\delta M = 12.77$$

*Calculate Mean Deviation (from arithmetic Average) for the following values. Also, calculate its Coefficient. 4800, 4600, 4400, 4200, 4000.*

**(Answer: Mean Deviation = 240,  
Co-efficient of Mean Deviation = 0.54)**

*Calculate Mean Deviation (from median) for the following values.*

<b>No. of Accidents</b>	<b>Persons having said no. of Accidents</b>	<b>No. of Accidents</b>	<b>Persons having said no. of Accidents</b>
<b>0</b>	<b>15</b>	<b>7</b>	<b>2</b>
<b>1</b>	<b>16</b>	<b>8</b>	<b>1</b>
<b>2</b>	<b>21</b>	<b>9</b>	<b>2</b>
<b>3</b>	<b>10</b>	<b>10</b>	<b>2</b>
<b>4</b>	<b>17</b>	<b>11</b>	<b>0</b>
<b>5</b>	<b>8</b>	<b>12</b>	<b>2</b>
<b>6</b>	<b>4</b>		

# STANDARD DEVIATION.

*Standard Deviation is the square root of the arithmetic average of the squares of the deviations measured from the mean.*

$$\sigma = \sqrt{\frac{\sum d^2}{N}}$$

*Calculate the standard deviation of the heights of 10 students given below:*

*Height: 160, 160, 161, 162, 163, 163, 163, 164, 164, 170  
(in cms)*

<b>X</b>	<b>d = X - 163</b>	<b>d<sup>2</sup></b>
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
<b>ΣX = 1630</b>		<b>Σ d<sup>2</sup> = 74</b>

$$\sigma = \sqrt{\frac{\sum d^2}{N}}$$

$$\sqrt{\frac{74}{10}}$$

$$\sqrt{7.4}$$

$$= 2.72$$

**cms**



# STANDARD DEVIATION – Method 2.

*In this method we need not calculate the deviations. The formula used is as follows:*

$$\sigma = \sqrt{\frac{\sum X^2 - (\sum X)^2 / N}{N}}$$

*Calculate the standard deviation in the previous example using the above method.*

<b>X</b>	<b>X<sup>2</sup></b>
160	25600
160	25600
161	25921
162	26244
163	26569
163	26569
163	26569
164	26896
164	26896
170	28900
<b>ΣX =</b> <b>1630</b>	<b>ΣX<sup>2</sup> = 265764</b>

$$\sigma = \sqrt{\frac{\sum X^2 - (\sum X)^2 / N}{N}}$$

$$\sqrt{\frac{265764 - (1630)^2 / 10}{10}}$$

$$= 2.72 \text{ cms}$$

# STANDARD DEVIATION

## Assumed Mean Method

*In this method the formula used is as follows:*

$$\sigma = \sqrt{\frac{\sum dX^2}{N} - \left(\frac{\sum dX}{N}\right)^2}$$

$$\sigma = \sqrt{\frac{\sum dX^2 - N(\bar{X} - A)^2}{N}}$$

$$\sigma = \sqrt{\frac{\sum dX^2}{N} - (\bar{X} - A)^2}$$

*Calculate the standard deviation in the previous example using the above method.*

# STANDARD DEVIATION - DISCRETE SERIES

*In case of discrete series the formula used is as follows:*

$$\sigma = \sqrt{\frac{\sum f d^2}{N}}$$

**Calculate the standard deviation from the following data:**

<b>Size of item</b>	<b>Frequency</b>	<b>Size of item</b>	<b>Frequency</b>
<b>6</b>	<b>3</b>	<b>10</b>	<b>8</b>
<b>7</b>	<b>6</b>	<b>11</b>	<b>5</b>
<b>8</b>	<b>9</b>	<b>12</b>	<b>4</b>
<b>9</b>	<b>13</b>		

<b>X</b>	<b>f</b>	<b>fX</b>	<b>d = X - 9</b>	<b>d<sup>2</sup></b>	<b>fd<sup>2</sup></b>
6	3	18	-3	9	27
7	6	42	-2	4	24
8	9	72	-1	1	9
9	13	117	0	0	0
10	8	80	1	1	8
11	5	55	2	4	20
12	4	48	3	9	36
	$\Sigma f = 48$	$\Sigma fX = 432$			$\Sigma fd^2 = 124$

# Practice Example 1

*Calculate standard deviation for the following distribution:*

<b>Values</b>	<b>10</b>	<b>20</b>	<b>30</b>	<b>40</b>	<b>50</b>	<b>60</b>	<b>70</b>
<b>Frequency y</b>	<b>1</b>	<b>5</b>	<b>12</b>	<b>22</b>	<b>17</b>	<b>9</b>	<b>4</b>

$$\sigma = \sqrt{\frac{\sum fdX^2}{N} - \left(\frac{fdX}{N}\right)^2}$$

**(Answer: S.D. = 13.26)**

# Practice Example 2

*The following table gives the number of finished articles turned out per day by different number of workers in a factory. Find the standard deviation of the daily output of finished articles.*

<b>No. of Articles</b>	<b>No. of Workers</b>	<b>No. of Articles</b>	<b>No. of Workers</b>
<b>18</b>	<b>3</b>	<b>23</b>	<b>17</b>
<b>19</b>	<b>7</b>	<b>24</b>	<b>13</b>
<b>20</b>	<b>11</b>	<b>25</b>	<b>8</b>
<b>21</b>	<b>14</b>	<b>26</b>	<b>5</b>
<b>22</b>	<b>18</b>	<b>27</b>	<b>4</b>

**(Answer: S.D. = 2.2 articles)**

# CONTINUOUS SERIES – Example

*Calculate the standard deviation for the following table giving the age distribution of 542 members of the LOK SABHA.*

Age	No. of Members	Age	No. of Members
<b>20-30</b>	<b>3</b>	<b>60-70</b>	<b>140</b>
<b>30-40</b>	<b>61</b>	<b>70-80</b>	<b>51</b>
<b>40-50</b>	<b>132</b>	<b>80-90</b>	<b>2</b>
<b>50-60</b>	<b>153</b>		

**(Answer: S.D. = 11.9 years)**



Let us ASSUME the MEAN AGE of the members as 45 years. Then we can use the following formula to calculate standard deviation:

$$\sigma = \sqrt{\frac{\sum fdx^2}{N} - \left(\frac{fdx}{N}\right)^2}$$

Age	X	f	dx	dx <sup>2</sup>	fdx	fdx <sup>2</sup>
20 – 30	25	3	-20	400	-60	1200
30 – 40	35	61	-10	100	-610	6100
40 – 50	45	132	0	0	0	0
50 – 60	55	153	10	100	1530	15300
60 – 70	65	140	20	400	2800	56000
70 – 80	75	51	30	900	1530	45900
80 – 90	85	2	40	1600	80	3200
		<b>542</b>			<b>5270</b>	<b>127700</b>

# Practice Example 1

*The following data relate to the age of a group of Govt. employees Calculate the standard deviation:*

<b>Age</b>	<b>50- 55</b>	<b>45- 50</b>	<b>40- 45</b>	<b>35- 40</b>	<b>30- 35</b>	<b>25- 30</b>	<b>20- 25</b>
<b>No. of Employee s</b>	<b>25</b>	<b>30</b>	<b>40</b>	<b>45</b>	<b>80</b>	<b>110</b>	<b>170</b>

**(Answer: S.D. = 9 years approx)**

# Practice Example 2

*The following table relates to the profits and losses of 100 firms. Calculate the standard deviation of profits.*

<b>Profits</b>	<b>Number of Firms</b>
<b>5000 to 6000</b>	<b>8</b>
<b>4000 to 5000</b>	<b>12</b>
<b>3000 to 4000</b>	<b>30</b>
<b>2000 to 3000</b>	<b>10</b>
<b>1000 to 2000</b>	<b>5</b>
<b>0 to 1000</b>	<b>5</b>
<b>-1000 to 0</b>	<b>6</b>
<b>-2000 to -1000</b>	<b>8</b>
<b>-3000 to -2000</b>	<b>9</b>

**SD**  
**Rs. 2791**

# Practice Example 3

*Calculate the standard deviation of the following data:*

<b>Age Under (Years)</b>	<b>10</b>	<b>20</b>	<b>30</b>	<b>40</b>	<b>50</b>	<b>60</b>	<b>70</b>	<b>80</b>
<b>No. of Persons dying</b>	<b>15</b>	<b>30</b>	<b>53</b>	<b>75</b>	<b>100</b>	<b>110</b>	<b>115</b>	<b>125</b>

**(Answer: S.D. = 19.7 years)**

# CORRELATION

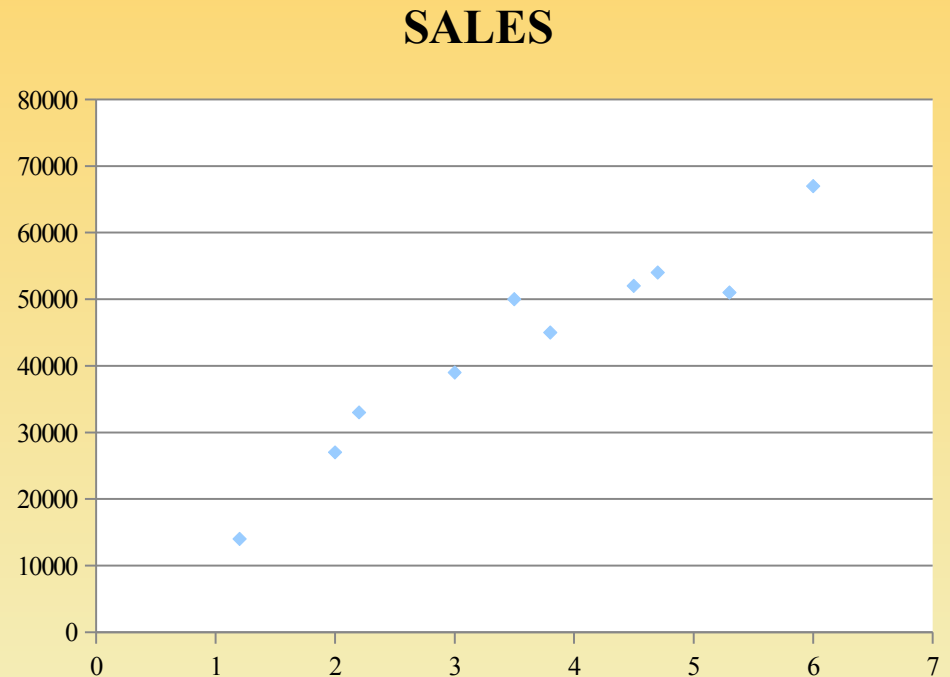
For Internal Circulation and Academic  
Purpose Only

# CORRELATION

- **Correlation** analysis is used to Measure strength of the association between two variables
- Only concerned with strength of the relationship
- No causal effect is implied

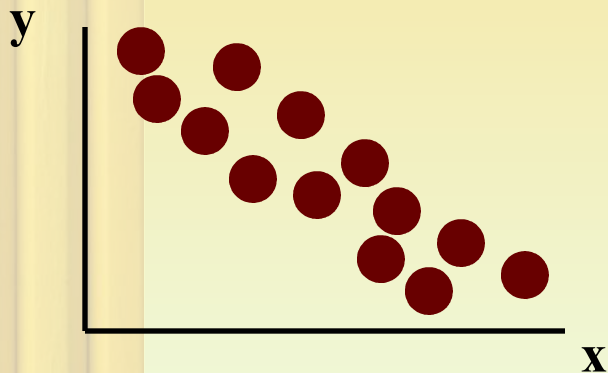
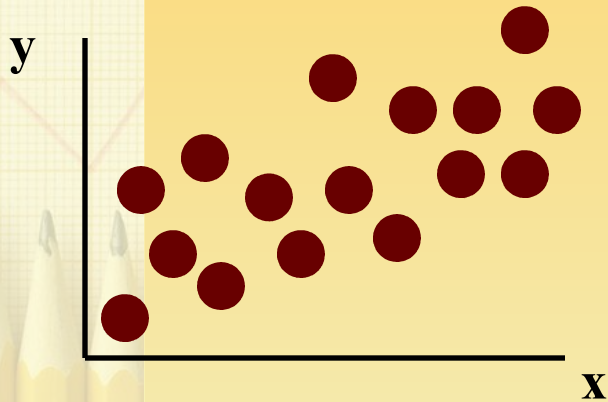
# EXAMPLE

YEAR	ADV. BUDGET	SALES
2001	2	27000
2002	3	39000
2003	4.5	52000
2004	6	67000
2005	2.2	33000
2006	1.2	14000
2007	3.5	50000
2008	4.7	54000
2009	5.3	51000
2010	3.8	45000

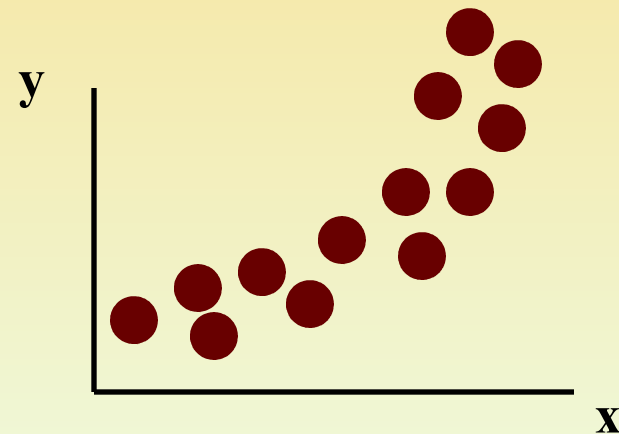
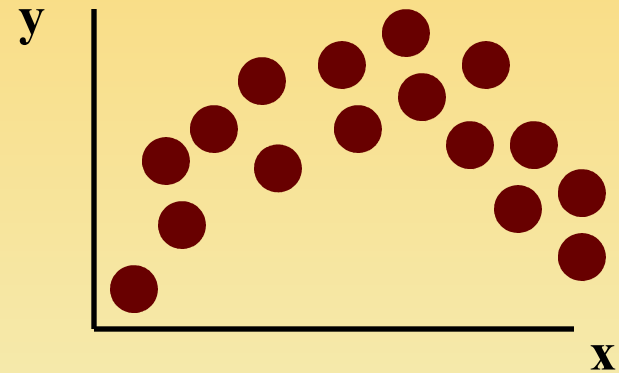


# Scatter Plots

Linear relationships



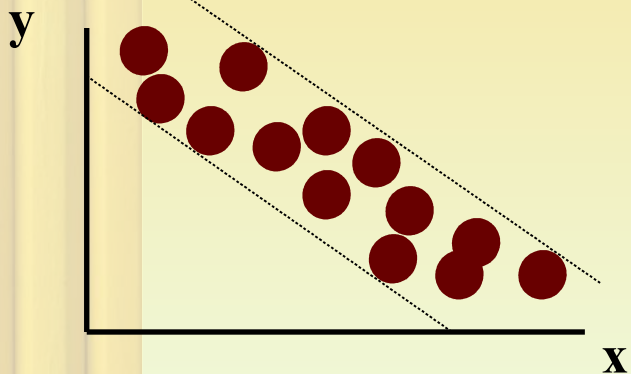
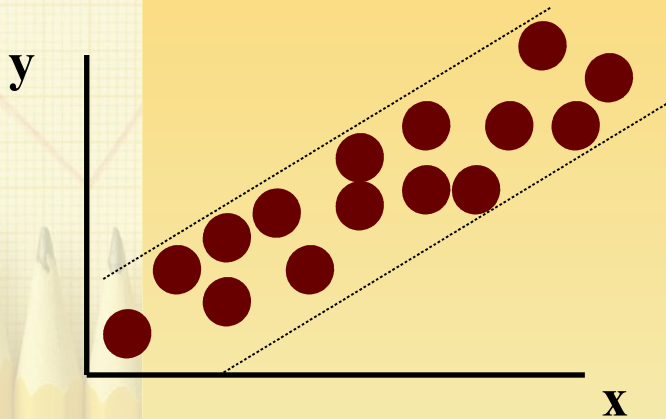
Curvilinear relationships



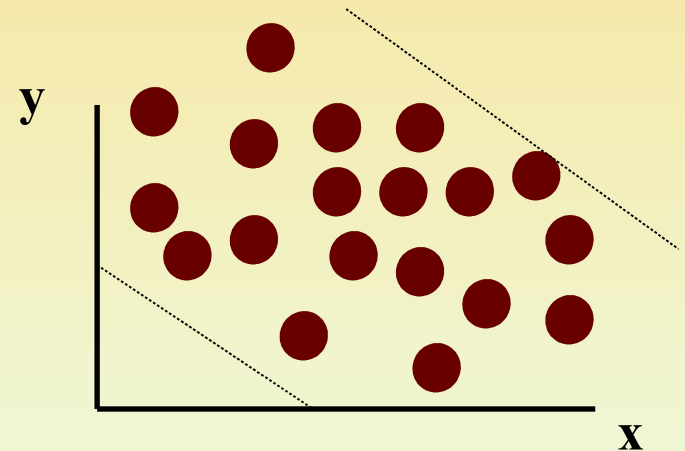
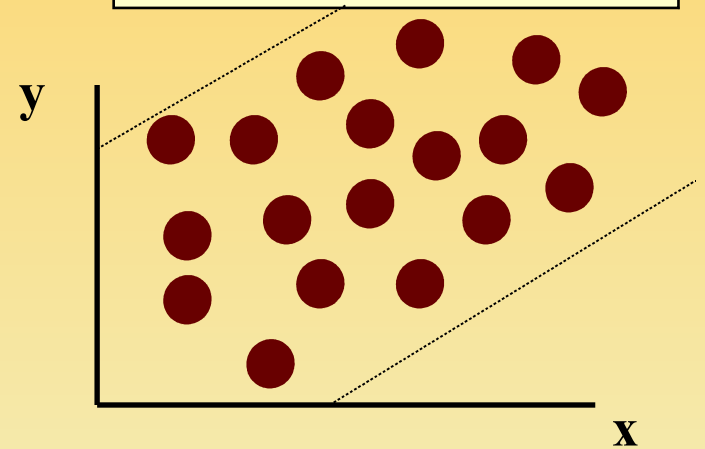


# Scatter Plots

Strong relationships

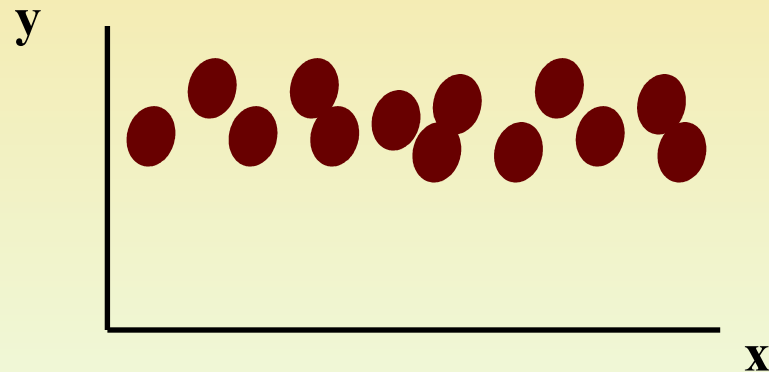
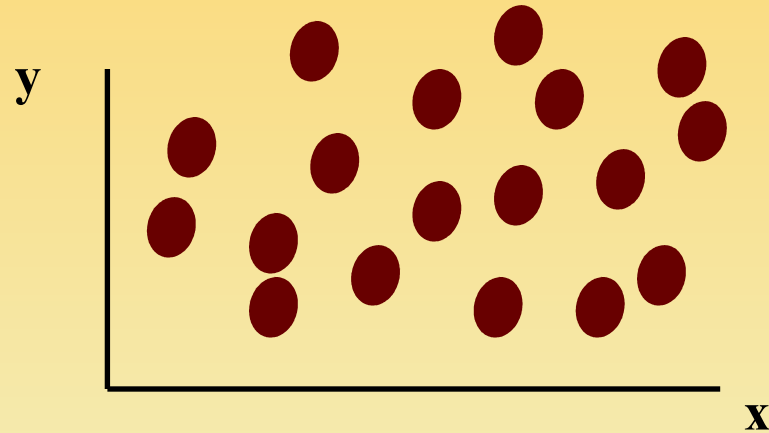


Weak relationships



# Scatter Plots

No relationship



For Internal Circulation and Academic Purpose Only

# Correlation Coefficient

- The population correlation coefficient  $\rho$  (**rho**) measures the strength of the association between the variables
- The sample correlation coefficient  $r$  is an estimate of  $\rho$  and is used to measure the strength of the linear relationship in the sample observations

# Features of $\rho$ and $r$

- **Range between -1 and 1**
- **The closer to -1, the stronger the negative linear relationship**
- **The closer to 1, the stronger the positive linear relationship**
- **The closer to 0, the weaker the linear relationship**

# Features of $\rho$ and $r$

Value of  $r$

Interpretation

**+ or - 1**

**Perfect correlation between variables**

**+ or - 0.9 to 0.99**

**Very high degree of correlation**

**+ or - 0.7 to 0.9**

**High degree of correlation**

**+ or - 0.5 to 0.7**

**Moderate degree of correlation**

**+ or - 0.25 to 0.5**

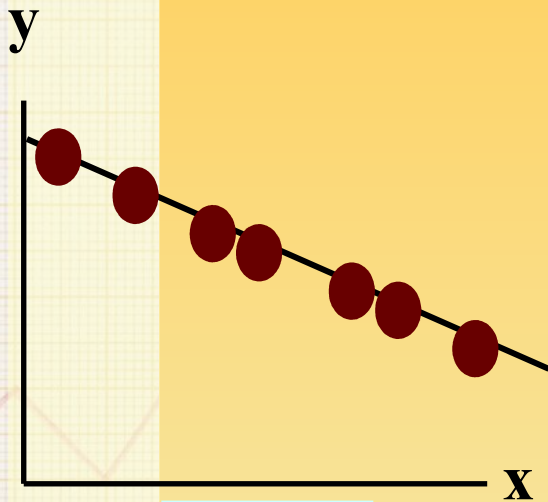
**Low degree of correlation**

**+ or - 0 to 0.25**

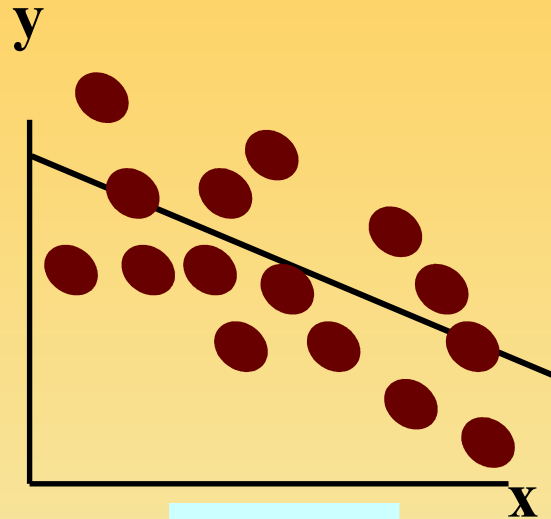
**Very low degree of correlation**

**0**

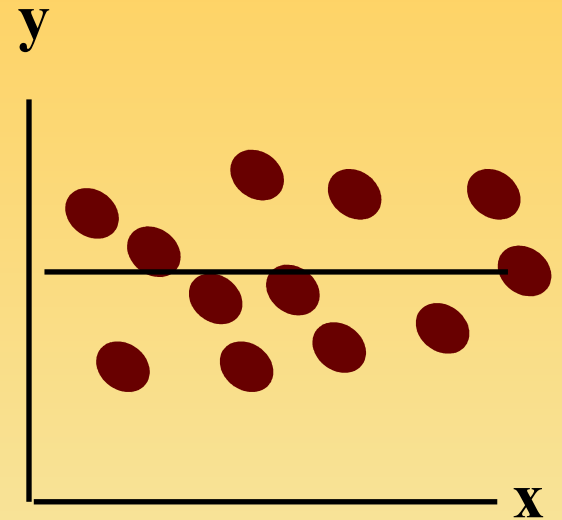
**No correlation**



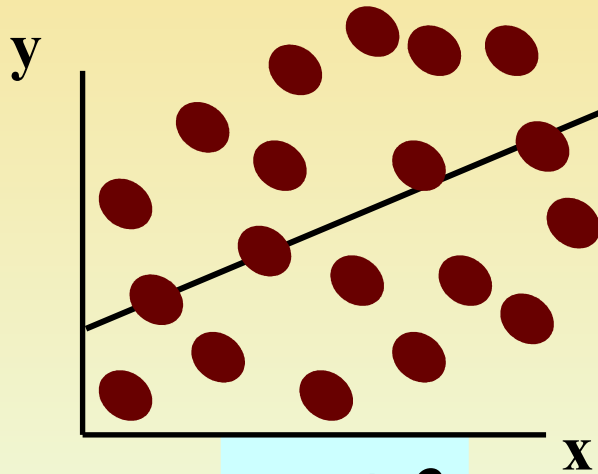
**$r = -1$**



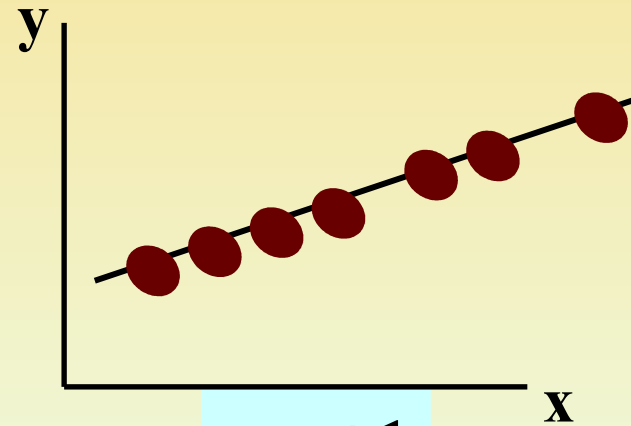
**$r = -0.6$**



**$r = 0$**



**$r = +0.3$**



**$r = +1$**

# Karl Pearson's coefficient of correlation - r

$$r = \frac{\sum xy}{\sqrt{\sum x^2 \sum y^2}}$$

where,

**r = Correlation Coefficient**

$$x = (X - \bar{X})$$

$$y = (Y - \bar{Y})$$

**From the data given below find out the Pearson's correlation coefficient for the given data and comment on the nature of correlation.**

<b>ROLL NUMBER OF STUDENT</b>	<b>MARKS IN MATHS (OUT OF 100)</b>	<b>MARKS IN STATISTICS (OUT OF 100)</b>
1	50	60
2	70	85
3	40	52
4	30	38
5	80	90



## Step 1

Calculating the Mean of first variable i.e.  $\bar{X}$

## Step 2

Calculating the Mean of second variable i.e.  $\bar{Y}$

## Step 3

Calculating the value of  $x = X - \bar{X}$

## Step 4

Calculating the value of  $y = Y - \bar{Y}$

## **Step 5**

**Calculating the  $\sum x y$**

## **Step 6**

**Calculating the  $\sum x^2$**

## **Step 7**

**Calculating the  $\sum y^2$**

## **Step 8**

**Putting the values in the formula to calculate  $r$**

<b>SR. NO.</b>	<b>X</b>	<b>Y</b>	<b><math>x = X -</math> <b>54</b></b>	<b><math>x^2</math></b>	<b><math>y = Y -</math> <b>65</b></b>	<b><math>y^2</math></b>	<b>XY</b>
<b>1</b>	<b>50</b>	<b>60</b>	<b>-4</b>	<b>16</b>	<b>-5</b>	<b>25</b>	<b>20</b>
<b>2</b>	<b>70</b>	<b>85</b>	<b>16</b>	<b>256</b>	<b>20</b>	<b>400</b>	<b>320</b>
<b>3</b>	<b>40</b>	<b>52</b>	<b>-14</b>	<b>196</b>	<b>-13</b>	<b>169</b>	<b>182</b>
<b>4</b>	<b>30</b>	<b>38</b>	<b>-24</b>	<b>576</b>	<b>-27</b>	<b>729</b>	<b>648</b>
<b>5</b>	<b>80</b>	<b>90</b>	<b>26</b>	<b>676</b>	<b>25</b>	<b>625</b>	<b>650</b>
<b>SUM (<math>\Sigma</math>)</b>	<b>270</b>	<b>325</b>		<b>172 0</b>		<b>194 8</b>	<b>182 0</b>
<b>MEAN</b>	<b>54</b>	<b>65</b>					

$$r = \frac{1820}{\sqrt{1720 \times 1948}}$$

$$r = + 0.99$$

### Comments:

- The Pearson's correlation coefficient is +0.99 which indicates that there is a very high degree of positive correlation between the marks obtained in Maths and marks obtained in Statistics.
- In other words a student who scores high marks in Maths also scores high marks in Statistics whereas a student who scores low marks in Maths also score low marks in Statistics.

## Example

*Making use of the below data calculate the coefficient of correlation.*

<b>CASE</b>	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>	<b>F</b>	<b>G</b>	<b>H</b>
<b>X1</b>	<b>10</b>	<b>6</b>	<b>9</b>	<b>10</b>	<b>12</b>	<b>13</b>	<b>11</b>	<b>9</b>
<b>X2</b>	<b>9</b>	<b>4</b>	<b>6</b>	<b>9</b>	<b>11</b>	<b>13</b>	<b>8</b>	<b>4</b>

**ANSWER:  $r = + 0.896$**

## Example

*Making use of the below data calculate the coefficient of correlation. (HINT: Since  $r$  is a pure number, changing the scale of series does not affect its values).*

<b>CASE</b>	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>	<b>F</b>	<b>G</b>
<b>X</b>	<b>10000</b>	<b>20000</b>	<b>30000</b>	<b>40000</b>	<b>50000</b>	<b>60000</b>	<b>70000</b>
<b>Y</b>	<b>30000</b>	<b>50000</b>	<b>60000</b>	<b>80000</b>	<b>100000</b>	<b>110000</b>	<b>130000</b>

**ANSWER:  $r = + 0.997$**

## Example

*From the data given below find out the Pearson's correlation coefficient for the given data and comment on the nature of correlation.*

Date	BSE SENSEX closing	Gold prices (per 10 gms)
5th Sep 2012	17250	29000
6th Sep 2012	16800	29900
7th Sep 2012	17100	29500
8th Sep 2012	17500	31000

Answer = -0.95  
For Internal Circulation and Academic Purpose Only

# Direct Method Karl Pearson's coefficient of correlation - r

In this method we need not calculate the deviations of items from mean.

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



## Example

*Making use of the below data calculate the coefficient of correlation by Direct Method.*

<b>CASE</b>	<b>A</b>	<b>B</b>	<b>C</b>	<b>D</b>	<b>E</b>	<b>F</b>	<b>G</b>	<b>H</b>
<b>X1</b>	<b>10</b>	<b>6</b>	<b>9</b>	<b>10</b>	<b>12</b>	<b>13</b>	<b>11</b>	<b>9</b>
<b>X2</b>	<b>9</b>	<b>4</b>	<b>6</b>	<b>9</b>	<b>11</b>	<b>13</b>	<b>8</b>	<b>4</b>

**ANSWER:  $r = + 0.896$**

# Shortcut or Assumed Mean Method for calculation of - r

When actual means are in fractions, the calculations of the correlation coefficient become complicated. So assume a mean and use the below formula:

$$r = \frac{N \sum dx dy - \{(\sum dx) \times (\sum dy)\}}{\sqrt{N \sum dx^2 - (\sum dx)^2} \sqrt{N \sum dy^2 - (\sum dy)^2}}$$

**Calculate the coefficient of correlation between X and Y from the following data. Assume 69 and 112 as the mean value for series X and Y respectively.**

<b>X</b>	<b>78</b>	<b>89</b>	<b>99</b>	<b>60</b>	<b>59</b>	<b>79</b>	<b>68</b>	<b>61</b>
<b>Y</b>	<b>12</b>	<b>13</b>	<b>15</b>	<b>11</b>	<b>10</b>	<b>13</b>	<b>12</b>	<b>10</b>
	<b>5</b>	<b>7</b>	<b>6</b>	<b>2</b>	<b>7</b>	<b>6</b>	<b>3</b>	<b>8</b>

<b>X</b>	<b>dx (X- 69)</b>	<b>dx<sup>2</sup></b>	<b>Y</b>	<b>dy (Y- 112)</b>	<b>dy<sup>2</sup></b>	<b>dx dy</b>
<b>78</b>	<b>9</b>	<b>81</b>	<b>125</b>	<b>13</b>	<b>169</b>	<b>117</b>
<b>89</b>	<b>20</b>	<b>400</b>	<b>137</b>	<b>25</b>	<b>625</b>	<b>500</b>
<b>99</b>	<b>30</b>	<b>900</b>	<b>156</b>	<b>44</b>	<b>1936</b>	<b>1320</b>
<b>60</b>	<b>-9</b>	<b>81</b>	<b>112</b>	<b>0</b>	<b>0</b>	<b>0</b>
<b>59</b>	<b>-10</b>	<b>100</b>	<b>107</b>	<b>-5</b>	<b>25</b>	<b>50</b>
<b>79</b>	<b>10</b>	<b>100</b>	<b>136</b>	<b>24</b>	<b>576</b>	<b>240</b>
<b>68</b>	<b>-1</b>	<b>1</b>	<b>123</b>	<b>11</b>	<b>121</b>	<b>-11</b>
<b>61</b>	<b>-8</b>	<b>64</b>	<b>108</b>	<b>-4</b>	<b>16</b>	<b>32</b>
<b><math>\Sigma X =</math> 593</b>	<b><math>\Sigma dx</math> =41</b>	<b><math>\Sigma dx^2 =</math> 1727</b>	<b><math>\Sigma Y =</math> 1004</b>	<b><math>\Sigma dy =</math> 108</b>	<b><math>\Sigma dy^2 =</math> 3468</b>	<b><math>\Sigma dx dy =</math> 2248</b>

$$r = \frac{N \sum dx dy - \{(\sum dx) x (\sum dy)\}}{\sqrt{N \sum dx^2 - (\sum dx)^2} \sqrt{N \sum dy^2 - (\sum dy)^2}}$$

$$r = \frac{8 \times 2248 - \{(41) \times (108)\}}{\sqrt{8 \times 1727 - (41)^2} \sqrt{8 \times 3468 - (108)^2}}$$

$$r = 0.97$$

**Calculate Karl Pearson's correlation coefficient from the advertisement cost and sales as per data given below:**

<b>Advertisement Cost</b>	<b>39</b>	<b>65</b>	<b>62</b>	<b>90</b>	<b>82</b>	<b>75</b>	<b>25</b>	<b>98</b>	<b>36</b>	<b>78</b>
<b>sales</b>	<b>47</b>	<b>53</b>	<b>58</b>	<b>86</b>	<b>62</b>	<b>68</b>	<b>60</b>	<b>91</b>	<b>51</b>	<b>84</b>

**(Answer  $r = +0.7804$ )**

**A computer while calculating correlation coefficient between variables X and Y from 25 pairs of observations obtained the following results:**

$$N = 25 \quad \Sigma X = 125 \quad \Sigma X^2 = 650$$

$$\Sigma Y = 100 \quad \Sigma Y^2 = 460 \quad \Sigma XY = 508$$

**It was, however, discovered at the time of checking that two pairs of observation were not correctly copied. They were taken as (6, 14) and (8, 6) while the correct values were (8, 12) and (6, 8). Prove that the correct value of correlation coefficient should be  $\frac{2}{3}$ .**

**Following are the results of B Com exam:  
Calculate the coefficient of correlation between  
age and successful candidates in examination.**

<b>Age of students</b>	<b>13-14</b>	<b>14-15</b>	<b>15-16</b>	<b>16-17</b>	<b>17-18</b>	<b>18-19</b>	<b>19-20</b>	<b>20-21</b>	<b>21-22</b>	<b>22-23</b>
<b>Students appeared for exam</b>	<b>200</b>	<b>300</b>	<b>100</b>	<b>50</b>	<b>150</b>	<b>400</b>	<b>250</b>	<b>150</b>	<b>25</b>	<b>75</b>
<b>Successful students</b>	<b>124</b>	<b>180</b>	<b>65</b>	<b>34</b>	<b>99</b>	<b>252</b>	<b>145</b>	<b>81</b>	<b>12</b>	<b>33</b>

**(Answer  $r = -0.7747$ )**



# Calculation of PE – Probable Error

Probable Error helps in interpreting the value of Coefficient of Correlation. It tells us about the reliability of the value of Coefficient of Correlation. It is obtained by the following formula:

$$P.E.r = 0.6745 \frac{1 - r^2}{\sqrt{N}}$$

- If the value of  $r$  is less than PE then  $r$  is not at all significant i.e. there is no evidence of correlation.
- **If value of  $r$  is 6 times the value PE** then  $r$  is significant i.e. there is a certain evidence of correlation.
- By adding and subtracting the PE from  $r$  we get the upper and lower limits within which the coefficient of correlation of population is expected to lie.

## EXAMPLE

Calculate Probable Error and the limits of correlation in population given that  $r=0.8$  and number of pairs in observed sample is 16.

$$P.E. r = 0.06$$

For Internal Circulation and Academic  
Purpose Only

The following table gives the distribution of items of production and also the relatively defective items among them, according to size groups. Find the correlation coefficient between **size** and **defect quality** and its probable error.

<b>Size Group</b>	<b>15-16</b>	<b>16-17</b>	<b>17-18</b>	<b>18-19</b>	<b>19-20</b>	<b>20-21</b>
<b>Number of Items</b>	200	270	340	360	400	300
<b>Number of defective Items</b>	150	162	170	180	180	114

# Calculation of Coefficient of Determination

The square of the coefficient of correlation is known as **Coefficient of Determination**.

$$r^2$$

It tells us about what amount of variation in dependent variable has been explained by independent variable.

Assume that the coefficient of correlation between rainfall and per acre yield of rice is 0.8. Find out the coefficient of determination and comment on its value.

$$\text{Coefficient of determination} = r^2 = (0.8)^2 = 0.64$$

### Comments:

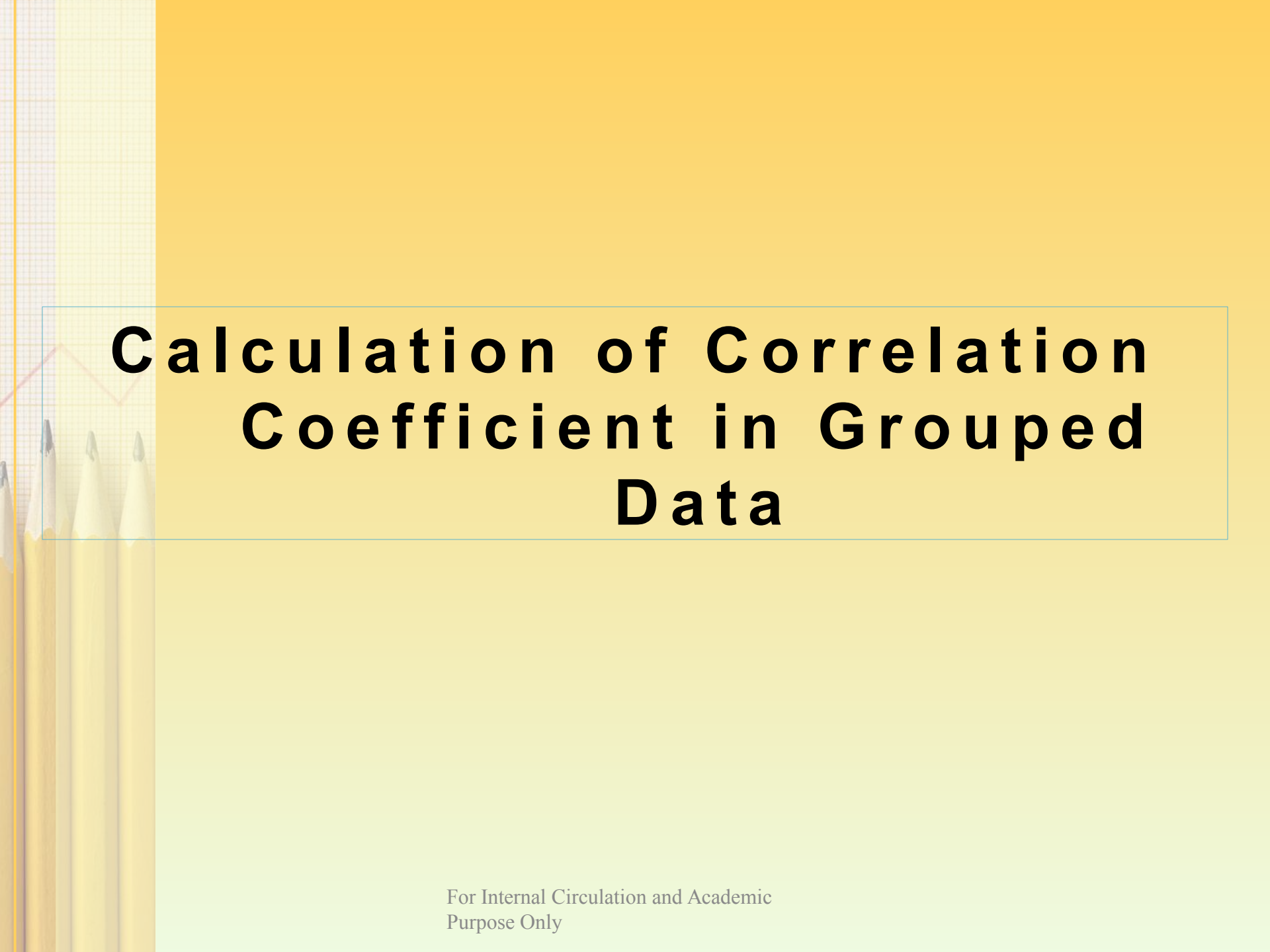
*0.64 Coefficient of Determination means that 64% variation in per acre yield of rice (Dependent Variable) is explained by rainfall (Independent Variable). 36% variation is unexplained by rainfall; it may be due to other factors such as use of fertilizers, soil and seed quality, etc.*

Following table gives the results of an examination. Calculate Karl Pearson's Correlation Coefficient and its probable error? Also comment if the value of correlation coefficient obtained is significant or not.

<b>Age</b>	<b>13-14</b>	<b>14-15</b>	<b>15-16</b>	<b>16-17</b>	<b>17-18</b>
<b>percent age of failures</b>	<b>39</b>	<b>40</b>	<b>43</b>	<b>43</b>	<b>36</b>
<b>Age</b>	<b>18-19</b>	<b>19-20</b>	<b>20-21</b>	<b>21-22</b>	
<b>percent age of failures</b>	<b>39</b>	<b>48</b>	<b>44</b>	<b>56</b>	

$$r = + 0.658 \text{ and PEr} = 0.127$$

For Internal Circulation and Academic  
Purpose Only



# Calculation of Correlation Coefficient in Grouped Data

**Calculate the coefficient of correlation for the following data.**

<b>AGE OF HUSBANDS</b>	<b>AGE OF WIVES</b>					
	<b>10-20</b>	<b>20-30</b>	<b>30-40</b>	<b>40-50</b>	<b>50-60</b>	<b>TOTAL</b>
<b>15-25</b>	<b>6</b>	<b>3</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>9</b>
<b>25-35</b>	<b>3</b>	<b>15</b>	<b>11</b>	<b>-</b>	<b>-</b>	<b>29</b>
<b>35-45</b>	<b>-</b>	<b>11</b>	<b>14</b>	<b>7</b>	<b>-</b>	<b>32</b>
<b>45-55</b>	<b>-</b>	<b>-</b>	<b>6</b>	<b>12</b>	<b>3</b>	<b>21</b>
<b>55-65</b>	<b>-</b>	<b>-</b>	<b>-</b>	<b>3</b>	<b>6</b>	<b>9</b>
<b>TOTAL</b>	<b>9</b>	<b>29</b>	<b>31</b>	<b>22</b>	<b>9</b>	<b>100</b>



# FORMULA

$$r = \frac{N \sum f dx dy - \{(\sum f dx) \times (\sum f dy)\}}{\sqrt{N \sum f dx^2 - (\sum f dx)^2} \sqrt{N \sum f dy^2 - (\sum f dy)^2}}$$

			X	10-20	20-30	30-40	40-50	50-60	Total			
			MV	15	25	35	45	55				
			dX (X-35)	-20	-10	0	10	20				
Y	MV	dY Y-40		-2	-1	0	1	2		fdY	fdY <sup>2</sup>	fdxd y
15-25	20	-20	-2	6	3	-	-	-	9	-18	36	30
25-35	30	-10	-1	3	15	11	-	-	29	-29	29	21
35-45	40	0	0	-	11	14	7	-	32	0	0	0
45-55	50	10	1	-	-	6	12	3	21	21	21	18
55-65	60	20	2	-	-	-	3	6	9	18	36	30
<b>Total</b>				9	29	31	22	9	100	$\Sigma fdY = -8$	$\Sigma fdY^2 = 122$	
<b>fdX</b>				-18	-29	0	22	18	$\Sigma fdX = -7$			

For Internal Circulation and Academic Purpose Only

$$r = \frac{100 \times 99 (-7 \times -8)}{\sqrt{100 \times 123 - (-7)^2} \sqrt{100 \times 122 - (-8)^2}}$$

$$r = +0.8073$$

# RANK CORRELATION (SPEARMAN'S CORRELATION COEFFICIENT)

$$R = 1 - \frac{6 \sum D^2}{N(N^2 - 1)}$$

## When Ranks are given:

- Take the differences of two ranks ( $R_1 - R_2$ ) and denote these differences by  $D$ .
- Square these differences and obtain the total  $\sum D^2$
- Apply the formula given above.

# EXAMPLE

The ranking of 10 students in two subjects A and B are given below. Calculate the Spearman's Correlation Coefficient.

<b>A</b>	<b>6</b>	<b>5</b>	<b>3</b>	<b>10</b>	<b>2</b>	<b>4</b>	<b>9</b>	<b>7</b>	<b>8</b>	<b>1</b>
<b>B</b>	<b>3</b>	<b>8</b>	<b>4</b>	<b>9</b>	<b>1</b>	<b>6</b>	<b>10</b>	<b>7</b>	<b>5</b>	<b>2</b>

<b>R1</b>	<b>R2</b>	<b>(R1 - R2) = D</b>	<b>D<sup>2</sup></b>
<b>6</b>	<b>3</b>	<b>3</b>	<b>9</b>
<b>5</b>	<b>8</b>	<b>-3</b>	<b>9</b>
<b>3</b>	<b>4</b>	<b>-1</b>	<b>1</b>
<b>10</b>	<b>9</b>	<b>1</b>	<b>1</b>
<b>2</b>	<b>1</b>	<b>1</b>	<b>1</b>
<b>4</b>	<b>6</b>	<b>-2</b>	<b>4</b>
<b>9</b>	<b>10</b>	<b>-1</b>	<b>1</b>
<b>7</b>	<b>7</b>	<b>0</b>	<b>0</b>
<b>8</b>	<b>5</b>	<b>3</b>	<b>9</b>
<b>1</b>	<b>2</b>	<b>-1</b>	<b>1</b>
		<b>∑D<sup>2</sup> =</b>	<b>36</b>

$$R = 1 - \frac{6 \sum D^2}{N(N^2 - 1)}$$

$$R = 1 - \frac{6 \times 36}{10(10^2 - 1)}$$

$$R = 0.782$$

Two ladies were asked to rank 7 different types of lipsticks. The ranks given by them are as follows: Calculate Spearman's rank correlation coefficient.

Lipsticks	A	B	C	D	E	F	G
Neelu	2	1	4	3	5	7	6
Neena	1	3	2	4	5	6	7

**(Answer: 0.786)**

Ten competitors in a beauty contest are ranked by three judges in the following order.

<b>1st judge</b>	1	6	5	10	3	2	4	9	7	8
<b>2nd judge</b>	3	5	8	4	7	10	2	1	6	9
<b>3rd judge</b>	6	4	9	8	1	2	3	10	5	7

Use the rank correlation coefficient to determine which pair of judges has the nearest approach to common tastes in beauty.

**(Answer: Pair of 1 and 3,  $R=0.636$ )**



## Where Ranks are NOT given:

- Assign the ranks by taking either the highest value or lowest value as 1 and find out the R.

Calculate the Spearman's coefficient of correlation between marks given to ten students by judges X and Y in a certain competitive examination.

Student No.	1	2	3	4	5	6	7	8	9	10
Marks by judge X	52	53	42	60	45	41	37	38	25	27
Marks by judge Y	65	68	43	38	77	48	35	30	25	50

## Where Ranks are EQUAL:

- In some cases it may be found necessary to rank two or more entries as equal. In such a case it is customary to give each individual an average rank.
- ***Thus, if two entries are ranked equal at 5th place they are each given the rank  $(5+6)/2 = 5.5$  while if 3 are ranked equal at 5th place they are given the rank  $(5+6+7)/3 = 6$ .***
- Where equal ranks are assigned to some entries an adjustment in the formula for calculating the rank coefficient of correlation is made.
- **The adjustment consists of adding  $1/12 (m^3 - m)$  to the value of  $\sum D^2$ .**
- Here,  $m$  stands for number of items whose ranks are common.

## EXAMPLE

Obtain the rank correlation coefficient between the variables X and Y from the following pairs of observed values.

<b>X</b>	<b>50</b>	<b>55</b>	<b>65</b>	<b>50</b>	<b>55</b>	<b>60</b>	<b>50</b>	<b>65</b>	<b>70</b>	<b>75</b>
<b>Y</b>	<b>110</b>	<b>110</b>	<b>115</b>	<b>125</b>	<b>140</b>	<b>115</b>	<b>130</b>	<b>120</b>	<b>115</b>	<b>160</b>

## Number of repetitions in series X

50 is repeated 3 times (m=3)

55 is repeated 2 times (m=2)

65 is repeated 2 times (m=2)

## Number of repetitions in series Y

115 is repeated 3 times (m=3)

110 is repeated 2 times (m=2)

$$R = 1 - \frac{6 \left[ \sum D^2 + \frac{1}{12} (m^3 - m) + \frac{1}{12} (m^3 - m) + \frac{1}{12} (m^3 - m) + \frac{1}{12} (m^3 - m) + \frac{1}{12} (m^3 - m) \right]}{N(N^2 - 1)}$$

## EXAMPLE

Obtain the rank correlation coefficient between the variables X and Y from the following pairs of observed values.

<b>X</b>	<b>15</b>	<b>10</b>	<b>20</b>	<b>28</b>	<b>12</b>	<b>10</b>	<b>16</b>	<b>18</b>
<b>Y</b>	<b>16</b>	<b>14</b>	<b>10</b>	<b>12</b>	<b>11</b>	<b>15</b>	<b>18</b>	<b>12</b>

**ANSWER:      $R = - 0.369$**

## EXAMPLE

Obtain the rank correlation coefficient between the variables X and Y from the following pairs of observed values.

<b>X</b>	<b>40</b>	<b>50</b>	<b>60</b>	<b>60</b>	<b>80</b>	<b>50</b>	<b>70</b>	<b>60</b>
<b>Y</b>	<b>80</b>	<b>120</b>	<b>160</b>	<b>170</b>	<b>130</b>	<b>200</b>	<b>210</b>	<b>130</b>

**ANSWER:      $R = 0.429$**

## EXAMPLE

Obtain the rank correlation coefficient from the following data.

Sr No.	1	2	3	4	5	6	7	8	9	10
Rank Diff.	-2	?	-1	+3	+2	0	-4	+3	+3	-2

**ANSWER:  $R = 0.636$**

## EXAMPLE

The coefficient of rank correlation of marks in two subjects for a group of 10 students was found to be 0.5.

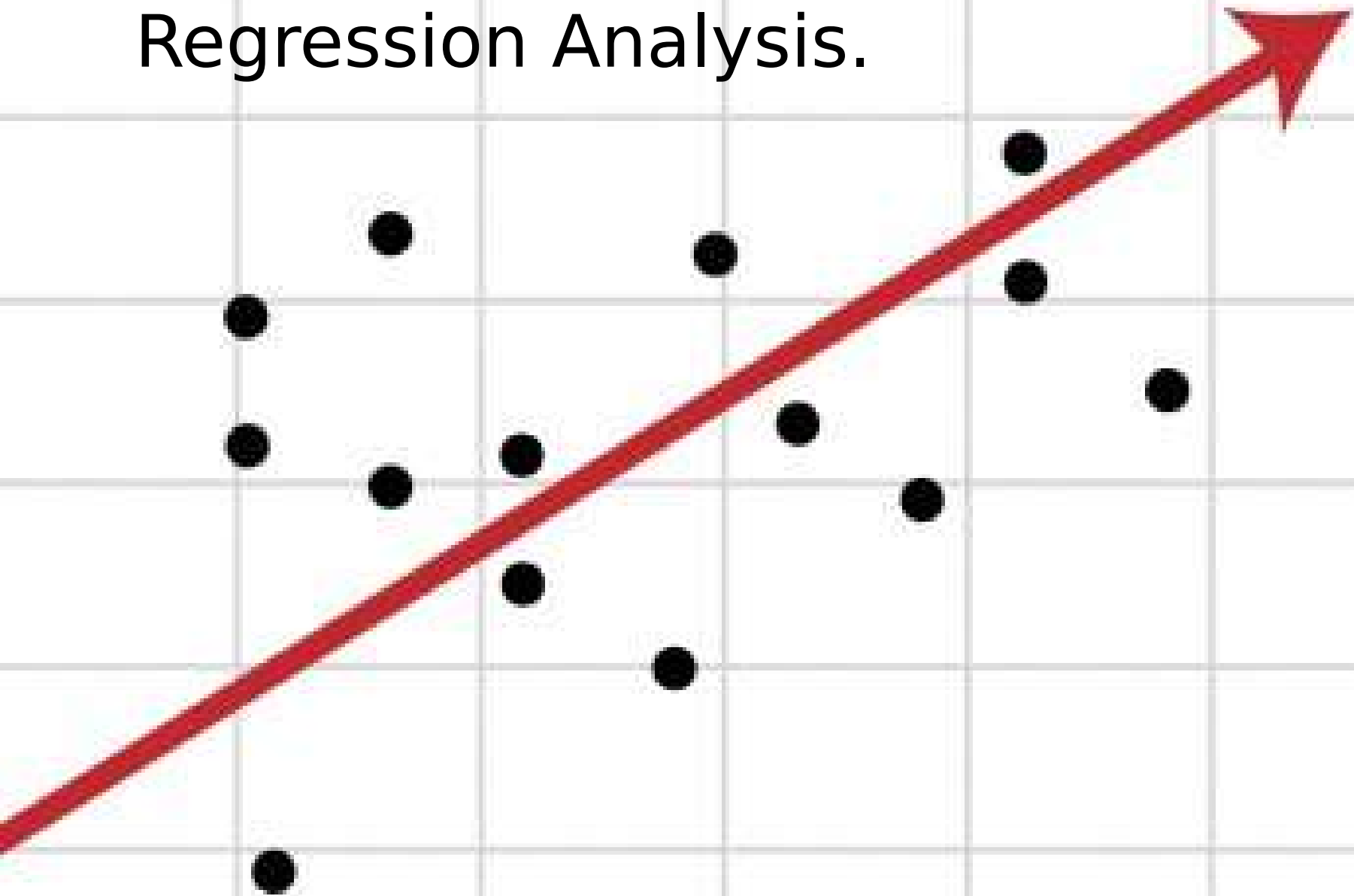
It was later noticed that the difference in ranks in two subjects of one student was wrongly taken as 3 instead of 7.

Find the correct value of Rank Correlation.

**ANSWER:  $R = 0.258$**



# Regression Analysis.



# Correlation & Regression coefficients.

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

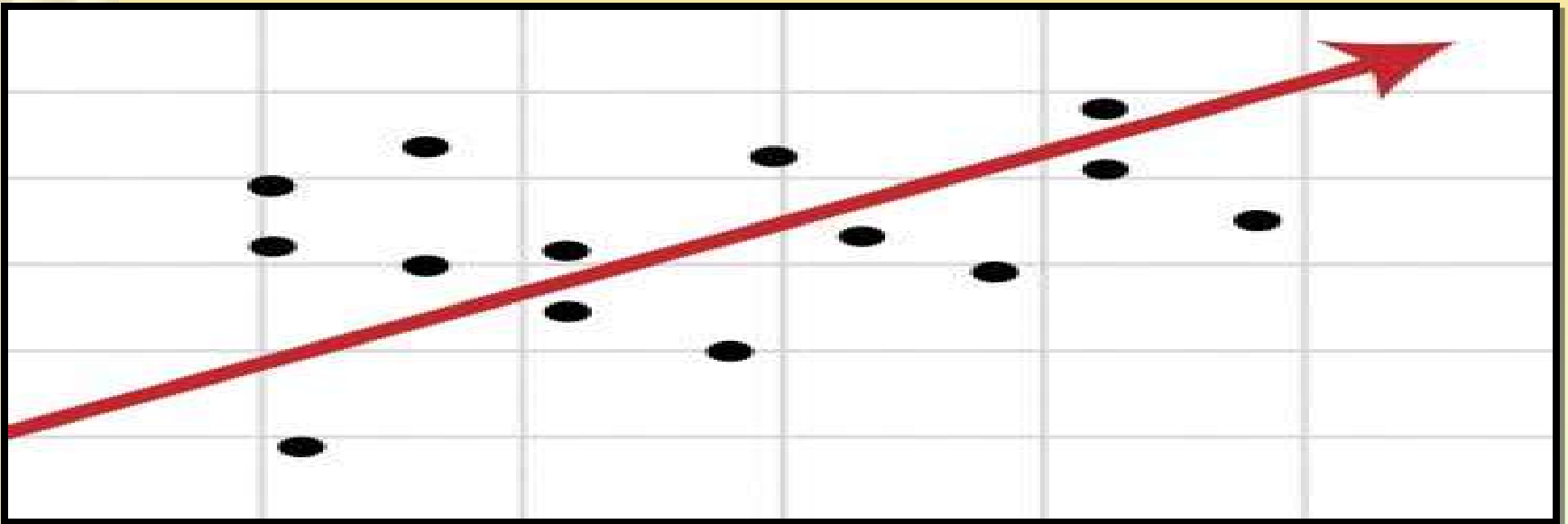
Find the value of the correlation coefficient if regression coefficients of Y on X and X on Y are 0.46 and 0.8 respectively.

**Answer :  $r = 0.606$**

# Method of LEAST SQUARES.

The method of least squares is a mathematical technique.

It is used to obtain the equation of a line which best fits the given data.



# Method of LEAST SQUARES.

Equation of a straight line is  $y = a + bx$

Normal Equations for obtaining the values of a and b are as follows

$$(i) \quad \sum y = Na + b \sum x$$

$$(ii) \quad \sum xy = a \sum x + b \sum x^2$$

# EXAMPLE.

Fit a straight line of Y on X from the following data:

<b>X</b>	<b>0</b>	<b>1</b>	<b>2</b>	<b>3</b>	<b>4</b>	<b>5</b>	<b>6</b>
<b>Y</b>	<b>2</b>	<b>1</b>	<b>3</b>	<b>2</b>	<b>4</b>	<b>3</b>	<b>5</b>

# SOLUTION.

	<b>X</b>	<b>Y</b>	<b>X<sup>2</sup></b>	<b>XY</b>
	<b>0</b>	<b>2</b>	<b>0</b>	<b>0</b>
	<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>
	<b>2</b>	<b>3</b>	<b>4</b>	<b>6</b>
	<b>3</b>	<b>2</b>	<b>9</b>	<b>6</b>
	<b>4</b>	<b>4</b>	<b>16</b>	<b>16</b>
	<b>5</b>	<b>3</b>	<b>25</b>	<b>15</b>
	<b>6</b>	<b>5</b>	<b>36</b>	<b>30</b>
<b>SU M</b>	<b>21</b>	<b>20</b>	<b>91</b>	<b>74</b>

$$\sum y = Na + b \sum x$$

$$\sum xy = a \sum x + 21b$$
$$20 = 7a + 21b$$

$$74 = 21a + 21b$$
$$20 = 7a + 21b$$

**Multiply by 3**

$$60 = 21a + 63b$$

# SOLUTION.

$$74 = 21a + 91b$$

minus

$$60 = 21a + 63b$$

$$14 = 28b$$

$$b = 0.5$$

Putting value of  $b = 0.5$  in above equations;

$$a = 1.357$$

$$Y = 1.357 + 0.5 X$$

**Equation of the line which fits the given data.**

# EXAMPLE.

Fit a straight line of Y on X from the following data:

<b>X</b>	<b>0</b>	<b>2</b>	<b>4</b>	<b>6</b>	<b>8</b>	<b>10</b>	<b>12</b>
<b>Y</b>	<b>4</b>	<b>2</b>	<b>6</b>	<b>4</b>	<b>8</b>	<b>6</b>	<b>10</b>



# EXAMPLE.

Fit a straight line of X on Y from the following data:

<b>X</b>	<b>0</b>	<b>2</b>	<b>4</b>	<b>6</b>	<b>8</b>	<b>10</b>	<b>12</b>
<b>Y</b>	<b>4</b>	<b>2</b>	<b>6</b>	<b>4</b>	<b>8</b>	<b>6</b>	<b>10</b>

**Normal Equations for obtaining the values of a and b are as follows**

$$(i) \sum X = Na + b \sum Y$$

$$(ii) \sum XY = a \sum Y + b \sum Y^2$$

# EXAMPLE.

From the following data obtain the Regression equations of Y on X and X on Y

<b>X</b>	<b>6</b>	<b>2</b>	<b>10</b>	<b>4</b>	<b>8</b>
<b>Y</b>	<b>9</b>	<b>11</b>	<b>5</b>	<b>8</b>	<b>7</b>

**Answer:**

$$Y = 11.9 - 0.65 X$$

$$X = 16.4 - 1.3 Y$$

# Derivation of Lines of Regression directly from the data. (DEVIATION FROM ACTUAL MEAN)

The equation of regression line of X on Y is

$$X - \bar{X} = b_{xy}(Y - \bar{Y})$$

$$X - \bar{X} = \frac{r \sigma_x}{\sigma_y} (Y - \bar{Y})$$

The equation of regression line of Y  
on X is

$$Y - \bar{Y} = b_{yx}(X - \bar{X})$$

$$Y - \bar{Y} = \frac{r \sigma_y}{\sigma_x} (X - \bar{X})$$

Derivation of Lines of Regression directly from the data. (DEVIATION FROM ACTUAL MEAN)

**Regression Coefficient of X on Y is**

$$b_{xy} = \frac{\sum xy}{\sum y^2}$$

**Regression Coefficient of Y on X**

$$b_{yx} = \frac{\sum xy}{\sum x^2}$$

**In these formulas x and y are deviations from mean.**

Regression Coefficients directly from the data.

The regression Coefficient of X on Y is

$$b_{xy} = \frac{r \sigma x}{\sigma y} = \frac{\sum XY - N \bar{X} \bar{Y}}{\sum Y^2 - N (\bar{Y})^2}$$

The regression Coefficient of Y on X  
is

$$b_{yx} = \frac{r \sigma y}{\sigma x} = \frac{\sum XY - N \bar{X} \bar{Y}}{\sum X^2 - N (\bar{X})^2}$$

# EXAMPLE.

From the following data obtain the Regression equations using deviation of means method.

<b>X</b>	<b>6</b>	<b>2</b>	<b>10</b>	<b>4</b>	<b>8</b>
<b>Y</b>	<b>9</b>	<b>11</b>	<b>5</b>	<b>8</b>	<b>7</b>

**Answer:**

$$Y = 11.9 - 0.65 X$$

$$X = 16.4 - 1.3 Y$$

## X on Y

X	x	x <sup>2</sup>	Y	y	y <sup>2</sup>	xy
6	0	0	9	1	1	0
2	-4	16	11	3	9	-12
10	4	16	5	-3	9	-12
4	-2	4	8	0	0	0
8	2	4	7	-1	1	-2
30	0	40	40	0	20	-26

$$X - \bar{X} = b_{xy}(Y - \bar{Y})$$

$$b_{xy} = \frac{\sum xy}{\sum y^2}$$

$$b_{xy} = \frac{-26}{20} = -1.3$$

$$X - 6 = -1.3(Y - 8)$$

$$X - 6 = -1.3Y + 10.4$$

$$X = -1.3Y + 16.4$$

$$\bar{X} = \frac{\sum X}{N} = \frac{30}{5} = 6$$

$$\bar{Y} = \frac{\sum Y}{N} = \frac{40}{5} = 8$$

# EXAMPLE.

Following data relate to the scores of 9 salesmen of a company in an intelligence test and weekly sales in thousands. If the intelligence test score of a salesman is 65 what would be his weekly expected sales?

<b>Test Score</b>	<b>50</b>	<b>60</b>	<b>50</b>	<b>60</b>	<b>80</b>	<b>50</b>	<b>80</b>	<b>40</b>	<b>70</b>
<b>Weekly Sales</b>	<b>30</b>	<b>60</b>	<b>40</b>	<b>50</b>	<b>60</b>	<b>30</b>	<b>70</b>	<b>50</b>	<b>60</b>

**Answer: Rs. 53750**



# EXAMPLE.

Following table shows ages (X) and blood pressure (Y) of 8 persons. Find the expected blood pressure of a person whose age is 49 years.

<b>X</b>	<b>52</b>	<b>63</b>	<b>45</b>	<b>36</b>	<b>72</b>	<b>65</b>	<b>47</b>	<b>25</b>
<b>Y</b>	<b>62</b>	<b>53</b>	<b>51</b>	<b>25</b>	<b>79</b>	<b>43</b>	<b>60</b>	<b>33</b>

**Answer: 49.5**

# EXAMPLE.

In a correlation study the following values are obtained. Find the two regression equations that are associated with above values.

	<b>X</b>	<b>Y</b>
<b>Mean</b>	<b>65</b>	<b>67</b>
<b>Std. Deviation</b>	<b>2.5</b>	<b>3.5</b>
<b>Correlation Coefficient</b>	<b>0.8</b>	

**Answer:**

$$X = 26.72 + 0.57 Y$$

$$Y = -5.8 + 11.2 X$$

# Properties of Regression Coefficients.

1. Both regression coefficients will have the same sign i.e. positive or negative
2. Square root of the product of two regression coefficients will give Correlation coefficient.
3. If one regression coefficient is more than 1 then the other will be less than 1.
4. The regression coefficients will have the same sign as that of the correlation coefficient.
5. Arithmetic mean of Regression coefficient is greater than correlation coefficient.
6. The two regression lines intersect at coordinates which are mean of series X and series Y.

Derivation of Lines of Regression directly from the data. (DEVIATION FROM ASSUMED MEAN)

**Regression Coefficient of X on Y is**

$$b_{xy} = r \frac{\sigma x}{\sigma y} = \frac{N \sum dx dy - (\sum dx X \sum dy)}{N \sum dy^2 - (\sum dy)^2}$$

**Regression Coefficient of Y on X**

**is**

$$b_{yx} = r \frac{\sigma y}{\sigma x} = \frac{N \sum dx dy - (\sum dx X \sum dy)}{N \sum dx^2 - (\sum dx)^2}$$

**In these formulas dx and dy are deviations from ASSUMED mean.**

# EXAMPLE.

From the following data obtain the Regression equations by taking deviation of series X from 5 and of series Y from 7.

<b>X</b>	<b>6</b>	<b>2</b>	<b>10</b>	<b>4</b>	<b>8</b>
<b>Y</b>	<b>9</b>	<b>11</b>	<b>5</b>	<b>8</b>	<b>7</b>

**Answer:**  
 $Y = 11.9 - 0.65 X$   
 $X = 16.4 - 1.3 Y$

# EXAMPLE.

A panel of two judges P and Q graded eight dance performances as below. However, Judge Q was absent during the 7<sup>th</sup> performance. What might have been his ranking if judge Q had been present?

<b>P</b>	<b>46</b>	<b>42</b>	<b>44</b>	<b>40</b>	<b>43</b>	<b>41</b>	<b>37</b>	<b>45</b>
<b>Q</b>	<b>40</b>	<b>38</b>	<b>36</b>	<b>35</b>	<b>39</b>	<b>37</b>	<b>?</b>	<b>41</b>

**Answer:**  
 **$Y = +0.75 X + 5.75$**   
**33.5 marks**

# EXAMPLE.

From the following regression equations find the mean values of X and Y series.

$$8X - 10Y = -66$$

$$40X - 18Y = 214$$

Hint: Regression lines (X on Y and Y on X) cut each other at the point of mean.

**Answer:**

**Mean of X series = 13**

**Mean of Y series = 17**

# EXAMPLE.

In a partially destroyed laboratory record of analysis of correlation data, only the following results are legible:

Variance of  $X = 9$

Regression equations

$$8X - 10Y + 66 = 0$$

$$40X - 18Y = 214$$

Find out :

1. The mean values of  $X$  and  $Y$
2. Coefficient of correlation between  $X$  and  $Y$
3. Standard Deviation of  $Y$



# Solution.

To find the mean values of X and Y we will solve the given equations:

Regression equations given:

$$(i) 8X - 10Y = -66$$

$$(ii) 40X - 18Y = 214$$

Multiplying equation (i) by 5 we get

$$(i) 40X - 50Y = -330$$

Subtracting equation (ii) from equation (i) we get:

$$-32 Y = -544$$

$$**Y = -544 / -32 = 17 (Mean of Y = 17)**$$

# Solution.

Substituting the value of  $Y=17$  in equation (i)

$$(i) 8X - 10Y = -66$$

$$8X - 10(17) = -66$$

$$8X - 170 = -66$$

$$8X = 104$$

$$\mathbf{X = 104 / 8 = 13 \text{ (Mean of X = 13)}}$$

To find the Correlation coefficient ( $r$ ) we must find the regression coefficients ( $b_{xy}$  and  $b_{yx}$ ). However, we don't know which equation is  $X$  on  $Y$  and which equation is  $Y$  on  $X$ . So let's assume equation (i) as  $X$  on  $Y$

# Solution.

Assuming equation (i) as X on Y we will try to find  $b_{xy}$

$$(i) 8X - 10Y = -66$$

$$8X = -66 + 10Y$$

$$X = (-66/8) + (10/8) Y$$

$$X = a + b_{xy} Y$$

$$b_{xy} = 10 / 8 = 1.25$$

From equation (ii) we can find out  $b_{yx}$

$$(ii) 40X - 18Y = 214$$

$$Y = -(214/18) + (40/18) X$$

$$b_{yx} = 40 / 18 = 2.22$$

## Solution.

Here we can see that both the regression coefficients are greater than 1 which is not possible. Therefore, our assumption that equation (i) is equation of X on Y is wrong. Equation (i) is actually equation of Y on X. So we can write the equation as:

$$-10 Y = - 8X - 66$$

$$Y = (8/10) X + 6.6 \text{ which means that } b_{yx} = 0.8$$

From equation (ii) i.e. X on Y we get

$$40X = 214 + 18Y \text{ or } X = (214/40) + (18/40) Y \text{ which means that } b_{xy} = 18/40 =$$

$$0.45$$

## Solution.

We know that Correlation coefficient is square root of product of Regression coefficients:

$$r = \text{square root of } 0.8 \times 0.45 = 0.6$$

We can further calculate standard deviation of Y using the formula of Regression Coefficient.

$$\text{Standard Deviation of Y} = 4$$

# STANDARD ERROR OF ESTIMATES.

- ❑ The standard error of estimate measures the accuracy of the estimated figures.
- ❑ Smaller the value of standard error, closer will be the dots to the regression line and better will be the estimates.
- ❑ If standard error of estimate is zero then there is no variation about the line and correlation is perfect.

$$S_{xy} = \sqrt{\frac{\sum(X - X_c)^2}{N}}$$

$$S_{xy} = \sigma_x \sqrt{1 - r^2}$$

$$S_{xy} = \sqrt{\frac{\sum X^2 - a\sum X - b\sum XY}{N}}$$

$$S_{yx} = \sqrt{\frac{\sum(Y - Y_c)^2}{N}}$$

$$S_{yx} = \sigma_y \sqrt{1 - r^2}$$

$$S_{yx} = \sqrt{\frac{\sum Y^2 - a\sum Y - b\sum XY}{N}}$$

## Uses of Regression analysis.

1. Regression line facilitates to predict the values of dependent variable from the given value of independent variable.
2. Standard Error facilitates to obtain a measure of error involved in using the regression line as a basis of estimation.
3. Regression coefficients help us to calculate coefficient of correlation and coefficient of determination,
4. Regression analysis is a highly useful tool in Economics and business.

# TIME SERIES & FORECASTING.

For Internal Circulation and Academic  
Purpose Only



- Components of Time Series.**
- Trend - Moving averages, semi-averages and least-squares.**
- Seasonal variation, cyclic variation and irregular variation.**
- Index numbers, calculation of seasonal indices.**
- Additive and multiplicative models.**
- Forecasting, Non linear trend – second degree parabolic trends**

# TIME SERIES.

A Time Series is a set of observations taken at specified times, usually at equal intervals.

Mathematically, a time series is defined by values  $Y_1, Y_2, \dots$  of a variable at times  $t_1, t_2, \dots$ . Thus  $Y$  is a function of  $t$  symbolized by  $Y = F(t)$ .

# UTILITY OF TIME SERIES ANALYSIS.

Helps in understanding past behaviour.

Helps in planning future operations.

Helps in evaluating current accomplishments.

Facilitates comparison.

# COMPONENTS OF A TIME SERIES.

Secular Trend	T
Seasonal Variations	S
Cyclical Variations	C
Irregular Variations	I

$$Y = T + S + C + I$$

Additive Model

or

$$Y = T \times S \times C \times I$$

Multiplicative Model

# Measurement of Trend

Freehand or Graphical method

Semi-average method

Moving average method

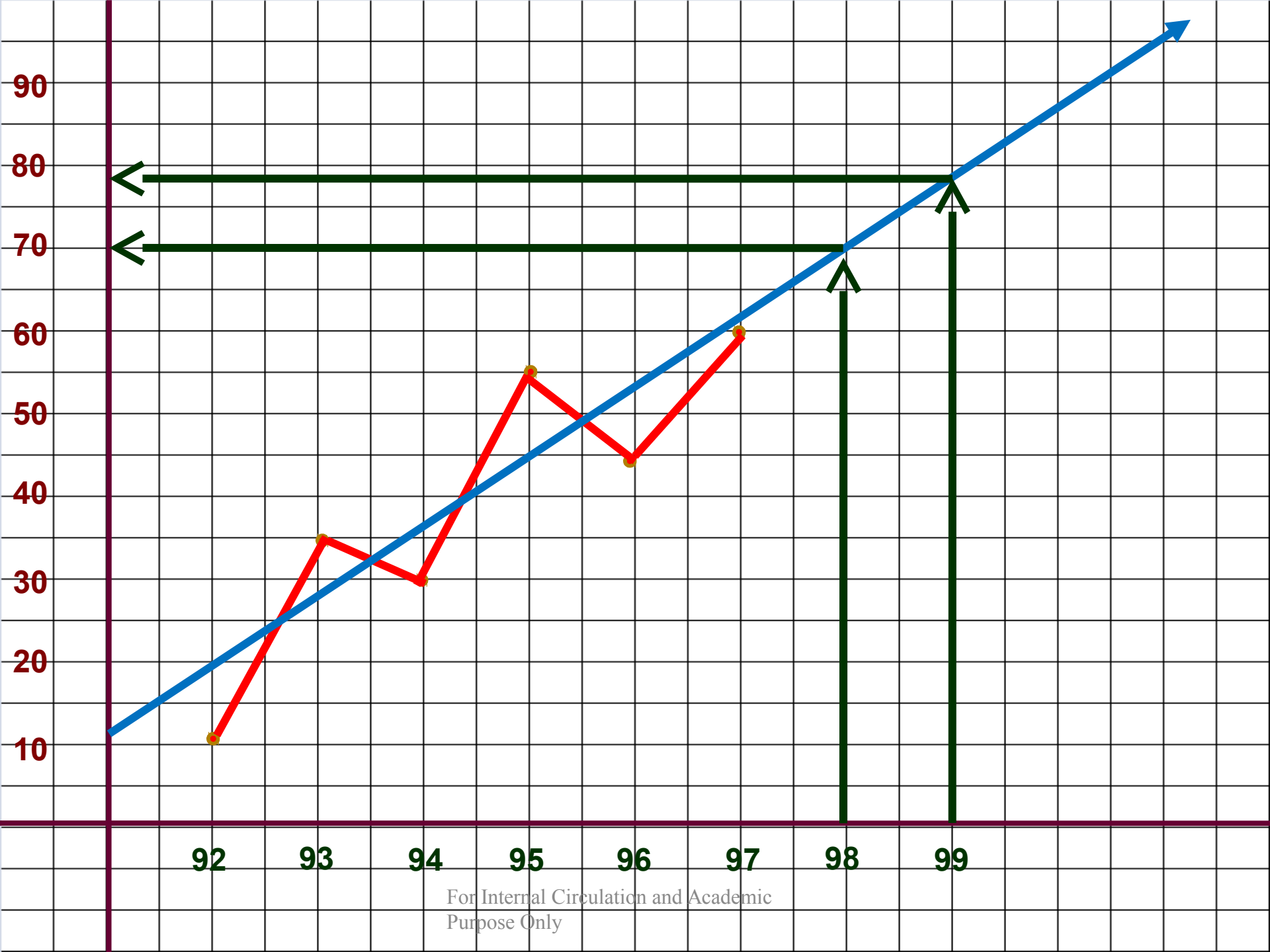
Least squares method.

# Freehand or Graphical method

1. Plot the time series on a graph paper.
2. Examine carefully the direction of dots.
3. Draw a straight line according to personal judgement.

Fit a trend line to the following data using the Freehand method and predict values for 1998 & 1999.

<b>YEAR</b>	<b>SUGAR PRODUCTION (Million Tonnes)</b>
<b>1992</b>	<b>10</b>
<b>1993</b>	<b>35</b>
<b>1994</b>	<b>30</b>
<b>1995</b>	<b>55</b>
<b>1996</b>	<b>45</b>
<b>1997</b>	<b>60</b>



92

93

94

95

96

97

98

99

For Internal Circulation and Academic Purpose Only



Fit a trend line to the following data using the Freehand method and predict values for 2009 & 2010.

<b>YEAR</b>	<b>Sales (Millions)</b>
<b>2001</b>	<b>5</b>
<b>2002</b>	<b>15</b>
<b>2003</b>	<b>10</b>
<b>2004</b>	<b>25</b>
<b>2005</b>	<b>30</b>
<b>2006</b>	<b>20</b>
<b>2007</b>	<b>35</b>
<b>2008</b>	<b>45</b>

# Method of SEMI AVERAGES

1. Divide the data in two equal parts. In case of odd years, omit the middle year.
2. Obtain the average of each part.
3. Plot the two points against the midpoint of class interval on a graph.
4. Joint the two points to get a trend line.

Fit a trend line to the following data using the semi averages method and predict values for 2009 & 2010.

<b>YEAR</b>	<b>Sales (Millions)</b>
<b>1993</b>	<b>102</b>
<b>1994</b>	<b>105</b>
<b>1995</b>	<b>114</b>
<b>1996</b>	<b>110</b>
<b>1997</b>	<b>108</b>
<b>1998</b>	<b>116</b>
<b>1999</b>	<b>112</b>

# Method of MOVING AVERAGES

There can be two ways to calculate moving averages.

1. 3 year, 5 year or 7 year moving averages. These are called odd year moving averages.

Or

2. 2 year, 4 year, 6 year or 8 year moving averages.

There is a slight difference in these two ways.

Calculate the 3 year moving averages of the production figures given below.

<b>YEAR</b>	<b>PRODUCTION</b>	<b>YEAR</b>	<b>PRODUCTION</b>
<b>1985</b>	<b>15</b>	<b>1993</b>	<b>63</b>
<b>1986</b>	<b>21</b>	<b>1994</b>	<b>70</b>
<b>1987</b>	<b>30</b>	<b>1995</b>	<b>74</b>
<b>1988</b>	<b>36</b>	<b>1996</b>	<b>82</b>
<b>1989</b>	<b>42</b>	<b>1997</b>	<b>90</b>
<b>1990</b>	<b>46</b>	<b>1998</b>	<b>95</b>
<b>1991</b>	<b>50</b>	<b>1999</b>	<b>102</b>
<b>1992</b>	<b>56</b>		

Construct 5 year moving averages of the number of students studying in a college.

<b>YEAR</b>	<b>No. of students</b>	<b>YEAR</b>	<b>No. of students</b>
<b>1990</b>	<b>332</b>	<b>1995</b>	<b>405</b>
<b>1991</b>	<b>317</b>	<b>1996</b>	<b>410</b>
<b>1992</b>	<b>357</b>	<b>1997</b>	<b>427</b>
<b>1993</b>	<b>392</b>	<b>1998</b>	<b>405</b>
<b>1994</b>	<b>402</b>	<b>1999</b>	<b>438</b>

Calculate the trend values by taking 4 year moving averages.

<b>YEAR</b>	<b>VALUE</b>	<b>YEAR</b>	<b>VALUE</b>
<b>1984</b>	<b>12</b>	<b>1991</b>	<b>100</b>
<b>1985</b>	<b>25</b>	<b>1992</b>	<b>82</b>
<b>1986</b>	<b>39</b>	<b>1993</b>	<b>65</b>
<b>1987</b>	<b>54</b>	<b>1994</b>	<b>49</b>
<b>1988</b>	<b>70</b>	<b>1995</b>	<b>34</b>
<b>1989</b>	<b>87</b>	<b>1996</b>	<b>20</b>
<b>1990</b>	<b>105</b>	<b>1997</b>	<b>7</b>

# WEIGHTED MOVING AVERAGES

Generally weighted moving average is used to forecast trend figures.

WMA gives higher weightage to recent figures.



Calculate the trend values using 3 year WMA for the following data. Weights are to be assigned in order 1, 2, 3.

<b>YEAR</b>	<b>SALES</b>	<b>YEAR</b>	<b>SALES</b>
<b>2001</b>	<b>10</b>	<b>2008</b>	<b>18</b>
<b>2002</b>	<b>12</b>	<b>2009</b>	<b>20</b>
<b>2003</b>	<b>12</b>	<b>2010</b>	<b>18</b>
<b>2004</b>	<b>14</b>	<b>2011</b>	<b>24</b>
<b>2005</b>	<b>16</b>	<b>2012</b>	<b>28</b>
<b>2006</b>	<b>18</b>		
<b>2007</b>	<b>22</b>		

<b>Year</b>	<b>Sales</b>	<b>WT</b>	<b>Wtd Sales</b>	<b>3 Y WMT</b>	<b>3 Y WMA</b>
<b>01</b>	<b>10</b>	<b>1</b>	<b>10</b>	<b>----</b>	<b>----</b>
<b>02</b>	<b>12</b>	<b>2</b>	<b>24</b>	<b>70</b>	<b>11.66</b>
<b>03</b>	<b>12</b>	<b>3</b>	<b>36</b>	<b>74</b>	<b>12.33</b>
<b>04</b>	<b>14</b>	<b>1</b>	<b>14</b>	<b>82</b>	<b>13.66</b>
<b>05</b>	<b>16</b>	<b>2</b>	<b>32</b>	<b>100</b>	<b>16.66</b>
<b>06</b>	<b>18</b>	<b>3</b>	<b>54</b>	<b>108</b>	<b>18</b>
<b>07</b>	<b>22</b>	<b>1</b>	<b>22</b>	<b>112</b>	<b>18.66</b>
<b>08</b>	<b>18</b>	<b>2</b>	<b>36</b>	<b>118</b>	<b>19.66</b>
<b>09</b>	<b>20</b>	<b>3</b>	<b>60</b>	<b>114</b>	<b>19</b>
<b>10</b>	<b>18</b>	<b>1</b>	<b>18</b>	<b>126</b>	<b>21</b>
<b>11</b>	<b>24</b>	<b>2</b>	<b>48</b>	<b>150</b>	<b>25</b>
<b>12</b>	<b>28</b>	<b>3</b>	<b>84</b>	<b>----</b>	<b>----</b>

For Internal Circulation and Academic Purpose Only

Calculate the trend values using 5 year WMA for the following data. Weights are to be assigned in order 1, 2, 2, 3, 3.

<b>YEAR</b>	<b>SALES</b>	<b>YEAR</b>	<b>SALES</b>
<b>1990</b>	<b>18</b>	<b>1997</b>	<b>32</b>
<b>1991</b>	<b>20</b>	<b>1998</b>	<b>28</b>
<b>1992</b>	<b>21</b>	<b>1999</b>	<b>36</b>
<b>1993</b>	<b>26</b>	<b>2000</b>	<b>34</b>
<b>1994</b>	<b>22</b>	<b>2001</b>	<b>35</b>
<b>1995</b>	<b>24</b>	<b>2002</b>	<b>44</b>
<b>1996</b>	<b>30</b>	<b>2003</b>	<b>46</b>
		<b>2004</b>	<b>42</b>

# LEAST SQUARES METHOD

EQUATION OF STRAIGHT TREND LINE

$$Y = a + bX$$

Normal Equations for obtaining the values of a and b are as follows

$$(i) \sum Y = Na + b \sum X$$

$$(ii) \sum XY = a \sum X + b \sum X^2$$

N = Number of years,

X = Converted value for years.

$$(i) \sum Y = Na + b \sum X$$

$$(ii) \sum XY = a \sum X + b \sum X^2$$

If we take the middle year as year of origin  
then  $\sum X = 0$ .

**Then  $a = (\sum Y / N) = \text{Mean of } Y$**

AND

Putting the value of  $\sum X = 0$  in equation (ii).

**Then  $b = (\sum XY / \sum X^2)$**

Fit a straight line trend for the following series and Estimate the values for 1997

<b>YEAR</b>	<b>Production</b>
<b>1990</b>	<b>60</b>
<b>1991</b>	<b>72</b>
<b>1992</b>	<b>75</b>
<b>1993</b>	<b>65</b>
<b>1994</b>	<b>80</b>
<b>1995</b>	<b>85</b>
<b>1996</b>	<b>95</b>

$$Y = 76 + 4.857 X$$

$$Y_{1997} = 95.428$$

Fit a straight line trend for the following series and Estimate the values for 1998

<b>YEAR</b>	<b>Production</b>
<b>1989</b>	<b>38</b>
<b>1990</b>	<b>40</b>
<b>1991</b>	<b>65</b>
<b>1992</b>	<b>72</b>
<b>1993</b>	<b>69</b>
<b>1994</b>	<b>60</b>
<b>1995</b>	<b>87</b>
<b>1996</b>	<b>95</b>

$$Y = 65.75 + 3.667 X$$

$$Y_{1997} = 106.087$$

# CALCULATION OF SEASONAL INDEX

There are 4 methods of computing seasonal component of time series:

1. Simple Average Method
2. Ratio to Trend Method
3. Ratio to Moving Average Method
4. Link Relative Method

We will study only the first method...



# Simple Average Method

1. Find the Quarterly totals.
2. Find Quarterly averages for each quarter.
3. Find grand average of quarterly averages.
4. Find the seasonal index of each quarter by dividing its quarterly average by grand average.

The given table shows trend free figures of quarterly sales made by a mega mall. Find the seasonal indices.

<b>YEAR</b>	<b>I</b>	<b>II</b>	<b>III</b>	<b>IV</b>
<b>2003</b>	<b>39</b>	<b>20</b>	<b>60</b>	<b>85</b>
<b>2004</b>	<b>45</b>	<b>23</b>	<b>62</b>	<b>90</b>
<b>2005</b>	<b>60</b>	<b>32</b>	<b>76</b>	<b>100</b>
<b>2006</b>	<b>47</b>	<b>35</b>	<b>65</b>	<b>85</b>

The following time series data on consumption of cold drinks contains only seasonal and irregular variations. Construct indices for seasonal variations using simple arithmetic mean.

<b>YEAR</b>	<b>I</b>	<b>II</b>	<b>III</b>	<b>IV</b>
<b>2003</b>	<b>39</b>	<b>20</b>	<b>60</b>	<b>85</b>
<b>2004</b>	<b>45</b>	<b>23</b>	<b>62</b>	<b>90</b>
<b>2005</b>	<b>60</b>	<b>32</b>	<b>76</b>	<b>100</b>
<b>2006</b>	<b>47</b>	<b>35</b>	<b>65</b>	<b>85</b>

Following data gives monthly production figures. Find monthly seasonal indices.

<b>Yr</b>	<b>01</b>	<b>02</b>	<b>03</b>	<b>04</b>	<b>05</b>
<b>Jan</b>	<b>31</b>	<b>34</b>	<b>36</b>	<b>39</b>	<b>42</b>
<b>Feb</b>	<b>28</b>	<b>30</b>	<b>32</b>	<b>34</b>	<b>37</b>
<b>Mar</b>	<b>27</b>	<b>29</b>	<b>28</b>	<b>34</b>	<b>37</b>
<b>Apr</b>	<b>25</b>	<b>26</b>	<b>26</b>	<b>31</b>	<b>33</b>
<b>May</b>	<b>23</b>	<b>24</b>	<b>25</b>	<b>29</b>	<b>31</b>
<b>Jun</b>	<b>21</b>	<b>23</b>	<b>25</b>	<b>27</b>	<b>29</b>
<b>Jul</b>	<b>22</b>	<b>24</b>	<b>28</b>	<b>28</b>	<b>30</b>
<b>Aug</b>	<b>24</b>	<b>26</b>	<b>30</b>	<b>30</b>	<b>33</b>
<b>Sep</b>	<b>26</b>	<b>28</b>	<b>34</b>	<b>32</b>	<b>35</b>
<b>Oct</b>	<b>30</b>	<b>32</b>	<b>36</b>	<b>36</b>	<b>39</b>
<b>Nov</b>	<b>32</b>	<b>34</b>	<b>39</b>	<b>38</b>	<b>42</b>
<b>Dec</b>	<b>34</b>	<b>36</b>	<b>40</b>	<b>41</b>	<b>45</b>

# References and Suggested Readings

Fundamentals of Statistics by S.C. Gupta

Statistics Methods by S.P.Gupta