

## Index

<b>Unit – 4 → Basic Statistics .....</b>	<b>3</b>
1) Method 1 → Measure of Central Tendency .....	3
2) Method 2 → Measure of Dispersion.....	14
3) Method 3 → Covariance .....	21
4) Method 4 → Correlation Coefficient .....	22
5) Method 5 → Rank Correlation Coefficient .....	24
6) Method 6 → Linear Regression.....	28
7) Method 7 → Curve Fitting.....	31
8) Method 8 → Fitting of Trend by Moving Average Method .....	36



## Unit – 4 $\rightsquigarrow$ Basic Statistics

### Method 1 $\rightsquigarrow$ Measure of Central Tendency

#### Introduction

- Statistics is the branch of science where we plan, gather and analyze information about a particular collection of objects under investigation.
- Statistics techniques are used in every other field of science, engineering and humanity, ranging from computer science to industrial engineering to sociology and psychology.
- For any statistical problem the initial information collection from the sample may look messy, and hence confusing. This initial information needs to be organized first before we make any sense out of it.

#### Central Tendency

- The central tendency of a distribution is an estimate of the **center** of a distribution of values.
- There are three measures to estimate central tendency which is
  - (1) Mean( $\bar{x}$ )
  - (2) Median(M)
  - (3) Mode(Z)
  - (4) Quartiles
  - (5) Percentiles

1.1 Mean

- The mean means **average**.
- Mean is denoted by “ $\bar{x}$ ” and read as x bar.
- Table of different formulae of mean.

Method	Ungrouped Data	Discrete Grouped Data	Continuous Grouped Data
Direct Method	$\frac{\sum x_i}{n}$	$\frac{\sum f_i x_i}{\sum f_i}$	
Assumed Mean Method	$A + \frac{\sum d_i}{n}$	$A + \frac{\sum f_i d_i}{\sum f_i}$	
Step Deviation Method	-----	-----	$A + \frac{\sum f_i u_i}{\sum f_i} \times c$

- $n$  = total number of observations
- In case of continuous frequency distribution,  
 $x_i$  = mid value of the respective class.
- In case of **assumed mean method**, A can be any value from  $x_i$ .
- Use below formula to calculate  $d_i$  &  $u_i$

$$d_i = x_i - A ; u_i = \frac{x_i - A}{c}$$

Example of Method-1.1: Examples of Mean

C	1	Find the mean of data 10.2, 9.5, 8.3, 9.7, 9.5, 11.1, 7.8, 8.8, 9.5, 10.  <b>Answer: 9.44</b>														
C	2	Find the mean for following data: <table border="1"><tr><td>Marks obtained</td><td>20</td><td>9</td><td>25</td><td>50</td><td>40</td><td>80</td></tr><tr><td>Number of students</td><td>6</td><td>4</td><td>16</td><td>7</td><td>8</td><td>2</td></tr></table> <b>Answer: 32.23</b>	Marks obtained	20	9	25	50	40	80	Number of students	6	4	16	7	8	2
Marks obtained	20	9	25	50	40	80										
Number of students	6	4	16	7	8	2										
C	3	Find the mean using direct method, assumed mean method and step deviation method: <table border="1"><tr><td>Marks</td><td>0 – 10</td><td>10 – 20</td><td>20 – 30</td><td>30 – 40</td><td>40 – 50</td></tr><tr><td>No. of students</td><td>5</td><td>10</td><td>40</td><td>20</td><td>25</td></tr></table> <b>Answer: 30</b>	Marks	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50	No. of students	5	10	40	20	25		
Marks	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50											
No. of students	5	10	40	20	25											
C	4	Find the missing frequency $f_1$ and $f_2$ in the table given below, it is being given that the mean of the given frequency distribution is 50. <table border="1"><tr><td>Class</td><td>0 – 20</td><td>20 – 40</td><td>40 – 60</td><td>60 – 80</td><td>80 – 100</td><td>Total</td></tr><tr><td>f</td><td>17</td><td><math>f_1</math></td><td>32</td><td><math>f_2</math></td><td>19</td><td>100</td></tr></table> <b>Answer: <math>f_1 = 18</math>, <math>f_2 = 14</math></b>	Class	0 – 20	20 – 40	40 – 60	60 – 80	80 – 100	Total	f	17	$f_1$	32	$f_2$	19	100
Class	0 – 20	20 – 40	40 – 60	60 – 80	80 – 100	Total										
f	17	$f_1$	32	$f_2$	19	100										
C	5	A co-operative bank has two branches employing 50 and 70 workers respectively. The average salaries paid by two respective branches are 360 and 390 rupees per month. Calculate the mean of the salaries of all the employees.  <b>Answer: 377.5</b>														

## Unit 4 – Basic Statistics

### 1.2 Median

- The median is the value found at the **exact middle** of the set of values.
- Median is denoted by capital letter “**M**”.
- To compute the median, list all observations in ascending order and then locate the value in the center of the sample.
- Table of formula of median.

Data	Formula
Ungrouped Data	If n is <b>odd</b> , then $M = \left( \frac{n + 1}{2} \right)^{\text{th}} \text{ observation}$
Discrete Grouped Data	If n is <b>even</b> , then $M = \frac{\left( \frac{n}{2} \right)^{\text{th}} \text{ observation} + \left( \frac{n}{2} + 1 \right)^{\text{th}} \text{ observation}}{2}$
Continuous Grouped Data	$M = L + \left( \frac{\frac{n}{2} - F}{f} \right) \times c$

Where,

Median class = Class whose cumulative frequency with property  $\min \left\{ cf \mid cf \geq \frac{n}{2} \right\}$

L = Lower boundary point of the median class

n = Total number of observation (sum of the frequencies)

F = Cumulative frequency of the class preceding the median class

f = The frequency of the median class

Example of Method-1.2: Median

C	1	Find the median of following data: 20, 25, 30, 15, 17, 35, 26, 18, 40, 45, 50.  <b>Answer: 26</b>														
C	2	The given observations have been arranged in ascending order. If the median of the data is 63, find the value of x for the following data: 29, 32, 48, 50, x, x + 2, 72, 78, 84, 95.  <b>Answer: x = 62</b>														
C	3	Calculate the median for the following data: <table border="1"><tr><td>Marks</td><td>20</td><td>9</td><td>25</td><td>50</td><td>40</td><td>80</td></tr><tr><td>No. of students</td><td>6</td><td>4</td><td>16</td><td>7</td><td>8</td><td>2</td></tr></table> <b>Answer: 25</b>	Marks	20	9	25	50	40	80	No. of students	6	4	16	7	8	2
Marks	20	9	25	50	40	80										
No. of students	6	4	16	7	8	2										
C	4	The following table gives marks obtained by 50 students in statistics. Find the median. <table border="1"><tr><td>Marks</td><td>0 – 10</td><td>10 – 20</td><td>20 – 30</td><td>30 – 40</td><td>40 – 50</td></tr><tr><td>No. of students</td><td>16</td><td>12</td><td>18</td><td>3</td><td>1</td></tr></table> <b>Answer: 17.5</b>	Marks	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50	No. of students	16	12	18	3	1		
Marks	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50											
No. of students	16	12	18	3	1											
C	5	The median of 60 observations (following data) is 28.5. Find x and y. <table border="1"><tr><td>Marks</td><td>0 – 10</td><td>10 – 20</td><td>20 – 30</td><td>30 – 40</td><td>40 – 50</td><td>50 – 60</td></tr><tr><td>No. of students</td><td>5</td><td>x</td><td>20</td><td>15</td><td>y</td><td>5</td></tr></table> <b>Answer: x = 8,      y = 7</b>	Marks	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50	50 – 60	No. of students	5	x	20	15	y	5
Marks	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50	50 – 60										
No. of students	5	x	20	15	y	5										

## Unit 4 – Basic Statistics

### 1.3 Mode

- The mode is the **most frequently** occurring value in the set.
- Mode is denoted by capital letter “**Z**”.
- The mode is not necessarily unique, like mean and median. we can have data with two modes (bi-modal) or more than two modes (multi-modal).
- Table of formula of mode.

Data	Formula
<b>Ungrouped Data</b>	Most repeated observation among given data
<b>Discrete Grouped Data</b>	Highest frequency among given data
<b>Continuous Grouped Data</b>	$Z = L + \left( \frac{f_1 - f_0}{2f_1 - f_0 - f_2} \right) \times c$

Where,

Modal class = A class with highest frequency

L = Lower boundary of modal class

c = Class length

$f_1$  = Frequency of the modal class

$f_0$  = Frequency of the class before the modal class

$f_2$  = Frequency of the class after the modal class

### Relation Between Mean, Median and Mode

- $Z = 3M - 2\bar{x}$ ; where  $\bar{x}$  = Mean, M = Median, Z = Mode



Example of Method-1.3: Mode

C	1	If mean is 16 and median is 20. Calculate the mode.  <b>Answer: 28</b>																								
C	2	Find the mode of following data: (a) 2, 4, 2, 5, 7, 2, 8, 9. (b) 2, 8, 4, 6, 10, 12, 4, 8, 14, 16.  <b>Answer: (a) 2, (b) 4 &amp; 8</b>																								
C	3	Find the mode of following data: <table border="1"><tr><td>x</td><td>11</td><td>22</td><td>33</td><td>44</td></tr><tr><td>f</td><td>15</td><td>20</td><td>19</td><td>10</td></tr></table> <b>Answer: 22</b>	x	11	22	33	44	f	15	20	19	10														
x	11	22	33	44																						
f	15	20	19	10																						
C	4	Find the mode of following data: <table border="1"><tr><td>Class</td><td>0 – 10</td><td>10 – 20</td><td>20 – 30</td><td>30 – 40</td><td>40 – 50</td></tr><tr><td>f</td><td>3</td><td>5</td><td>7</td><td>10</td><td>12</td></tr><tr><td></td><td>50 – 60</td><td>60 – 70</td><td>70 – 80</td><td>80 – 90</td><td>90 – 100</td></tr><tr><td></td><td>15</td><td>12</td><td>6</td><td>2</td><td>8</td></tr></table> <b>Answer: 55</b>	Class	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50	f	3	5	7	10	12		50 – 60	60 – 70	70 – 80	80 – 90	90 – 100		15	12	6	2	8
Class	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50																					
f	3	5	7	10	12																					
	50 – 60	60 – 70	70 – 80	80 – 90	90 – 100																					
	15	12	6	2	8																					
C	5	Find the mode of following data: <table border="1"><tr><td>Class</td><td>200 – 220</td><td>220 – 240</td><td>240 – 260</td><td>260 – 280</td></tr><tr><td>f</td><td>7</td><td>15</td><td>21</td><td>19</td></tr><tr><td></td><td>280 – 300</td><td>300 – 320</td><td>320 – 340</td><td></td></tr><tr><td></td><td>6</td><td>4</td><td>2</td><td></td></tr></table> <b>Answer: 255</b>	Class	200 – 220	220 – 240	240 – 260	260 – 280	f	7	15	21	19		280 – 300	300 – 320	320 – 340			6	4	2					
Class	200 – 220	220 – 240	240 – 260	260 – 280																						
f	7	15	21	19																						
	280 – 300	300 – 320	320 – 340																							
	6	4	2																							
C	6	Obtain the mean, mode and median for the following information: <table border="1"><tr><td>x</td><td>&lt; 10</td><td>&lt; 20</td><td>&lt; 30</td><td>&lt; 40</td><td>&lt; 50</td><td>&lt; 60</td></tr><tr><td>f</td><td>12</td><td>30</td><td>57</td><td>77</td><td>94</td><td>100</td></tr></table> <b>Answer: <math>\bar{x}</math> = 28, M = 27.407, Z = 25.625</b>	x	< 10	< 20	< 30	< 40	< 50	< 60	f	12	30	57	77	94	100										
x	< 10	< 20	< 30	< 40	< 50	< 60																				
f	12	30	57	77	94	100																				

## Unit 4 – Basic Statistics

### 1.4 Quartiles

→ Quartiles are measures which divide a series into **four** equal parts using three quartiles namely  $Q_1$ ,  $Q_2$  and  $Q_3$ .

The quartile  $Q_1$  is known as first quartile or lower quartile,

The quartile  $Q_2$  is known as second quartile or **median** quartile,

The quartile  $Q_3$  is known as third quartile or upper quartile.

→ To compute the quartile, list all observations in ascending order.

→ Table of formula of quartile.

Data	Formula
<b>Ungrouped Data</b>	$Q_1 = \left(\frac{n+1}{4}\right)^{\text{th}}$ observation
	$Q_2 = \text{Median}$
<b>Discrete Grouped Data</b>	$Q_3 = 3\left(\frac{n+1}{4}\right)^{\text{th}}$ observation
<b>Continuous Grouped Data</b>	$Q_1 = L + \left(\frac{\frac{n}{4} - F}{f}\right) \times c, \quad Q_2 = L + \left(\frac{\frac{n}{2} - F}{f}\right) \times c$ $Q_3 = L + \left(\frac{\frac{3n}{4} - F}{f}\right) \times c$

Where,

Quartile class  $Q_k$

= Class whose cumulative frequency with property  $\min \left\{ cf \mid cf \geq \frac{kn}{4} \right\}$

$L$  = Lower boundary point of the quartile class

$n$  = Total number of observation (sum of the frequencies)

$F$  = Cumulative frequency of the class preceding the quartile class

$f$  = The frequency of the quartile class

Example of Method-1.4: Quartiles

C	1	Find the quartiles of the data: 4, 6, 7, 8, 10, 23, 34.  <b>Answer: 6, 8, 23</b>																				
C	2	Find the quartile $Q_1$ , and $Q_3$ . <table><tr><td>x</td><td>2</td><td>4</td><td>6</td><td>8</td><td>10</td><td>12</td></tr><tr><td>f</td><td>4</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table> <b>Answer: 4, 12</b>	x	2	4	6	8	10	12	f	4	1	2	3	4	5						
x	2	4	6	8	10	12																
f	4	1	2	3	4	5																
C	3	Compute $Q_1$ , and $Q_3$ for the data relating to age in years of 543 members in a village. <table><tr><td>x</td><td>20</td><td>30</td><td>40</td><td>50</td><td>60</td><td>70</td><td>80</td></tr><tr><td>f</td><td>3</td><td>61</td><td>132</td><td>153</td><td>140</td><td>51</td><td>3</td></tr></table> <b>Answer: 40, 60</b>	x	20	30	40	50	60	70	80	f	3	61	132	153	140	51	3				
x	20	30	40	50	60	70	80															
f	3	61	132	153	140	51	3															
C	4	Calculate the quartiles $Q_1$ , and $Q_3$ for the following data. <table><tr><td>Class</td><td>30 – 32</td><td>32 – 34</td><td>34 – 36</td><td>36 – 38</td></tr><tr><td>f</td><td>12</td><td>18</td><td>16</td><td>14</td></tr><tr><td></td><td>38 – 40</td><td>40 – 42</td><td>42 – 44</td><td></td></tr><tr><td></td><td>12</td><td>8</td><td>6</td><td></td></tr></table> <b>Answer: 33.06, 38.75</b>	Class	30 – 32	32 – 34	34 – 36	36 – 38	f	12	18	16	14		38 – 40	40 – 42	42 – 44			12	8	6	
Class	30 – 32	32 – 34	34 – 36	36 – 38																		
f	12	18	16	14																		
	38 – 40	40 – 42	42 – 44																			
	12	8	6																			
C	5	Calculate the quartiles $Q_1$ , and $Q_3$ for the following data. <table><tr><td>Marks</td><td>&lt; 20</td><td>20 – 30</td><td>30 – 40</td><td>40 &lt;</td></tr><tr><td>Number of Students</td><td>14</td><td>20</td><td>28</td><td>18</td></tr></table> <b>Answer: 23, 39.28</b>	Marks	< 20	20 – 30	30 – 40	40 <	Number of Students	14	20	28	18										
Marks	< 20	20 – 30	30 – 40	40 <																		
Number of Students	14	20	28	18																		

## Unit 4 – Basic Statistics

### 1.5 Percentiles

- Percentiles are measures which divide a series into **hundred** equal parts.
- There are ninety-nine percentiles which is
  - The percentile  $P_1$  is known as first percentile,
  - The percentile  $P_2$  is known as second percentile and so on.
- To compute the percentile, list all observations in ascending order.
- Table of formula of percentile.

Data	Formula
Ungrouped Data	$P_1 = \left( \frac{n+1}{100} \right)^{\text{th}}$ observation
Discrete Grouped Data	$P_2 = 2 \left( \frac{n+1}{100} \right)^{\text{th}}$ observation $\vdots$ $P_{99} = 99 \left( \frac{n+1}{100} \right)^{\text{th}}$ observation
Continuous Grouped Data	$P_1 = L + \left( \frac{\frac{n}{100} - F}{f} \right) \times c$ $P_2 = L + 2 \left( \frac{\frac{n}{100} - F}{f} \right) \times c$ $\vdots$ $P_{99} = L + 99 \left( \frac{\frac{n}{100} - F}{f} \right) \times c$

## Unit 4 – Basic Statistics

Where,

Percentile class  $P_k$

= Class whose cumulative frequency with property  $\min \left\{ cf \mid cf \geq \frac{kn}{100} \right\}$

$L$  = Lower boundary point of the percentile class

$n$  = Total number of observation (sum of the frequencies)

$F$  = Cumulative frequency of the class preceding the percentile class

$f$  = The frequency of the percentile class

### Example of Method-1.5: Percentiles

C	1	Find 20 <sup>th</sup> Percentile of the data: 4, 6, 7, 8, 10, 23, 34, 55, 60.  <b>Answer: 6</b>														
C	2	Find the 30 <sup>th</sup> Percentile of the following data. <table border="1"><tr><td>x</td><td>2</td><td>4</td><td>6</td><td>8</td><td>10</td><td>12</td></tr><tr><td>f</td><td>4</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr></table> <b>Answer: 6</b>	x	2	4	6	8	10	12	f	4	1	2	3	4	5
x	2	4	6	8	10	12										
f	4	1	2	3	4	5										
C	3	Find the 51 <sup>th</sup> Percentile of the following data: <table border="1"><tr><td>Marks</td><td>15 – 20</td><td>20 – 25</td><td>25 – 30</td><td>30 – 35</td><td>35 – 40</td></tr><tr><td>No. of students</td><td>2</td><td>12</td><td>15</td><td>20</td><td>25</td></tr></table> <b>Answer: 32.185</b>	Marks	15 – 20	20 – 25	25 – 30	30 – 35	35 – 40	No. of students	2	12	15	20	25		
Marks	15 – 20	20 – 25	25 – 30	30 – 35	35 – 40											
No. of students	2	12	15	20	25											

## Method 2 $\rightsquigarrow$ Measure of Dispersion

### Dispersion

- Dispersion refers to the **spread** of the values around the central tendency.
- For Example:  
–5, 0, 5 and – 50, 0, 50 both have the same mean 0 but clearly the data given in the second case much more widely dispersed than those in the first case.
- So, measures of central tendency are not sufficient for having some idea about dispersion.
- Measures of dispersion gives the idea about the degree to which numerical data tend to spread about an average life.
- There are certain measures of dispersion which is,
  - (1) Range
  - (2) Interquartile Range
  - (3) Standard Deviation
  - (4) Mean Deviation

### Range

- Range is simply the highest value **minus** the lowest value of a set of data values.
- For Example:  
Range of –5, 0, 5 is 10  
**Reason:** Range = Highest value – lowest value  
$$= 5 - (-5)$$
$$= 10$$

### Interquartile Range

- The difference between the upper and lower quartile is known as the interquartile range.  
i.e., Interquartile range = upper Quartile – lower Quartile =  $Q_3 - Q_1$ .

### Standard Deviation

- Standard deviation is a measure that is used to quantify the amount of variation or dispersion of a set of data values.
- It is denoted by " $\sigma$ " and read as "sigma".

## Unit 4 – Basic Statistics

→ Table of different formulae of standard deviation.

Method	Ungrouped Data	Discrete Grouped Data	Continuous Grouped Data
Direct Method	$\sqrt{\frac{\sum x_i^2}{n} - \left(\frac{\sum x_i}{n}\right)^2}$	$\sqrt{\frac{\sum f_i x_i^2}{\sum f_i} - \left(\frac{\sum f_i x_i}{\sum f_i}\right)^2}$	
Assumed Mean Method	$\sqrt{\frac{\sum d_i^2}{n} - \left(\frac{\sum d_i}{n}\right)^2}$	$\sqrt{\frac{\sum f_i d_i^2}{\sum f_i} - \left(\frac{\sum f_i d_i}{\sum f_i}\right)^2}$	
Step Deviation Method	-----	-----	$\sqrt{\frac{\sum f_i u_i^2}{\sum f_i} - \left(\frac{\sum f_i u_i}{\sum f_i}\right)^2} \times c$

### Variance

- Variance is **expectation** of the squared deviation.
- It informally measures how far a set of (random) numbers are spread out from their mean.
- It is denoted by capital letter “**V**” and defined as  $V = \sigma^2$ .

### Coefficient of Variation

- The Coefficient of Variation is the **ratio** of the standard deviation to the mean and shows the extent of variability in relation to the mean of the population.
- Coefficient of Variance is defined as

$$C.V. = \frac{\sigma}{\bar{x}} \times 100$$

- If C.V. is high, then it is less consistent. Similarly, if C.V. is less, then it is more consistent.
- The higher the Coefficient of Variation, the greater the dispersion.

### Mean Deviation

- The mean deviation is defined as a statistical measure that is used to calculate the average deviation from the mean value of the given data set.
- In a simple word, the mean deviation is used to calculate how far the values fall from the middle of the data set.

→ Table of different formulae of mean deviation:

Method	Ungrouped Data	Grouped Data
M.D. about Mean	$\frac{\sum  x_i - \bar{x} }{n}$	$\frac{\sum f_i  x_i - \bar{x} }{\sum f_i}$
M.D. about Median	$\frac{\sum  x_i - M }{n}$	$\frac{\sum f_i  x_i - M }{\sum f_i}$
M.D. about Mode	$\frac{\sum  x_i - Z }{n}$	$\frac{\sum f_i  x_i - Z }{\sum f_i}$

### Example of Method-3: Measure of Dispersion

C	1	Determine the interquartile range value for 2, 3, 5, 7, 11, 13, 17, 19, 23, 29.  <b>Answer: 11</b>																
C	2	Find the interquartile range for the following distribution: <table border="1"><tr><td>x</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr><tr><td>f</td><td>5</td><td>10</td><td>12</td><td>5</td><td>3</td></tr></table> <b>Answer: 10</b>	x	1	2	3	4	5	f	5	10	12	5	3				
x	1	2	3	4	5													
f	5	10	12	5	3													
C	3	Find the standard deviation for the following data: 6, 7, 10, 12, 13, 4, 8, 12.  <b>Answer: 3.0414</b>																
C	4	Find the standard deviation and variance for the following distribution: <table border="1"><tr><td>x</td><td>0 – 10</td><td>10 – 20</td><td>20 – 30</td><td>30 – 40</td><td>40 – 50</td><td>50 – 60</td><td>60 – 70</td></tr><tr><td>f</td><td>6</td><td>14</td><td>10</td><td>8</td><td>1</td><td>3</td><td>8</td></tr></table> <b>Answer: <math>\sigma = 19.6214</math>,      <math>V = 384.9993</math></b>	x	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50	50 – 60	60 – 70	f	6	14	10	8	1	3	8
x	0 – 10	10 – 20	20 – 30	30 – 40	40 – 50	50 – 60	60 – 70											
f	6	14	10	8	1	3	8											



## Unit 4 – Basic Statistics

C	5	<p>The arithmetic means of runs scored by three batsmen A, B and C, in the same series of 10 innings, are 50, 48 and 12 respectively. The standard deviations of their runs are 15, 12 and 2 respectively. Who is the most consistent of the three?</p> <p><b>Answer: Batsman C is more consistent.</b></p>																						
C	6	<p>Two machines A, B are used to fill a mixture of cement concrete in a beam. Find the standard deviation of each machine &amp; comment on the performances of two machines.</p> <table><tr><td>A</td><td>32</td><td>28</td><td>47</td><td>63</td><td>71</td><td>39</td><td>10</td><td>60</td><td>96</td><td>14</td></tr><tr><td>B</td><td>19</td><td>31</td><td>48</td><td>53</td><td>67</td><td>90</td><td>10</td><td>62</td><td>40</td><td>80</td></tr></table> <p><b>Answer: <math>\sigma_A = 25.4950</math>, <math>\sigma_B = 24.4290</math></b></p> <p><b>There is less variability in the performance of the machine B.</b></p>	A	32	28	47	63	71	39	10	60	96	14	B	19	31	48	53	67	90	10	62	40	80
A	32	28	47	63	71	39	10	60	96	14														
B	19	31	48	53	67	90	10	62	40	80														
C	7	<p>Find mean deviation about the mean, median and mode for the following data:</p> <table><tr><td>x</td><td>5</td><td>10</td><td>15</td><td>20</td><td>25</td></tr><tr><td>f</td><td>7</td><td>4</td><td>6</td><td>3</td><td>5</td></tr></table> <p><b>Answer: <math>MD(\bar{x}) = 6.32</math>, <math>MD(M) = 6.2</math>, <math>MD(Z) = 9</math></b></p>	x	5	10	15	20	25	f	7	4	6	3	5										
x	5	10	15	20	25																			
f	7	4	6	3	5																			
C	8	<p>Find mean deviation about the mean, median and mode for the following data:</p> <table><tr><td>Class</td><td>5 – 25</td><td>25 – 45</td><td>45 – 65</td><td>65 – 85</td><td>85 – 105</td></tr><tr><td>f</td><td>12</td><td>8</td><td>14</td><td>20</td><td>6</td></tr></table> <p><b>Answer: <math>MD(\bar{x}) = 21.33</math>, <math>MD(M) = 21.904</math>, <math>MD(Z) = 23.466</math></b></p>	Class	5 – 25	25 – 45	45 – 65	65 – 85	85 – 105	f	12	8	14	20	6										
Class	5 – 25	25 – 45	45 – 65	65 – 85	85 – 105																			
f	12	8	14	20	6																			

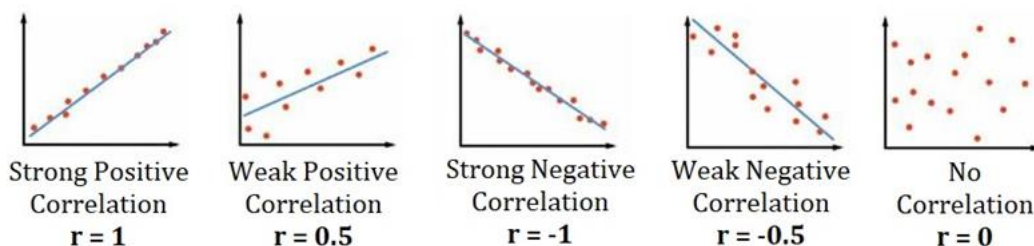
## Unit 4 – Basic Statistics

### Measure of Association

- In statistics, various factors or coefficients used to quantify a relationship between two or more variables.
- Some measures of association are Pearson's correlation coefficient, the Spearman rank correlation coefficient.
- Measures of association are used in various fields of research.
- A measure of association determined by correlation analysis and regression analysis.
  - Correlation and regression are the most commonly used techniques for investigating the relationship between two quantitative variables.
  - Correlation refers to the relationship of two or more variables. Regression establishes a functional relationship between the variables.
  - The coefficient of correlation is a relative measure whereas the regression coefficient is an absolute figure.

### Correlation

- Two variables are known as **correlated** if a change in one variable affects a change in the other variable. Such a data connecting two variables is called bivariate data.
- For Example:
  - Relationship between heights and weights.
  - Relationship between price and demand of commodity.
  - Relationship between age of husband and age of wife.
- When two variables are correlated with each other, it is important to know the amount or extent of correlation between them.
- The numerical measure of correlation or degree of relationship existing between two variables is known as the coefficient of correlation.
- It is denoted by **r** and it is always lying between  $-1$  and  $1$ .



## Unit 4 – Basic Statistics

- The value of  $r$  is  $\pm 0.9$  or  $\pm 0.8$  etc. shows high degree of relationship between the variables while  $\pm 0.2$  or  $\pm 0.1$  etc. shows low degree of correlation.

### Types of Correlation

- Correlation is classified into four types.

- Positive Correlation
- Negative Correlation
- Linear Correlation
- Nonlinear Correlation

### Positive Correlation

- If both the variables vary in same direction, then such correlation is known as **positive correlation**.
- In other words, if the value of one variable increases, the value of other variables also increases, or, if the value of one variable decreases, the value of other variables also decreases.
- For Example:
- The correlation between heights and weights of group of persons is a positive correlation.

<b>Height(cm)</b>	150	152	155	160	162	165
<b>Weight(kg)</b>	60	62	64	65	67	69

### Negative Correlation

- If both the variables vary in opposite direction, then such correlation is known as **negative correlation**.
- In other words, if the value of one variable increases, the value of other variables decreases, or, if the value of one variable decreases, the value of other variables increases.
- For Example:
- The correlation between the price and demand of a commodity is a negative correlation.

<b>Price (₹ per unit)</b>	10	8	6	5	4	1
<b>Demand(units)</b>	100	200	300	400	500	600

## Unit 4 – Basic Statistics

### Linear Correlation

- If the ratio of change between two variables is constant, then such correlation is known as **linear correlation**.
- If such variables are plotted on a graph paper, a straight line is obtained.
- For Example:

<b>Milk (l)</b>	5	10	15	20	25	30
<b>Curd (kg)</b>	2	4	6	8	10	12

### Nonlinear Correlation

- If the ratio of change between two variables is not constant, then such correlation known as **nonlinear correlation**.
- If such variables are plotted on a graph paper, a curve is obtained.
- For Example:

<b>Advertising expenses (₹ in lacs)</b>	3	6	9	12	15
<b>Curd (kg)</b>	10	12	15	15	16

### Methods of Studying Correlation

- There are two different methods of studying correlation:
  - Graphical Methods
    - (1) Scatter Diagram
    - (2) Simple Graph
  - Mathematical Methods
    - (1) Karl Pearson's coefficient of correlation
    - (2) Spearman's rank coefficient of correlation

### Method 3 $\rightsquigarrow$ Covariance

#### Covariance

→ Covariance is a measure of the relationship between two random variables X and Y.

→ It is defined as

$$\text{cov}(X, Y) = \frac{1}{n} \cdot \sum (x - \bar{x})(y - \bar{y})$$

Where,  $\bar{x}$  and  $\bar{y}$  are the mean of the X and Y respectively.

→ Also, we can find covariance using following formula:

$$\text{cov}(X, Y) = \frac{1}{n} \cdot \sum xy - \frac{1}{n^2} \cdot \sum x \sum y$$

#### Example of Method-3: Covariance

C	1	Find the covariance of the following data: (15,44), (20,43), (25,45), (30,37), (40,34), (50,37).  <b>Answer: – 39.17</b>												
C	2	Determine $\sum xy$ if $n = 5$ , $\text{cov}(X,Y) = 7.2$ , $\sum x = 25$ , $\sum y = 30$  <b>Answer: 186</b>												
C	3	Compute the covariance between x and y using the following data: <table border="1"><tr><td>x</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr><tr><td>y</td><td>3</td><td>2</td><td>5</td><td>4</td><td>6</td></tr></table> <b>Answer: 1.6</b>	x	1	2	3	4	5	y	3	2	5	4	6
x	1	2	3	4	5									
y	3	2	5	4	6									

## Method 4 $\rightarrow$ Correlation Coefficient

### Karl Pearson's Coefficient of Correlation

$\rightarrow$  The coefficient of correlation is the measure of correlation between two random variables X and Y, and is denoted by r. It is defined as below:

$$r = \frac{\text{cov}(X, Y)}{\sigma_x \sigma_y} \quad \dots \dots (1)$$

Where,

cov(X, Y) is the covariance of variables X and Y.

$\sigma_x$  &  $\sigma_y$  are standard deviation of X and Y respectively.

$\rightarrow$  This expression is known as Karl Pearson's coefficient of correlation.

$\rightarrow$  We have,

$$\text{cov}(X, Y) = \frac{1}{n} \cdot \sum (x - \bar{x})(y - \bar{y})$$

$$\sigma_x = \sqrt{\frac{\sum (x - \bar{x})^2}{n}}$$

$$\sigma_y = \sqrt{\frac{\sum (y - \bar{y})^2}{n}}$$

$\rightarrow$  Substitute the values of cov(X, Y),  $\sigma_x$  and  $\sigma_y$  in equation (1), We get

$$r = \frac{\sum (x - \bar{x})(y - \bar{y})}{\sqrt{\sum (x - \bar{x})^2 \cdot \sum (y - \bar{y})^2}} \quad \dots \dots (2)$$

$\rightarrow$  Equation (2) can be further reduced to below equation.

$$r = \frac{n \cdot \sum xy - \sum x \sum y}{\sqrt{n \cdot \sum x^2 - (\sum x)^2} \sqrt{n \cdot \sum y^2 - (\sum y)^2}} \quad \dots \dots (3)$$

Example of Method-4: Correlation Coefficient

C	1	Compute the coefficient of correlation between x and y using the following data: <table><tr><td>x</td><td>2</td><td>4</td><td>5</td><td>6</td><td>8</td><td>11</td></tr><tr><td>y</td><td>18</td><td>12</td><td>10</td><td>8</td><td>7</td><td>5</td></tr></table> <b>Answer: <math>r = -0.9203</math></b>	x	2	4	5	6	8	11	y	18	12	10	8	7	5							
x	2	4	5	6	8	11																	
y	18	12	10	8	7	5																	
C	2	Calculate Karl Pearson's correlation coefficient between age and playing habits: <table><tr><td>Age</td><td>20</td><td>21</td><td>22</td><td>23</td><td>24</td><td>25</td></tr><tr><td>No. of students</td><td>500</td><td>400</td><td>300</td><td>240</td><td>200</td><td>160</td></tr><tr><td>Regular players</td><td>400</td><td>300</td><td>180</td><td>96</td><td>60</td><td>24</td></tr></table> <b>Answer: <math>r = -0.9823</math></b>	Age	20	21	22	23	24	25	No. of students	500	400	300	240	200	160	Regular players	400	300	180	96	60	24
Age	20	21	22	23	24	25																	
No. of students	500	400	300	240	200	160																	
Regular players	400	300	180	96	60	24																	
C	3	Determine the coefficient of correlation if $n = 10, \bar{x} = 5.5, \bar{y} = 4,$ $\sum x^2 = 385, \sum y^2 = 192, \quad \sum (x + y)^2 = 947.$  <b>Answer: <math>r = -0.6812</math></b>																					
C	4	Given that $n = 25, \sum x = 125, \sum x^2 = 650, \sum y = 100, \sum y^2 = 460$ and $\sum xy = 508$ . Later on, it was found that two of the points (8, 12)  <b>and (6, 8) were wrongly entered as (6, 14) and (8, 6). Prove that <math>r = \frac{2}{3}</math>.</b>																					

## Method 5 $\Rightarrow$ Rank Correlation Coefficient

### Rank Correlation

- Let a group of  $n$  individuals be arranged in order of merit with respect to some characteristics. The same group would give a different order(rank) for different characteristics.
- Considering the orders corresponding to two characteristics A and B, the correlation between these  $n$  pairs of ranks is known as **rank correlation** in the characteristics A and B for that group of individuals.
- It is denoted by  $\rho$ .

### Spearman's Rank Correlation Coefficient

- Edward Spearman's formula for rank correlation coefficient ( $\rho$ ) is given by

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

Where,  $d$  = Difference of Ranks

### Tied Rank

- If there is a tie between the ranks, then it is known as tied rank.
- Formula for rank correlation coefficient ( $\rho$ ) is,

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

Where,  $d$  = Difference of Ranks

$m_i$  = number of times a data repeats = 2, 3, ...

- Note:
  - If any data  $x_i$  is repeated 2 times, then  $m_1 = 2$ .
  - If any data  $x_j$  is repeated 3 times, then  $m_2 = 3$ .
- In case of tie between individuals' ranks, **the rank is divided among equal individuals**.
- For Example:
  - If there is tie with two items at 4<sup>th</sup> rank, then give average rank 4.5 as rank to both items.

$$\text{Average} = \frac{4 + 5}{2} = 4.5$$



### Example of Method-5: Rank Correlation Coefficient

C	1	<p>In a college, IT department has arranged one competition for IT students to develop an efficient program to solve a problem. Ten students took part in the competition and ranked by two judges given in the following table. Find the degree of agreement between the two judges using rank correlation coefficient.</p> <table><tr><td>1st judge</td><td>3</td><td>5</td><td>8</td><td>4</td><td>7</td><td>10</td><td>2</td><td>1</td><td>6</td><td>9</td></tr><tr><td>2nd judge</td><td>6</td><td>4</td><td>9</td><td>8</td><td>1</td><td>2</td><td>3</td><td>10</td><td>5</td><td>7</td></tr></table> <p><b>Answer: <math>\rho = -0.2970</math></b></p>	1st judge	3	5	8	4	7	10	2	1	6	9	2nd judge	6	4	9	8	1	2	3	10	5	7											
1st judge	3	5	8	4	7	10	2	1	6	9																									
2nd judge	6	4	9	8	1	2	3	10	5	7																									
C	2	<p>The competitions in a beauty contest are ranked by three judges:</p> <table><tr><td>1<sup>st</sup> judge</td><td>1</td><td>5</td><td>4</td><td>8</td><td>9</td><td>6</td><td>10</td><td>7</td><td>3</td><td>2</td></tr><tr><td>2<sup>nd</sup> judge</td><td>4</td><td>8</td><td>7</td><td>6</td><td>5</td><td>9</td><td>10</td><td>3</td><td>2</td><td>1</td></tr><tr><td>3<sup>rd</sup> judge</td><td>6</td><td>7</td><td>8</td><td>1</td><td>5</td><td>10</td><td>9</td><td>2</td><td>3</td><td>4</td></tr></table> <p>Use rank correlation to discuss which pair of judges has nearest approach to beauty.</p> <p><b>Answer: 2<sup>nd</sup> and 3<sup>rd</sup> judges has nearest approach</b></p> <p><b>[ <math>\rho_{12} = 0.5515, \rho_{23} = 0.7333, \rho_{13} = 0.0545</math> ]</b></p>	1 <sup>st</sup> judge	1	5	4	8	9	6	10	7	3	2	2 <sup>nd</sup> judge	4	8	7	6	5	9	10	3	2	1	3 <sup>rd</sup> judge	6	7	8	1	5	10	9	2	3	4
1 <sup>st</sup> judge	1	5	4	8	9	6	10	7	3	2																									
2 <sup>nd</sup> judge	4	8	7	6	5	9	10	3	2	1																									
3 <sup>rd</sup> judge	6	7	8	1	5	10	9	2	3	4																									
C	3	<p>Find the rank correlation coefficient and comment on its value:</p> <table><tr><td>Roll no.</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td><td>6</td><td>7</td><td>8</td><td>9</td></tr><tr><td>Marks in Math.</td><td>78</td><td>36</td><td>98</td><td>25</td><td>75</td><td>82</td><td>90</td><td>62</td><td>65</td></tr><tr><td>Marks in Chem.</td><td>84</td><td>51</td><td>91</td><td>60</td><td>68</td><td>62</td><td>86</td><td>58</td><td>53</td></tr></table> <p><b>Answer: <math>\rho = 0.8333</math></b></p>	Roll no.	1	2	3	4	5	6	7	8	9	Marks in Math.	78	36	98	25	75	82	90	62	65	Marks in Chem.	84	51	91	60	68	62	86	58	53			
Roll no.	1	2	3	4	5	6	7	8	9																										
Marks in Math.	78	36	98	25	75	82	90	62	65																										
Marks in Chem.	84	51	91	60	68	62	86	58	53																										
C	4	<p>Calculate coefficient of correlation by Spearman's method from following.</p> <table><tr><td>Sales</td><td>45</td><td>56</td><td>39</td><td>54</td><td>45</td><td>40</td><td>56</td><td>60</td><td>30</td><td>36</td></tr><tr><td>Cost</td><td>40</td><td>36</td><td>30</td><td>44</td><td>36</td><td>32</td><td>45</td><td>42</td><td>20</td><td>36</td></tr></table> <p><b>Answer: <math>\rho = 0.7636</math></b></p>	Sales	45	56	39	54	45	40	56	60	30	36	Cost	40	36	30	44	36	32	45	42	20	36											
Sales	45	56	39	54	45	40	56	60	30	36																									
Cost	40	36	30	44	36	32	45	42	20	36																									

## Unit 4 – Basic Statistics

---

<b>C</b>	<b>5</b>	<p>The coefficient of rank correlation of marks obtained by 10 students in English and Economics was found to be 0.6. It was later discovered that the difference in ranks in the two subjects obtained by one of the students was wrongly taken as 7 instead of 1. Find the correct coefficient of rank correlation.</p> <p><b>Answer: 0.8909</b></p>
----------	----------	--

## Unit 4 – Basic Statistics

---

### Regression

- Regression is defined as a method of estimating the value of one variable when the other is known and both are correlated.
- We use the general form regression line for these algebraic expressions. The algebraic expressions of the regression lines are known as **Regression Equations**.
- It is highly used in statistical estimation of demand curve, supply curve, production function, cost function, consumption function etc.

### Types of Regression

- Regression is classified into four types:
  - Simple Regression
  - Multiple Regression
  - Linear Regression
  - Nonlinear Regression

### Simple Regression

- The regression analysis for studying only two variables at a time is known as **simple regression**.

### Multiple Regression

- The regression analysis for studying more than two variables at a time is known as **multiple regression**.

### Linear Regression

- If the regression curve is a straight line, the regression is known as **linear regression**.

### Nonlinear Regression

- If the regression curve is not a straight line, the regression is known as **nonlinear regression**.

### Method of Studying Regression

- There are two methods of studying regression:
  - Method of scatter diagram
  - Method of least square
- We will use method of least square only to find out regression.

## Method 6 $\rightarrow$ Linear Regression

### Line of Regression (Linear Regression)

- $\rightarrow$  If the variables, which are highly correlated, are plotted on a graph then the points are around a straight line, the line is known as the **line of regression**.
- $\rightarrow$  There are two types of line of regression.
- Line of regression of y on x
  - Line of regression of x on y

### Line of Regression of y on x

- $\rightarrow$  It is the line which gives the **best estimate for the values of y** for given values of x.
- $\rightarrow$  The regression equation of y on x is given by

$$y - \bar{y} = b_{yx} (x - \bar{x})$$

$$\begin{aligned} \text{Where, } b_{yx} = \text{Regression Coefficient} &= r \frac{\sigma_y}{\sigma_x} = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum x^2 - (\sum x)^2} \\ &= \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (x - \bar{x})^2} \end{aligned}$$

$r$  = Correlation coefficient between x and y

$\bar{x}, \sigma_x$  = Mean & Standard Deviation of all  $x_i$

$\bar{y}, \sigma_y$  = Mean & Standard Deviation of all  $y_i$

### Line of Regression of x on y

- $\rightarrow$  It is the line which gives the **best estimate for the values of x** for given values of y.
- $\rightarrow$  The regression equation of x on y is given by

$$x - \bar{x} = b_{xy} (y - \bar{y})$$

$$\begin{aligned} \text{Where, } b_{xy} = \text{Regression Coefficient} &= r \frac{\sigma_x}{\sigma_y} = \frac{n \sum xy - (\sum x)(\sum y)}{n \sum y^2 - (\sum y)^2} \\ &= \frac{\sum (x - \bar{x})(y - \bar{y})}{\sum (y - \bar{y})^2} \end{aligned}$$

$r$  = Correlation coefficient between x and y

$\bar{x}, \sigma_x$  = Mean & Standard Deviation of all  $x_i$

$\bar{y}, \sigma_y$  = Mean & Standard Deviation of all  $y_i$

## Unit 4 – Basic Statistics

### Properties of Regression Coefficients

→ The coefficient of correlation is the geometric mean of the coefficients of regression.

$$\text{i.e., } r = \sqrt{b_{yx} \cdot b_{xy}}$$

→ The product of both  $b_{xy}$  and  $b_{yx}$  cannot be more than 1.

→ Both the regression coefficients will have the same sign. They are either both positive and both negative. It means,

If  $r < 0$ , then  $b_{yx} < 0$  &  $b_{xy} < 0$ .

If  $r > 0$ , then  $b_{yx} > 0$  &  $b_{xy} > 0$ .

→ The arithmetic mean of the regression coefficients is greater than or equal to the correlation coefficient.

$$\text{i.e., } \frac{b_{xy} + b_{yx}}{2} \geq r$$

### Example of Method-6: Linear Regression

C	1	Obtain the two lines of regression for the following data:																		
		<table><tr><td>Sales (No. of tablets)</td><td>190</td><td>240</td><td>250</td><td>300</td><td>310</td><td>335</td><td>300</td></tr><tr><td>Advertising expense (Rs.)</td><td>5</td><td>10</td><td>12</td><td>20</td><td>20</td><td>30</td><td>30</td></tr></table> <p><b>Answer: <math>y = 0.1766x - 30.4221</math> ; <math>x = 4.7357y + 189.0807</math></b></p>	Sales (No. of tablets)	190	240	250	300	310	335	300	Advertising expense (Rs.)	5	10	12	20	20	30	30		
Sales (No. of tablets)	190	240	250	300	310	335	300													
Advertising expense (Rs.)	5	10	12	20	20	30	30													
C	2	A study of amount of rainfall and quantity of air pollution removed is:																		
		<table><tr><td>Daily rainfall (0.01 cm)</td><td>4.3</td><td>4.5</td><td>5.9</td><td>5.6</td><td>6.1</td><td>5.2</td><td>3.8</td><td>2.1</td><td>7.5</td></tr><tr><td>Particulate removed unit</td><td>126</td><td>121</td><td>116</td><td>118</td><td>114</td><td>118</td><td>132</td><td>141</td><td>108</td></tr></table> <p>a. Find the equation of the regression line to predict the particulate removed from the amount of daily rainfall.</p> <p>b. Find the amount of particulate removed when daily rainfall is 4.8 units.</p> <p><b>Answer: a. <math>y = -6.3240x + 153.1755</math> ; b. 122.8203</b></p>	Daily rainfall (0.01 cm)	4.3	4.5	5.9	5.6	6.1	5.2	3.8	2.1	7.5	Particulate removed unit	126	121	116	118	114	118	132
Daily rainfall (0.01 cm)	4.3	4.5	5.9	5.6	6.1	5.2	3.8	2.1	7.5											
Particulate removed unit	126	121	116	118	114	118	132	141	108											

C	3	<p>The following data regarding the height(y) and weight(x) of 100 students are given: <math>\sum x = 15000</math>, <math>\sum y = 6800</math>, <math>\sum x^2 = 2272500</math>, <math>\sum y^2 = 463025</math>, <math>\sum xy = 1022250</math>. Find the equation of regression line of height on weight.</p> <p><b>Answer: <math>y = 0.1x + 53</math></b></p>									
C	4	<p>The data for advertising and sale given below:</p> <table border="1"> <thead> <tr> <th></th><th>Adv. Exp.(x) (Rs. lakh)</th><th>Sales(y) (Rs. lakh)</th></tr> </thead> <tbody> <tr> <td>Mean</td><td>10</td><td>90</td></tr> <tr> <td>Standard deviation</td><td>3</td><td>12</td></tr> </tbody> </table> <p>a. Correlation coefficient between prices is 0.8.            b. Calculate the two regression lines.            c. Find the likely sales when advertising expenditure is 15 lakhs.            d. What should be the advertising expenditure if the company wants to attain a sales target of 120 lakhs?</p> <p><b>Answer: b. <math>x = 0.2y - 8</math> ; <math>y = 3.2x + 58</math> ; c. 106 ; d. 16</b></p>		Adv. Exp.(x) (Rs. lakh)	Sales(y) (Rs. lakh)	Mean	10	90	Standard deviation	3	12
	Adv. Exp.(x) (Rs. lakh)	Sales(y) (Rs. lakh)									
Mean	10	90									
Standard deviation	3	12									
C	5	<p>A study of prices of a certain commodity at Raipur and Kanpur yields the below data:</p> <table border="1"> <thead> <tr> <th></th><th>Raipur (Rs)</th><th>Kanpur (Rs)</th></tr> </thead> <tbody> <tr> <td>Average price/kg</td><td>2.463</td><td>2.797</td></tr> <tr> <td>Standard deviation</td><td>0.326</td><td>0.207</td></tr> </tbody> </table> <p>Correlation coefficient between prices at Raipur and Kanpur is 0.774. Estimate the most likely price at Raipur corresponding to the price of 3.052 per kilo at Kanpur.</p> <p><b>Answer: 2.774</b></p>		Raipur (Rs)	Kanpur (Rs)	Average price/kg	2.463	2.797	Standard deviation	0.326	0.207
	Raipur (Rs)	Kanpur (Rs)									
Average price/kg	2.463	2.797									
Standard deviation	0.326	0.207									

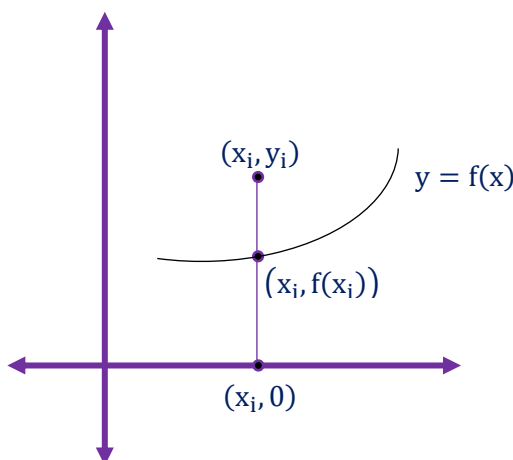
## Method 7 $\rightsquigarrow$ Curve Fitting

### Introduction

- We come across many situations where we often require to find a relationship between two or more variables. For example, weight and height of a person, demand and supply, expenditure depends on income, etc.
- This relation may be expressed by polynomial or exponential or logarithmic relationship. In order to determine such relationship, first it is requiring to collect the data showing corresponding values of the variables under consideration.
- Suppose  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  be the data showing corresponding values of the variables  $x$  and  $y$  under consideration. If we plot the above data points on a coordinate system, then the set of points so plotted form a scatter diagram.
- From this diagram, it is sometimes possible to visualize a smooth curve approximating the data. Such a curve is known as an **approximating curve**.
- In particular, if the data approximate well to a straight line, we say that a linear relationship exists between the variables. It is quite possible that the relationship of the form  $y = f(x)$  between two variables  $x$  and  $y$ , giving the approximating curve and which fit the given data of  $x$  and  $y$ , is known as **curve fitting**.

### Least Square Method

- Suppose that the data points are  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ , where  $x$  is independent and  $y$  is dependent variable.
- Let the fitting curve  $f(x)$  has the following deviations (or errors or residuals) from each data points. i.e.,  $d_1 = y_1 - f(x_1)$ .
- These  $d_i = y_i - f(x_i)$  are known as deviation, error or residual. Its value may be positive, negative or zero.



- To give equal weightage to each error, we square each of these and form their sum.

## Unit 4 – Basic Statistics

$$\begin{aligned} \text{i.e., } D &= d_1^2 + d_2^2 + \dots + d_n^2 = \sum_{i=1}^n d_i^2 \\ &= \sum_{i=1}^n [y_i - f(x_i)]^2 \end{aligned}$$

### 7.1 Curve Fitting by Straight Line

→ Let,  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  be the set of  $n$  values and let the relation between  $x$  and  $y$  be  **$y = a + bx$** .

→ We have,

$$D = \sum_{i=1}^n [y_i - f(x_i)]^2 = \sum_{i=1}^n (y_i - a - bx_i)^2$$

→ If  $D = 0$ , then all the  $n$  points will lie on  $y = f(x)$ .

→ If  $D \neq 0$ ,  $f(x)$  is chosen such that  $D$  is minimum.

→ This will be minimum at,

$$\frac{\partial D}{\partial a} = 0 \Rightarrow -2 \sum_{i=1}^n (y_i - a - bx_i) = 0$$

$$\Rightarrow \sum_{i=1}^n y_i - a \sum_{i=1}^n 1 - b \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \sum_{i=1}^n y_i - na - b \sum_{i=1}^n x_i = 0$$

$$\Rightarrow \sum_{i=1}^n y_i = na + b \sum_{i=1}^n x_i$$

$$\text{Similarly, by } \frac{\partial D}{\partial b} = 0 \Rightarrow \sum_{i=1}^n x_i y_i = a \sum_{i=1}^n x_i + b \sum_{i=1}^n x_i^2$$

→ We obtain following **Normal Equations** for the best fitting straight line  $y = a + bx$ .

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

$$\sum xy = a \sum x + b \sum x^2$$

→ These equations can be solved simultaneously to give the best value of  $a$  and  $b$  such that straight line is the best fit to the data.



Example of Method-7.1: Curve Fitting by Straight Line

C	1	Fit a straight line for the given pairs of (x,y) which are (1, 5), (2, 7), (3, 9), (4, 10), (5, 11).  <b>Answer: <math>y = 3.9 + 1.5x</math></b>												
C	2	Fit a straight-line $y = ax + b$ to the following data: <table><tr><td>x</td><td>-2</td><td>-1</td><td>0</td><td>1</td><td>2</td></tr><tr><td>y</td><td>1</td><td>2</td><td>3</td><td>3</td><td>4</td></tr></table> <b>Answer: <math>y = 0.7x + 2.6</math></b>	x	-2	-1	0	1	2	y	1	2	3	3	4
x	-2	-1	0	1	2									
y	1	2	3	3	4									
C	3	By method of least squares, fit a linear relation of the form $P = a + bW$ to the following data, P is the pull required to lift a weight W. Also estimate P, when W is 150. <table><tr><td>P</td><td>50</td><td>70</td><td>100</td><td>120</td></tr><tr><td>W</td><td>12</td><td>15</td><td>21</td><td>25</td></tr></table> <b>Answer: <math>P = -11.8005 + 5.3041W</math> ; <math>P(150) = 783.8145</math></b>	P	50	70	100	120	W	12	15	21	25		
P	50	70	100	120										
W	12	15	21	25										

## Unit 4 – Basic Statistics

### 7.2 Curve Fitting by Parabola

→ Let,  $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$  be the set of  $n$  values and let the relation between  $x$  and  $y$  be  $y = a + bx + cx^2$ .

→ We have,

$$D = \sum_{i=1}^n [y_i - f(x_i)] = \sum_{i=1}^n (y_i - a - bx_i - cx_i^2)^2$$

→ If  $D = 0$ , then all the  $n$  points will lie on  $y = f(x)$ . If  $D \neq 0$ ,  $f(x)$  is chosen such that  $D$  is minimum.

→ Differentiating  $S$  with respect to  $a, b, c$  and equating with zero (as done while fitting a linear curve). We obtain following **Normal Equations** for the best fitting  $y = a + bx + cx^2$  curve (parabola) of second degree.

$$\sum y = na + b \sum x + c \sum x^2$$

$$\sum xy = a \sum x + b \sum x^2 + c \sum x^3$$

$$\sum x^2 y = a \sum x^2 + b \sum x^3 + c \sum x^4$$

### Example of Method-7.2: Curve Fitting by Parabola

C	1	<p>Fit a polynomial of degree two using least square method for the following experimental data. Also, estimate <math>y(2.4)</math>.</p> <table><tr><td>x</td><td>1</td><td>2</td><td>3</td><td>4</td><td>5</td></tr><tr><td>y</td><td>5</td><td>12</td><td>26</td><td>60</td><td>97</td></tr></table> <p><b>Answer: <math>y = 10.4 - 11.0857x + 5.7143x^2</math> ; <math>y(2.4) = 16.7087</math></b></p>	x	1	2	3	4	5	y	5	12	26	60	97				
x	1	2	3	4	5													
y	5	12	26	60	97													
C	2	<p>Fit a second - degree parabola <math>y = ax^2 + bx + c</math> to the following data:</p> <table><tr><td>x</td><td>-3</td><td>-2</td><td>-1</td><td>0</td><td>1</td><td>2</td><td>3</td></tr><tr><td>y</td><td>12</td><td>4</td><td>1</td><td>2</td><td>7</td><td>15</td><td>30</td></tr></table> <p><b>Answer: <math>y = 2.1190x^2 + 2.9286x + 1.6667</math></b></p>	x	-3	-2	-1	0	1	2	3	y	12	4	1	2	7	15	30
x	-3	-2	-1	0	1	2	3											
y	12	4	1	2	7	15	30											

**C****3**

Fit a relation of the form  $R = a + bV + cV^2$  to the following data, where V is the velocity in km/hr. and R is the resistance in km/quintal. Estimate R when  $V = 90$ .

V	20	40	60	80	100	120
R	5.5	9.1	14.9	22.8	33.3	46.0

**Answer:  $R = 4.35 + 0.0024V + 0.0029V^2$  ;  $R(90) = 28.0560$**

## Method 8 $\rightarrow$ Fitting of Trend by Moving Average Method

### Introduction

- $\rightarrow$  Moving average method is a simple device of reducing fluctuations and obtaining trend values with a fair degree of accuracy.
- $\rightarrow$  In this method the average value of number of years (or months, weeks or days) is taken as the trend value for the middle point of the period of moving average.
- $\rightarrow$  The process of averaging smoothest the curve and reduces the fluctuation.
- $\rightarrow$  The first thing to be decided in the method is the period of the moving average. What it means is to take a decision about the number of consecutive items whose average would be calculated each time.

### Moving Averages of Odd Number of Years

- $\rightarrow$  Steps to find moving averages of  $k$  years (where  $k$  is an odd number)
  - (1) Find the sum of **first  $k$**  observations and also find their average (by dividing  **$k$** ) and place it against the  $\left(\frac{k+1}{2}\right)^{\text{th}}$  observation.
  - (2) **Leave** the first observation, find average of next  **$k$**  observations and place it against the  $\left(\frac{k+3}{2}\right)^{\text{th}}$  observation.
  - (3) Repeat above steps till last observation of the data is used.
- $\rightarrow$  i.e.,

Year	Observations	3 yearly moving average
1980	$x_1$	
1981	$x_2$	$\frac{1}{3}(x_1 + x_2 + x_3)$
1982	$x_3$	$\frac{1}{3}(x_2 + x_3 + x_4)$
1983	$x_4$	$\frac{1}{3}(x_3 + x_4 + x_5)$
1983	$x_5$	

## Unit 4 – Basic Statistics

### Moving Averages of Even Number of Years

→ Steps to find moving averages of  $k$  years (where  $k$  is an even number)

- (1) Find the sum of first  $k$  observations and also find their average (by dividing  $k$ ) and place it against between  $\left(\frac{k}{2}\right)^{\text{th}}$  &  $\left(\frac{k}{2} + 1\right)^{\text{th}}$  observation.
- (2) **Leave** the first observation, find the average of **next  $k$**  observations and place it against between  $\left(\frac{k}{2} + 1\right)^{\text{th}}$  &  $\left(\frac{k}{2} + 2\right)^{\text{th}}$  observation.
- (3) Repeat above steps till last observation of the data is used and we get  $k$  yearly moving average.
- (4) Find  $k$  yearly centered moving average (by dividing  $2$ ), which is the trend value of moving average.

→ i.e.,

Year	Observations	4 yearly moving average	4 yearly centered moving average
1980	$x_1$		
1981	$x_2$		
		$\frac{1}{4}(x_1 + x_2 + x_3 + x_4) = A_1$	
1982	$x_3$		$\frac{1}{2}(A_1 + A_2)$
		$\frac{1}{4}(x_2 + x_3 + x_4 + x_5) = A_2$	
1983	$x_4$		$\frac{1}{2}(A_2 + A_3)$
		$\frac{1}{4}(x_3 + x_4 + x_5 + x_6) = A_3$	
1984	$x_5$		

Example of Method-8: Fitting of Trend by Moving Average Method

C	1	Calculate 3 yearly moving averages of the following data.									
		Years	1971	1972	1973	1974	1975				
		Value	11	8	12	3	4				
		<b>Answer: 10.333, 7.666, 6.333</b>									
C	2	Calculate 4 yearly moving averages of number of students studying in a higher secondary school in a particular village from the following data.									
		Years	1991	1992	1993	1994	1995	1996	1997	1998	
		No. of Students	36	43	43	34	44	54	34	24	
		<b>Answer: 40, 42.37, 42.62, 40.25</b>									
C	3	Calculate 5 yearly moving averages of the following data.									
		Years	1971	1972	1973	1974	1975	1976	1977	1978	1979
		Value	25	49	55	12	52	40	88	12	44
		<b>Answer: 38.6, 41.6, 49.4, 40.8, 47.2</b>									
C	4	Calculate 6 yearly moving averages of the following data.									
		Years	2001	2002	2003	2004	2005	2006	2007	2008	2009
		Value	124	120	135	140	145	158	162	170	175
		<b>Answer: 140.166, 147.5, 155</b>									

\*\*\*\*\* End of the Unit \*\*\*\*\*