

Analysis of Variance

Structure

- 11.1 Introduction
 - Objectives
- 11.2 Problem Description
- 11.3 One-Way Analysis of Variance
- 11.4 Two-Way Analysis of Variance
- 11.5 Two-Way Analysis of Variance with m Observations per Cell

11.1 INTRODUCTION

In Block 2 of MST-005 (Statistical Techniques), you have learnt that we use the t-test for testing the hypothesis about the means of two independent populations/groups. But if we are interested in testing the means of more than two independent populations/groups, then we cannot apply the t-test. This is because if we use the t-test many times, the type I error increases while testing the hypothesis. In such situations, we use the analysis of variance (ANOVA). The analysis of variance is a method of splitting the total variation in data into different components that measure different sources of variations. The analysis of variance (ANOVA) is categorised as follows in accordance with the classification of the data:

- One-way ANOVA,
- Two-way ANOVA, and
- Two-way ANOVA with m observations per cell.

The analysis of variance (ANOVA) is based on the following assumptions:

1. populations under study are normally distributed with approximately equal variances,
2. samples are independently and randomly drawn from the populations,
3. various effects are additive in nature, and
4. error terms (e_{ij}) are independently and identically distributed with mean zero and variance σ_e^2 .

In this lab session, you will learn how to apply the one-way ANOVA, two-way ANOVA and two-way ANOVA with m observations per cell using MS Excel 2007. In the next lab session, you will study the design of experiments.

Prerequisite

- Lab Session 10 of MSTL-001 (Basic Statistics Lab).
- Block 2 of MST-005 (Statistical Techniques).

Objectives

After performing the activities of this session, you should be able to:

- prepare the spreadsheet in MS Excel 2007;
- apply the one-way ANOVA;
- apply the two-way ANOVA; and
- apply the two-way ANOVA with m observations per cell.

11.2 PROBLEM DESCRIPTION

In this lab session, we consider three problems to illustrate the applications of different types of ANOVA:

1. A researcher wishes to find out whether the average waiting time for a patient to meet a doctor in the emergency room at four hospitals (A, B, C and D) is equal. To study it, he/she takes a sample of patients in the emergency rooms of each hospital on a particular day and records the waiting time. It is the time measured from the instant the patient arrives in the emergency room until he/she is attended by a doctor. The data are given in Table 1.

Table 1: Waiting time

Waiting Time (in minutes)			
Hospital A	Hospital B	Hospital C	Hospital D
10	12	15	15
12	10	20	17
9	20	15	15
12	15	12	18
10	10	18	15
12	8	20	17
8	10	15	16
10	18	15	14
12	15	10	15
10	10	15	17
15	15	14	18
12	20	13	10
10	15	15	15
12	15	20	12
18	18	24	15
	10	15	18
	14	12	
	12	12	
	20		
	15		

Assuming that the waiting time is normally distributed in each hospital and the variances of all waiting time distributions are approximately equal:

- Formulate the null and alternative hypotheses.
 - Is there enough evidence that the average waiting time for a patient to meet a doctor in the emergency room in the hospitals is equal at 5% level of significance?
2. The yields of five varieties of wheat using four different types of fertilisers are recorded in the following table:

Varieties of Wheat	Types of Fertiliser			
	F ₁	F ₂	F ₃	F ₄
V ₁	52	55	45	54
V ₂	54	55	48	56
V ₃	50	53	52	60
V ₄	56	55	50	52
V ₅	58	52	50	58

Use a two-way ANOVA at 2% level of significance to examine:

- whether the effect of fertilisers on the yield is significantly different by assuming that the effect of each fertiliser on the yield of wheat is normally distributed with approximately equal variances.
 - whether there is a significant difference between the effect of five varieties of wheat on the yield by assuming that the effect of each variety of wheat on the yield of wheat is normally distributed with approximately equal variances.
3. A company organised a training programme for three categories of officers: sales managers, zonal managers and regional managers. The company also considered the education level of the employees. Based on their qualifications, the officers were divided into three categories: graduates, post graduates and doctorates. The company wishes to ascertain the effectiveness of the training programme on the employees across designation and qualification levels. The scores obtained from randomly selected employees across different categories are given in Table 3.

Table 3: Scores (out of 100) of the employees

Qualification	Designation		
	Sales Manager	Zonal Manager	Regional Manager
Graduation	62	70	76
	80	80	78
	84	81	80
	65	90	84
Post Graduation	86	90	75
	90	88	92
	80	95	82
	76	86	74
Doctorate	75	90	64
	84	90	72
	79	80	64
	75	74	70

It is known that the effect of each designation level on the effectiveness of the training is normally distributed with approximately equal variances.

Similarly, the effect of each qualification level on the effectiveness of the training is normally distributed with approximately equal variances. Apply the two-way ANOVA and test the significance for the qualification levels, designations and their interaction on the effectiveness of the training at 5% level of significance.

11.3 ONE-WAY ANALYSIS OF VARIANCE

In Unit 6 of MST-005, you have learnt about the one-way analysis of variance (ANOVA). It is used for testing the equality of means of more than two populations when the observations in an experiment are classified with respect to single criterion. That is, we are interested in studying the effect of various levels of one factor/treatment on the dependent variable. For example, we can apply ANOVA for testing the hypothesis that five varieties of wheat (factor) produce equal yield (dependent variable) on an average or for testing the stress level of employees in three different organisations, and so on. In one-way ANOVA, we divide the total source of variation into two components: the source of variation between groups (due to treatment or assignable causes) and within groups (due to error or random causes).

The procedure of the one-way ANOVA has been described in Unit 6 of MST-005. We briefly mention the main steps and formulae as follows:

Step 1: We first formulate the null hypothesis (H_0) and alternative hypothesis (H_1). If the factor/treatment has k levels and μ_i ($i = 1, 2, \dots, k$) is the average effect of the i^{th} level of the factor/treatment, we can formulate the null and alternative hypotheses as follows:

H_0 : The average effect of different levels of the factor/treatment is equal or the average effect of different levels of the factor/treatment is not significantly different

H_1 : The average effect of all levels of the factor/treatment is not equal or the average effect of all levels of the factor/treatment is significantly different

Symbolically,

$$H_0: \mu_1 = \mu_2 = \dots = \mu_k$$

$$H_1: \text{At least one } \mu_i \neq \mu_j \ (i \neq j = 1, 2, \dots, k)$$

Step 2: We calculate the correction factor (CF) and the raw sum of squares (RSS) using the formulae given below:

$$CF = \frac{G^2}{N} \quad \dots (1)$$

$$RSS = \sum_{i=1}^k \sum_{j=1}^{n_i} y_{ij}^2 \quad \dots (2)$$

where G – the grand total or sum of all observations,

N – the total number of observations,

n_i – the size of the i^{th} sample or i^{th} level of the factor,

y_{ij} – j^{th} observation of the i^{th} sample or i^{th} level of the factor.

Step 3: We calculate the total sum of squares (TSS), the sum of squares between samples or sum of squares due to the factor/treatment (SST), the sum of squares within samples or sum of squares due to error (SSE) using the formulae given below:

$$TSS = RSS - CF \quad \dots (3)$$

$$SST = \sum_{i=1}^k T_i^2 - CF \quad \dots (4)$$

$$SSE = TSS - SST \quad \dots (5)$$

where $T_{i\cdot} = \sum_{j=1}^{n_i} T_{ij}$ – the sum of the observations of i^{th} sample or i^{th} level.

Step 4: We find the degrees of freedom (df) as

The df for treatment/factor (v_1) = $k - 1$

The df for error (v_2) = $N - k$

Step 5: We obtain the various mean sums of squares as follows:

$$\text{Mean sum of squares due to treatment/factor (MSST)} = \frac{SST}{k-1} \dots (6)$$

$$\text{Mean sum of squares due to error (MSSE)} = \frac{SST}{N-k} \dots (7)$$

Step 6: We calculate the value of the test statistic F using the formula given below:

$$F = \frac{\text{MSST}}{\text{MSSE}} \dots (8)$$

Step 7: We take a decision about the null hypothesis as follows:

i) Using the critical region approach

We compare the calculated value of the test statistic (F_{cal}) with the critical value ($F_{(v_1, v_2), \alpha}$) at $\alpha\%$ level of significance. If $F_{\text{cal}} \geq F_{(v_1, v_2), \alpha}$, we reject H_0 and if $F_{\text{cal}} < F_{(v_1, v_2), \alpha}$, we do not reject H_0 .

ii) Using the p-value value approach

We calculate the p-value using the formula given below :

$$\text{p-value} = P[F \geq F_{\text{cal}}] \dots (9)$$

We compare the calculated p-value with the given level of significance (α). If the p-value is less than or equal to α , we reject the null hypothesis and if it is greater than α , we do not reject the null hypothesis.

Step 8: If the null hypothesis is rejected, i.e., there is a significant difference between the means of different levels of the factor/treatment, we do pair-wise comparison to know which pair of means differs significantly. We use the t-test for testing the equality of two treatment (population) means as follows:

i) We formulate the null and alternative hypothesis as given below :

$$H_0 : \mu_i = \mu_j$$

$$H_1 : \mu_i \neq \mu_j \quad (i \neq j; i = 1, 2, \dots, k \text{ and } j = 1, 2, \dots, k)$$

ii) We calculate the mean of each level of the factor/treatment, say,

$$\bar{T}_1, \bar{T}_2, \dots, \bar{T}_k .$$

iii) We take all possible combinations of two levels of the factor/treatment and then calculate the absolute difference of the means of each pair of the levels of the factor/treatment as $|\bar{T}_i - \bar{T}_j|$, where $i \neq j$ ($i = 1, 2, \dots, k$ and $j = 1, 2, \dots, k$).

- iv) We compute the critical difference (CD) of the means for each pair of the levels of the factor/treatment using the formula given below:

$$CD = t_{(\text{error df}), \alpha/2} \sqrt{\frac{1}{MSSE} \left(\frac{1}{n_i} + \frac{1}{n_j} \right)} \quad \dots (10)$$

where $t_{(\text{error df}), \alpha/2}$ is the critical (tabulated) value of the t-distribution at error degrees of freedom.

- v) We compare $|\bar{T}_i - \bar{T}_j|$ with the corresponding CD and if $|\bar{T}_i - \bar{T}_j|$ is greater than or equal to CD, we reject the null hypothesis
 $H_0 : \mu_i = \mu_j$, i.e., the difference between μ_i and μ_j is significant.
Otherwise, we do no reject the null hypothesis.

Steps in Excel

In Problem 1, we have to test whether the average waiting time for a patient to meet a doctor in the emergency room at four hospitals (A, B, C and D) is equal. Here we have to test the equality of more than two means. It is given that the waiting time is normally distributed in each hospital and the variances of all waiting time distributions are approximately equal. So we may use the one-way ANOVA.

We first have to set up the null and alternative hypotheses. If μ_1, μ_2, μ_3 and μ_4 denote the average waiting time for a patient to meet a doctor in the emergency room at Hospitals A, B, C and D, respectively, we can formulate the null and alternative hypotheses as follows:

$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4$ (claim)

H_1 : All means are not equal

The **Data Analysis ToolPak** in MS Excel 2007 gives direct results for the one-way ANOVA. So we explain the procedure for applying the one-way ANOVA using the **Data Analysis ToolPak** in MS Excel 2007 in the following steps:

Step 1: We enter the data (given in Table 1) in an Excel 2007 spreadsheet. We start by entering the heading of the data in Rows 1 and 2 in Excel sheet and the data from Row 3 onwards as shown in Fig. 11.1.

	A	B	C	D	E
1	Waiting time (in minutes)				
2	S.No.	Hospital A	Hospital B	Hospital C	Hospital D
3	1	10	12	15	15
4	2	12	10	20	17
5	3	9	20	15	15
6	4	12	15	12	18
7	5	10	10	18	15
8	6	12	8	20	17
9	7	8	10	15	16
10	8	10	18	15	14
11	9	12	15	10	15
12	10	10	10	15	17

Fig. 11.1: Partial screenshot of the spreadsheet for the given data.

Step 2: We click on **Data** tab → **Data Analysis** → **Anova: Single Factor** → **OK** as shown in Figs. 11.2a and b. A new dialog box opens (Fig. 11.2c).

Analysis of Variance

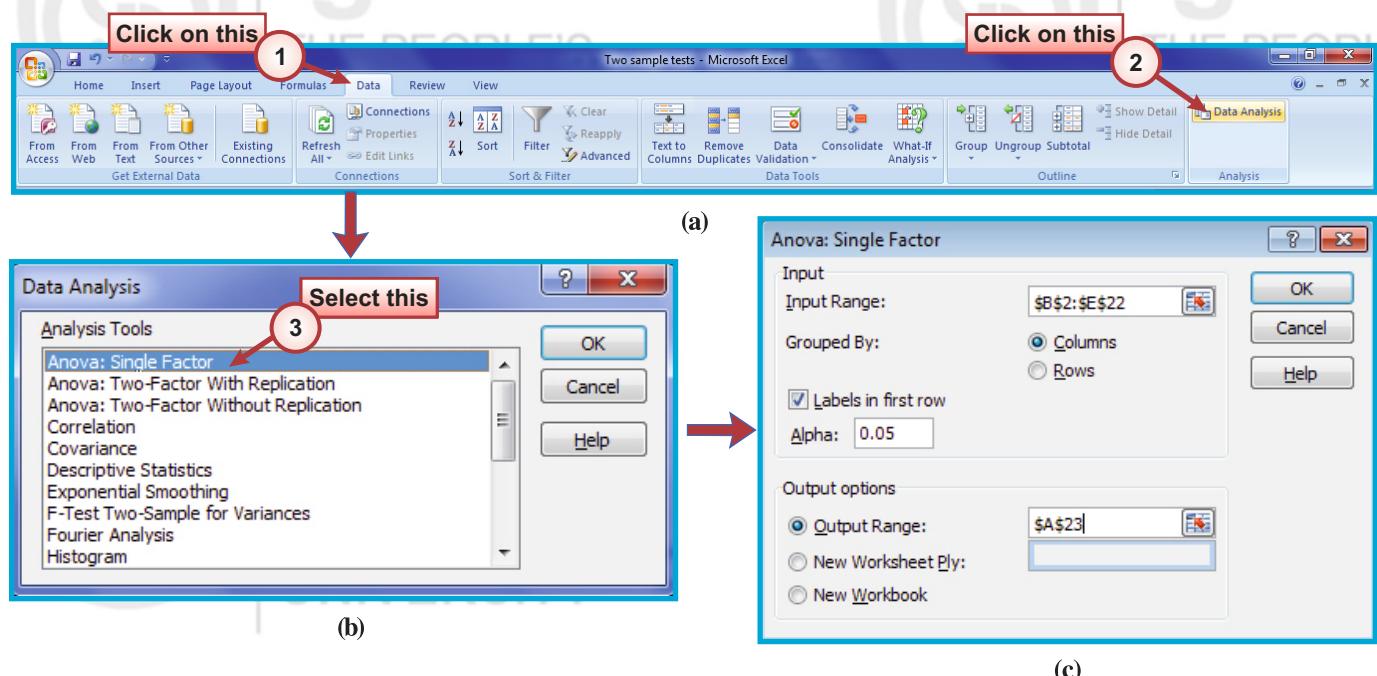


Fig. 11.2

Step 3: In the new dialog box (Fig. 11.2c), we

1. specify the data with labels in **Input Range**, i.e., Cells B2:E22,
2. click on **columns** since we have grouped the data column-wise, i.e., here we use Columns B, C, D and E for Hospitals A, B, C and D, respectively,
3. click on **Labels in first row** box since we have included data labels, i.e., Hospital A in Cell B2, Hospital B in Cell C2, Hospital C in Cell D2 and Hospital D in Cell E2,
4. type the value of the level of significance (α) in **Alpha**, i.e., 0.05,
5. click on **Output Range** and then select the cell in which we wish to put the results. Here we select Cell A23.

Finally, we click on **OK**. When we do so, we obtain the results shown in Fig. 11.3.

	A	B	C	D	E	F	G
22	20		15				
23	Anova: Single Factor						
24							
25	SUMMARY						
26	Groups	Count	Sum	Average	Variance		
27	Hospital A	15	172	11.4667	6.1238		
28	Hospital B	20	282	14.1	14.2		
29	Hospital C	18	280	15.5556	12.732		
30	Hospital D	16	247	15.4375	4.7958		
31							
32							
33	ANOVA						
34	Source of Variation	SS	df	MS	F	P-value	F crit
35	Between Groups	169.8239	3	56.6080	5.7143	0.0016	2.7459
36	Within Groups	643.9153	65	9.9064			
37							
38	Total	813.7391	68				

Fig. 11.3

Note from Fig. 11.3 that the **Data analysis ToolPak** of MS Excel gives results in two parts: SUMMARY and ANOVA.

SUMMARY has the following headings:

1. *Groups* – Under this heading, it gives the groups (treatments or levels of a treatment/factor) under study. In this problem, the hospitals are being studied. So it shows Hospitals A, B, C and D in Cells A27, A28, A29 and A30, respectively.
2. *Counts* – Under this heading, it counts the number of observations in each group. The number of observations of Hospitals A, B, C and D are shown in Cells B27, B28, B29 and B30, respectively.
3. *Sum* – Under this heading, it gives the sum of observations in each group. The sum of waiting times in Hospitals A, B, C and D are shown in Cells C27, C28, C29 and C30, respectively.
4. *Average* – Under this heading, it gives the average of observations in each group. The average waiting times in Hospitals A, B, C and D are shown in Cells D27, D28, D29 and D30, respectively.
5. *Variance* – Under this heading, it gives the variance of observations in each group. The variances of waiting times in Hospitals A, B, C and D are shown in Cells E27, E28, E29 and E30, respectively.

ANOVA has the following headings:

1. *Source of Variation* – Under this heading, it gives the source of variation in the given data. In one-way ANOVA, we divide the total source of variation into two components: source of variation between groups (due to treatment/factor or assignable causes) and within groups (due to error or random causes). Therefore, it shows source of variation between groups in Cell A35 and within groups in Cell A36.
2. *SS* – Under this heading, it gives the sum of squares of the observations between groups and within groups. Here it calculates the sum of squares of waiting times between hospitals and within hospitals in Cells B35 and B36, respectively, and the total of both in Cell B38.
3. *df* – Under this heading, it gives the degrees of freedom between groups and within groups. Here it shows df for source of variation between hospitals and within hospitals in Cell C35 and C36, respectively. The total df shows in Cell C38.
4. *MS* – Under this heading, it gives the mean sum of squares of the observations between groups and within groups. Here it calculates the mean sum of squares of waiting times between hospitals and within hospitals in Cells D35 and D36, respectively.
5. *F* – Under this heading, it calculates the value of the test statistic F. Here it gives its value in Cell E35.
6. *p-value* – It gives the p-value for the test in Cell F35.
7. *F crit* – Under this heading, it gives the critical (tabulated) value for the test at the given level of significance. Here it gives the value of $F_{(3,65), 0.05}$ in Cell G35 at 5% level of significance.

Decision using the critical region approach

Note that the calculated F value is 5.7143, which is greater than the critical value 2.7459. It means that F_{cal} lies in the rejection region. So we reject the null hypothesis. Since the null hypothesis is the claim, we reject the claim. Hence, we conclude that the samples provide us sufficient

evidence against the claim. So we may conclude that the average waiting time for a patient to meet a doctor in the emergency room at four hospitals is not equal at 5% level of significance.

Decision using the p-value approach

Since the p-value ($= 0.0016$) is less than α ($= 0.05$), we reject the null hypothesis.

Pair-wise comparison

Step 4: Since the null hypothesis is rejected, i.e., there is a significant difference between the average waiting times, we do the pair-wise comparison as follows:

1. Since the average waiting time in each hospital has already been calculated using the **Data Analysis ToolPak** in Cells D27:D30, we prepare the comparison table as shown in Fig. 11.4.

	A	B	C	D	E	F	G	H	I
40	Comparison table								
41	Pair of hospitals	n_i	n_j	\bar{T}_i	\bar{T}_j	$ \bar{T}_i - \bar{T}_j $	$t_{(65), 0.025}$	CD	Inference
42	A, B								
43	A, C								
44	A, D								
45	B, C								
46	B, D								
47	C, D								

Fig. 11.4

2. We type the number of observations (n_i and n_j) under the first hospital and the second hospital of each pair of hospitals in Cells B42:B47 and C42:C47, respectively, as shown in Fig. 11.5.

	A	B	C	D
40	Comparison table			
41	Pair of hospitals	n_i	n_j	\bar{T}_i
42	A, B	15	20	
43	A, C	15	18	
44	A, D	15	16	
45	B, C	20	18	
46	B, D	20	16	
47	C, D	18	16	

Fig. 11.5

3. We copy the average waiting time of the first hospital (\bar{T}_i) and the second hospital (\bar{T}_j) for each pair of hospitals in Cells D42:D47 and Cells E42:E47 from Cells D27:D30 (Fig. 11.6).

	A	B	C	D	E
40	Comparison table				
41	Pair of hospitals	n_i	n_j	\bar{T}_i	\bar{T}_j
42	A, B	15	20	11.467	14.1
43	A, C	15	18	11.467	15.556
44	A, D	15	16	11.467	15.438
45	B, C	20	18	14.1	15.556
46	B, D	20	16	14.1	15.438
47	C, D	18	16	15.556	15.438

Fig. 11.6

4. We calculate the absolute difference $|\bar{T}_i - \bar{T}_j|$ of each pair of hospitals by typing “=Abs(D42-E42)” in Cell F42 and then dragging down the cell up to Cell F47 as shown in Fig. 11.7.

	B	C	D	E	F	G
40						
41	n_i	n_j	\bar{T}_i	\bar{T}_j	$ \bar{T}_i - \bar{T}_j $	$t_{(65),0.025}$
42	15	20	11.467	14.1	2.6333	
43	15	18	11.467	15.556	4.0889	
44	15	16	11.467	15.438	3.9708	
45	20	18	14.1	15.556	1.4556	
46	20	16	14.1	15.438	1.3375	
47	18	16	15.556	15.438	0.1181	

Fig. 11.7

5. We now obtain the critical value of the t-test in Cell G42 as explained in Step 4 under the heading ‘Steps in Excel’ of Sec. 9.4 of Lab Session 9. Then we drag down the cell up to Cell G47. The output is shown in Fig. 11.8.

	B	C	D	E	F	G
40						
41	n_i	n_j	\bar{T}_i	\bar{T}_j	$ \bar{T}_i - \bar{T}_j $	$t_{(65),0.025}$
42	15	20	11.467	14.1	2.6333	1.9971
43	15	18	11.467	15.556	4.0889	1.9971
44	15	16	11.467	15.438	3.9708	1.9971
45	20	18	14.1	15.556	1.4556	1.9971
46	20	16	14.1	15.438	1.3375	1.9971
47	18	16	15.556	15.438	0.1181	1.9971

Fig. 11.8

6. We calculate the critical difference (CD) for each pair. For the first pair (A, B), we type “=G42*SQRT(\$D\$36*(1/B42+1/C42))” in Cell H42 and press **Enter**. Then we drag down the cell up to Cell H47. The output is shown in Fig. 11.9.

	C	D	E	F	G	H
40						
41	n_j	\bar{T}_i	\bar{T}_j	$ \bar{T}_i - \bar{T}_j $	$t_{(65),0.025}$	CD
42	20	11.467	14.1	2.6333	1.9971	2.1470
43	18	11.467	15.556	4.0889	1.9971	2.1976
44	16	11.467	15.438	3.9708	1.9971	2.2591
45	18	14.1	15.556	1.4556	1.9971	2.0422
46	16	14.1	15.438	1.3375	1.9971	2.1083
47	16	15.556	15.438	0.1181	1.9971	2.1598

Fig. 11.9

7. We compare $|\bar{T}_i - \bar{T}_j|$ with the corresponding CD and if $|\bar{T}_i - \bar{T}_j|$ is greater than or equal to CD, we reject the null hypothesis $H_0 : \mu_i = \mu_j$, i.e., the difference between μ_i and μ_j is significant. So we type “Significant” in the corresponding cell under the **Inference** column. Otherwise, we do no reject the null hypothesis and type “Insignificant” in the corresponding cell. The results are shown in Fig. 11.10.

	A	B	C	D	E	F	G	H	I
40	Comparison table								
41	Pair of hospitals	n_i	n_j	\bar{T}_i	\bar{T}_j	$ \bar{T}_i - \bar{T}_j $	$t_{(65),0.025}$	CD	Inference
42	A, B	15	20	11.467	14.1	2.6333	1.9971	2.1470	Significant
43	A, C	15	18	11.467	15.556	4.0889	1.9971	2.1976	Significant
44	A, D	15	16	11.467	15.438	3.9708	1.9971	2.2591	Significant
45	B, C	20	18	14.1	15.556	1.4556	1.9971	2.0422	Insignificant
46	B, D	20	16	14.1	15.438	1.3375	1.9971	2.1083	Insignificant
47	C, D	18	16	15.556	15.438	0.1181	1.9971	2.1598	Insignificant

Fig. 11.10

8. From the comparison table shown in Fig. 11.10, we observe that the average waiting time for a patient to meet a doctor in the emergency room in Hospital A is significantly different from that in Hospitals B, C and D. However, the average waiting time in Hospital B is not significantly different from that in Hospitals C and D. Also the average waiting time in Hospital C is not significantly different from that in Hospital D.

11.4 TWO-WAY ANALYSIS OF VARIANCE

In Unit 7 of MST-005, you have learnt that if we are interested in studying the simultaneous effect of two independent factors on the dependent variable, we use two-way ANOVA with one observation per cell or simply two-way ANOVA. For example, we may wish to study the simultaneous effects of five varieties of wheat (first criterion) and four different types of fertilisers (second criterion) on the yield (dependent variable) or test the stress level of employees in three different organisations in different regions, and so on. In such situations, we can also apply two separate one-way ANOVAs for each treatment/factor. However, it is more advantageous to use two-way ANOVA because the variance can be reduced by introducing the second factor.

In two-way ANOVA, the total variation in the data is divided into three components: variation due to the first criterion (factor), variation due to the second criterion (factor) and variation due to error.

The testing procedure of two-way ANOVA has been described in Unit 7 of MST-005. We briefly mention the main steps and formulae as follows:

Step 1: We first formulate the null hypothesis (H_0) and alternative hypothesis (H_1). In two-way ANOVA, we can test two hypotheses simultaneously: one for different levels of factor A and the other for different levels of factor B. If factor A has p levels, we can set up the null and alternative hypotheses as follows:

$$H_{0A} : \alpha_1 = \alpha_2 = \dots = \alpha_p$$

$$H_{1A} : \text{At least one } \alpha_i \neq \alpha_j \ (i \neq j = 1, 2, \dots, p)$$

Similarly, if factor B has q levels, we can set up the null and alternative hypotheses as follows:

$$H_{0B} : \beta_1 = \beta_2 = \dots = \beta_q$$

$$H_{1B} : \text{At least one } \beta_i \neq \beta_j \ (i \neq j = 1, 2, \dots, q)$$

Step 2: We calculate the correction factor (CF) and the raw sum of squares (RSS) using the formulae given below:

$$CF = \frac{G^2}{N} \quad \dots (11)$$

$$RSS = \sum_{i=1}^p \sum_{j=1}^q y_{ij}^2 \quad \dots (12)$$

where G – the grand total, i.e., sum of all observations,

N – the total number of observations, i.e., $N = pq$, and

y_{ij} – observation of the i^{th} level of factor A and the j^{th} level of factor B.

Step 3: We calculate the total sum of squares (TSS), the sum of squares between rows or sum of squares due to factor A (SSA), the sum of squares between columns or sum of squares due to factor B (SSB) and sum of squares due to error (SSE) as follows:

$$TSS = RSS - CF \quad \dots (13)$$

$$SSA = \frac{1}{q} \sum_{i=1}^p y_i^2 - CF \quad \dots (14)$$

$$SSB = \frac{1}{p} \sum_{j=1}^q y_j^2 - CF \quad \dots (15)$$

$$SSE = TSS - SSA - SSB \quad \dots (16)$$

where y_i – the sum of the observations of the i^{th} level of factor A.

y_j – the sum of the observations of the j^{th} level of factor B.

Step 4: We find the degrees of freedom (df) as

The df for factor A = $p - 1$

The df for factor B = $q - 1$

The df for error = $(p - 1)(q - 1)$

Step 5: We obtain the various mean sums of squares as follows:

$$\text{Mean sum of squares due to factor A (MSSA)} = \frac{SSA}{p-1} \quad \dots (17)$$

$$\text{Mean sum of squares due to factor B (MSSB)} = \frac{SSB}{q-1} \quad \dots (18)$$

$$\text{Mean sum of squares due to error (MSSE)} = \frac{SSE}{(p-1)(q-1)} \quad \dots (19)$$

Step 6: We calculate the value of the test statistics using the formulae given below:

$$F_A = \frac{MSSA}{MSSE} \quad \dots (20)$$

$$F_B = \frac{MSSB}{MSSE} \quad \dots (21)$$

Step 7: We take decisions about the null hypotheses for factor A and factor B as explained in Step 7 of the screened box in Sec 11.3.

Steps in Excel

In Problem 2, we have to test whether the effect of fertilisers on the yield is significantly different and whether there is a significant difference between the effect of five varieties of wheat on the yield. Here we have to test whether there is significant difference between the varieties of wheat (different levels of factor A) and different fertilisers (different levels of factor B). Since two factors are being studied, we use the two-way ANOVA.

We first have to set up the null and alternative hypotheses. Let α_i ($i = 1, 2, \dots, 5$) and β_j ($j = 1, 2, 3, 4$) denote the average effect of the i^{th} variety of wheat and the j^{th} type of fertiliser on the yield, respectively. We formulate the null and alternative hypotheses for varieties of wheat and types of fertilisers as follows:

$$H_{0V} : \alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = \alpha_5$$

H_{1V} : Average effect of all varieties of wheat on the yield are not equal (claim)

$$H_{0F} : \beta_1 = \beta_2 = \beta_3 = \beta_4$$

H_{1F} : Average effect of all types of fertiliser on the yield are not equal (claim)

We use the **Data Analysis ToolPak** in MS Excel 2007 to apply the two-way ANOVA as follows:

Step 1: We enter the data (given in Table 2) in an Excel 2007 spreadsheet as shown in Fig. 11.11.

	A	B	C	D	E
1	Varieties of wheat	Types of fertiliser			
2		F ₁	F ₂	F ₃	F ₄
3	V ₁	52	55	45	54
4	V ₂	54	55	48	56
5	V ₃	50	53	52	60
6	V ₄	56	55	50	52
7	V ₅	58	52	50	58
8					

Fig. 11.11: Partial screenshot of the spreadsheet for the given data.

Step 2: We click on **Data** tab → **Data Analysis** → **Anova: Two-Factor Without Replication** → **OK** as shown in Figs. 11.12a and b.
As a result, we get the new dialog box shown in Fig. 11.12c.

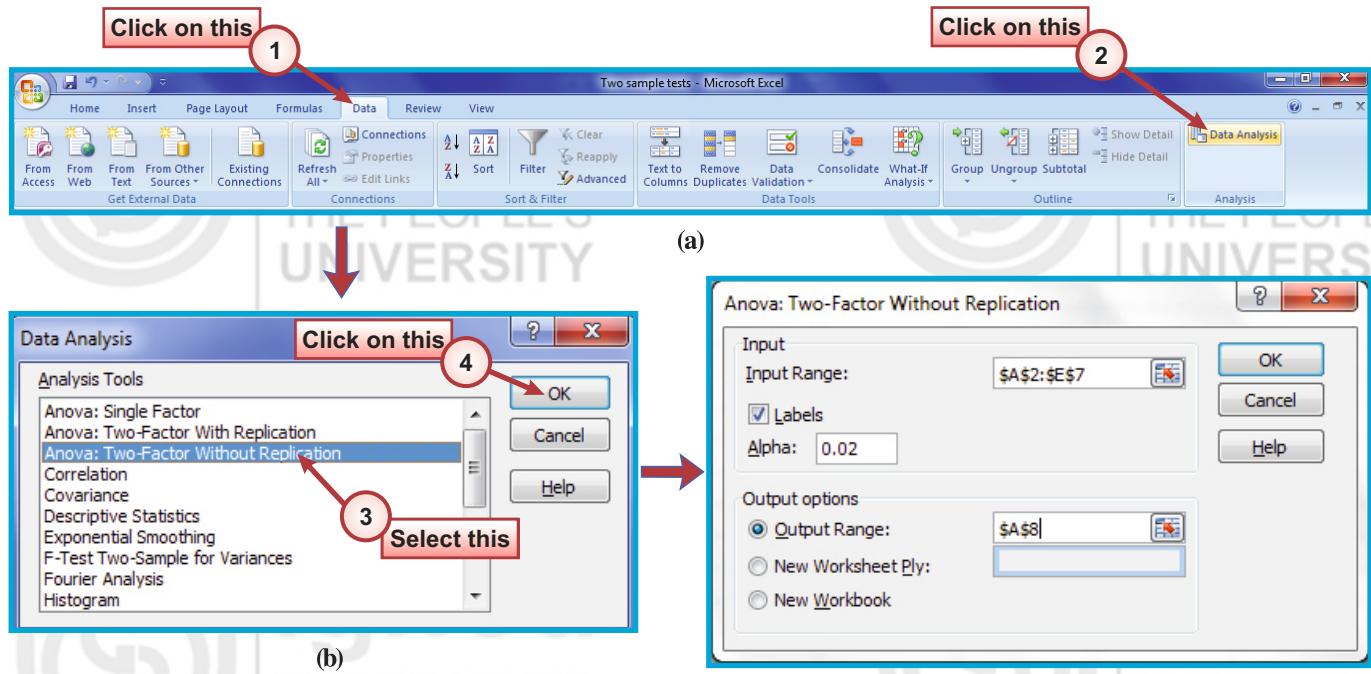


Fig. 11.12

Step 3: In the new dialog box, we

- specify data with label in **Input Range**, i.e. Cells A2:E7,
- click on the **Labels** box since we have included data labels, i.e., varieties of wheat in Cells A3, A4, A5, A6, A7 and types of fertiliser in Cells B2, C2, D2, E2,
- type the value of the level of significance (α) in **Alpha**, i.e., 0.02,
- click on the **Output Range** and select the cell in which we wish to put the results. Here we select Cell A8. Finally, we click on the **OK**. When we do so, we obtain the results shown in Fig. 11.13.

	A	B	C	D	E	F	G
8	Anova: Two-Factor Without Replication						
9							
10	SUMMARY	<i>Count</i>	<i>Sum</i>	<i>Average</i>	<i>Variance</i>		
11	V1	4	206	51.5	20.333		
12	V2	4	213	53.25	12.917		
13	V3	4	215	53.75	18.917		
14	V4	4	213	53.25	7.5833		
15	V5	4	218	54.5	17		
16							
17	F1	5	270	54	10		
18	F2	5	270	54	2		
19	F3	5	245	49	7		
20	F4	5	280	56	10		
21							
22							
23	ANOVA						
24	<i>Source of Variation</i>	<i>SS</i>	<i>df</i>	<i>MS</i>	<i>F</i>	<i>P-value</i>	<i>F crit</i>
25	Rows	19.5	4	4.8750	0.6062	0.6658	4.4187
26	Columns	133.75	3	44.5833	5.5440	0.0127	4.8145
27	Error	96.5	12	8.0417			
28	Total	249.75	19				
29							

Fig. 11.13

The explanation of these results is the same as that for one-way ANOVA. Let us now test the hypotheses.

Decision using the critical region approach

Note that the calculated F value for varieties of wheat (rows) is 0.6062, which is less than the corresponding critical value 4.4187. It means that F_{cal} lies in the non-rejection region. So we do not reject the null hypothesis. Since the alternative hypothesis is the claim, we reject the claim. Hence, we conclude that there is not enough evidence to claim that there is a significant difference between the five varieties of wheat on the yield at 2% level of significance.

The calculated F value for fertilisers (columns) is 5.5440, which is greater than the corresponding critical value 4.8145. It means that it lies in the rejection region. So we reject the null hypothesis. Since the alternative hypothesis is the claim, we do not reject the claim. Hence, we conclude that there is enough evidence to claim that there is a significant difference between the four fertilisers on the yield at 2% level of significance.

Decision using the p-value approach

Since the p-value for varieties of wheat ($= 0.6658$) is greater than $\alpha (= 0.02)$, we do not reject the null hypothesis.

Since the p-value for fertilisers ($= 0.0127$) is less than $\alpha (= 0.02)$, we reject the null hypothesis.

Step 4: Since the null hypothesis for fertilisers is rejected, i.e., the average effect of different fertilisers on the yield is significantly different, we do pair-wise comparison as explained in Step 4 under the heading ‘Steps in Excel’ of Sec. 11.3. The final output is shown in Fig. 11.14.

	A	B	C	D	E	F	G	H	I
32	Comparison table								
33	Pair of fertilisers	n_i	n_j	\bar{T}_i	\bar{T}_j	$ \bar{T}_i - \bar{T}_j $	$t_{(12), 0.01}$	CD	Inference
34	F ₁ , F ₂	5	5	54	54	0	2.681	4.8084	Insignificant
35	F ₁ , F ₃	5	5	54	49	5	2.681	4.8084	Significant
36	F ₁ , F ₄	5	5	54	56	2	2.681	4.8084	Insignificant
37	F ₂ , F ₃	5	5	54	49	5	2.681	4.8084	Significant
38	F ₂ , F ₄	5	5	54	56	2	2.681	4.8084	Insignificant
39	F ₃ , F ₄	5	5	49	56	7	2.681	4.8084	Significant

Fig. 11.14

From the comparison table shown in Fig. 11.14, we observe that the effect of fertilisers F_1 and F_2 on the yield is significantly different from that of fertiliser F_3 . However, the effect of fertiliser F_1 is not significantly different from that of fertilisers F_2 and F_4 . Also, the effect of fertiliser F_2 is not significantly different from that of fertiliser F_4 . However, the effect of fertiliser F_3 is significantly different from that of fertiliser F_4 .

11.5 TWO-WAY ANALYSIS OF VARIANCE WITH m OBSERVATIONS PER CELL

So far we have used MS Excel to apply one-way ANOVA and two-way ANOVA techniques. In the two-way ANOVA, we had one observation per cell. But if we are interested in studying the effect of two factors as well as the effect of interaction between the two factors on the dependent variable, then we use two-way ANOVA with more than one observation per cell. However, the number of observations in each cell should be equal. Otherwise, the analysis becomes very complicated. In such situations, we can also apply two separate one-way ANOVAs. However, it is more advantageous to use two-way ANOVA because in addition to knowing the effect of two factors separately, we can also study the effect of their interaction on the dependent variable. Further, in two way ANOVA, the error variance can be reduced by introducing the second factor.

The procedure for the two-way ANOVA with m observations per cell has been described in Unit 8 of MST-005. We briefly mention the main steps and formulae as follows:

Step 1: We first formulate the null hypothesis (H_0) and alternative hypothesis (H_1). In two-way ANOVA with m observations per cell, we can test three hypotheses simultaneously: one for different levels of factor A, the second for different levels of factor B and the third for their interaction. If factor A has p levels, we can take the null and alternative hypotheses as follows:

$$H_{0A} : \alpha_1 = \alpha_2 = \dots = \alpha_p$$

H_{1A} : Means of all levels of factor A are not equal

Similarly, if factor B has q levels, we can take the null and alternative hypotheses as follows:

$$H_{0B} : \beta_1 = \beta_2 = \dots = \beta_q$$

H_{1B} : Means of all levels of factor B are not equal

We can take the null and alternative hypotheses for interaction as follows:

$$H_{0AB} : \alpha_i\beta_j = 0 \quad (i \neq j; i=1, 2, \dots, p \text{ and } j=1, 2, \dots, q)$$

H_{1AB} : All interaction effects are not zero

Let there be m observations in each cell.

Step 2: We calculate the correction factor (CF) and the raw sum of squares (RSS) using the formulae given below:

$$CF = \frac{G^2}{N} \quad \dots (22)$$

$$RSS = \sum_{i=1}^p \sum_{j=1}^q \sum_{k=1}^m y_{ijk}^2 \quad \dots (23)$$

where G – the grand total, i.e., the sum of all observations,
 N – the total number of observations, i.e., $N = pqm$,
 y_{ijk} – the k^{th} observation of the i^{th} level of factor A and the j^{th} level of factor B.

Step 3: We calculate the various sums of squares using the formulae given below:

$$\text{TSS} = \text{RSS} - \text{CF} \quad \dots (24)$$

$$\text{SSA} = \frac{1}{qm} \sum_{i=1}^p y_{i..}^2 - \text{CF} \quad \dots (25)$$

$$\text{SSB} = \frac{1}{pm} \sum_{j=1}^q y_{.j}^2 - \text{CF} \quad \dots (26)$$

$$\text{SSAB} = \frac{1}{m} \sum_{i=1}^p \sum_{j=1}^q y_{ij.}^2 - \text{CF} - \text{SSA} - \text{SSB} \quad \dots (27)$$

$$\text{SSE} = \text{TSS} - \text{SSA} - \text{SSB} - \text{SSAB} \quad \dots (28)$$

where $y_{i..}$ – the sum of observations of the i^{th} level of factor A.

$y_{.j}$ – the sum of observations of the j^{th} level of factor B.

$y_{ij.}$ – the sum of observations of the i^{th} level of factor A and the j^{th} level of factor B.

Step 4: We obtain the various mean sums of squares as follows:

$$\text{Mean sum of squares due to factor A (MSSA)} = \frac{\text{SSA}}{p-1} \quad \dots (29)$$

$$\text{Mean sum of squares due to factor B (MSSB)} = \frac{\text{SSB}}{q-1} \quad \dots (30)$$

$$\text{Mean sum of squares due to interaction (MSSAB)} = \frac{\text{SSAB}}{(p-1)(q-1)} \quad \dots (31)$$

$$\text{Mean sum of squares due to error (MSSE)} = \frac{\text{SSE}}{pq(m-1)} \quad \dots (32)$$

Step 5: We calculate the value of the test statistics using the formulae given below:

$$F_A = \frac{\text{MSSA}}{\text{MSSE}} \quad \dots (33)$$

$$F_B = \frac{\text{MSSB}}{\text{MSSE}} \quad \dots (34)$$

$$F_{AB} = \frac{\text{MSSAB}}{\text{MSSE}} \quad \dots (35)$$

Step 6: We take decisions about the null hypothesis for factor A, factor B and their interaction AB as explained in Step 7 of the screened box of Sec. 11.3.

Steps in Excel

In Problem 3, we have to test the significance of two factors (designation and qualification) and their interaction simultaneously. So we use the two-way ANOVA with m observations per cell.

We first set up the null and alternative hypotheses. Let α_i ($i = 1, 2, 3$), β_j ($j = 1, 2, 3$) and $\alpha_i\beta_j$ denote the average effect of the i^{th} qualification level, j^{th} designation level and the effect of interaction between the i^{th} qualification level and j^{th} designation level, respectively, on the effectiveness of the training. We formulate the null and alternative hypotheses for qualification levels, designation levels and their interaction as follows:

$$H_{0Q} : \alpha_1 = \alpha_2 = \alpha_3$$

H_{1Q} : Average effect of all qualification levels are not equal (claim)

$$H_{0D} : \beta_1 = \beta_2 = \beta_3$$

H_{1D} : Average effect of all designations are not equal (claim)

$$H_{0QD} : \alpha_i\beta_j = 0 \quad (i=1,2,3 \text{ and } j=1,2,3)$$

H_{1QD} : All interaction effects are not zero (claim)

We apply the two-way ANOVA with m observations per cell using the **Data Analysis ToolPak** in MS Excel 2007 as follows:

Step 1: We enter the data (given in Table 3) in an Excel 2007 spreadsheet as shown in Fig. 11.15.

	A	B	C	D	E
1			Designation		
2			Sales manager	Zonal manager	Regional manager
Qualification	Graduate		62	70	76
			80	80	78
			84	81	80
			65	90	84
	Post graduate		86	90	75
			90	88	92
			80	95	82
	Doctorate		76	86	74
			75	90	64
			84	90	72
			79	80	64
			75	74	70

Fig. 11.15: Partial screenshot of the spreadsheet for the given data.

Step 2: We click on **Data** tab \rightarrow **Data Analysis** \rightarrow **Anova: Two-Factor With Replication** \rightarrow **OK** (Figs. 11.16a and b). A new dialog box opens (Fig. 11.16c).

The figure consists of three parts labeled (a), (b), and (c), illustrating the process of running a two-factor ANOVA with replication in Microsoft Excel 2007.

- (a) Excel ribbon:** Shows the 'Data' tab selected. Red boxes and numbers indicate the steps: 1 points to the 'Data' tab, 2 points to the 'Data Analysis' button, and 3 points to the 'Anova: Two-Factor With Replication' option in the Data Analysis dialog box.
- (b) Data Analysis dialog box:** Shows the 'Analysis Tools' list with 'Anova: Two-Factor With Replication' selected. Red box 4 points to the 'OK' button. A red box labeled 'Select this' points to the selected option.
- (c) Anova: Two-Factor With Replication dialog box:** Shows the input range as '\$B\$2:\$E\$14', rows per sample as '4', alpha as '0.05', and output range as '\$A\$15'. Red box 2 points to the 'OK' button.

Fig. 11.16

Step 3: In the new dialog box, we

1. specify data with labels in **Input Range**, i.e., Cells B2:E14,
2. type the number of rows appearing in each column factor (designation) in **Rows per Sample**. Here it is 4 since there are four rows in each designation,
3. type the value of level of significance (α) in **Alpha**, i.e., 0.05,
4. click on the **Output Range** and select the cell in which we wish to put the output results. Here we select Cell A18.

Finally, we click on **OK** and obtain the results shown in Fig. 11.17.

15 Anova: Two-Factor With Replication					
	SUMMARY	Sales managers	Zonal managers	Regional managers	Total
<i>Graduate</i>					
19 Count	4	4	4	12	
20 Sum	291	321	318	930	
21 Average	72.75	80.25	79.5	77.5	
22 Variance	118.25	66.9167	11.6667	66.0909	
23					
<i>Post graduates</i>					
25 Count	4	4	4	12	
26 Sum	332	359	323	1014	
27 Average	83	89.75	80.75	84.5	
28 Variance	38.6667	14.9167	68.9167	49.3636	
29					
<i>Doctorate</i>					
31 Count	4	4	4	12	
32 Sum	313	334	270	917	
33 Average	78.25	83.5	67.5	76.4167	
34 Variance	18.25	62.3333	17	74.9924	
35					
<i>Total</i>					
37 Count	12	12	12		
38 Sum	936	1014	911		
39 Average	78	84.5	75.9167		
40 Variance	66.9091	56.2727	65.5379		
41					
42					
43 ANOVA					
Source of Variation	SS	df	MS	F	P-value F crit
45 Sample	462.0556	2	231.0278	4.9872	0.0143 3.3541
46 Columns	481.0556	2	240.5278	5.1923	0.0124 3.3541
47 Interaction	363.1111	4	90.7778	1.9596	0.1293 2.7278
48 Within	1250.7500	27	46.3241		
49					
50 Total	2556.972	35			

Fig. 11.17

The explanation of these results is the same as that for one-way ANOVA. We now test the hypotheses.

Decision using the critical region approach

Note that the calculated F value for the qualification level (sample) is 4.9872 and it is greater than the corresponding critical value 3.3541. It means that it lies in the rejection region. So we reject the null hypothesis. Since the alternative hypothesis is the claim, we do not reject the claim. Hence, we conclude that there is not enough evidence against the claim. So we may conclude that there is a significant difference between the qualification levels on the effectiveness of the training at 5% level of significance.

The calculated F value for the designation (columns) is 5.1923, which is greater than the corresponding critical value 3.3541. It means that it

lies in the rejection region. So we reject the null hypothesis. Since the alternative hypothesis is the claim, we do not reject the claim. Hence, we conclude that there is not enough evidence against the claim. So we may conclude that there is a significant difference between the designation levels on the effectiveness of the training at 5% level of significance.

The calculated F value for interaction is 1.9596, which is less than the corresponding critical value 2.7278. It means that it lies in the non-rejection region. So we do not reject the null hypothesis. Since the alternative hypothesis is the claim, we reject the claim. Hence, we conclude that there is enough evidence against the claim. So we may conclude that there is not enough evidence to claim that all interactions are not zero at 5% level of significance.

Decision using the p-value approach

Since the p-value for qualification level ($= 0.0143$) is less than $\alpha (= 0.05)$, we reject the null hypothesis.

Since the p-value for designation ($= 0.0124$) is less than $\alpha (= 0.05)$, we reject the null hypothesis.

Since the p-value for interaction ($= 0.1293$) is greater than $\alpha (= 0.05)$, we do not reject the null hypothesis.

Step 4: Since the null hypothesis for qualification level as well as designation are rejected, we do pair-wise comparison as explained in Step 4 under the heading ‘Steps in Excel’ of Sec. 11.3 for both. The final output for the pair-wise comparison for qualification level is shown in Fig. 11.18.

	A	B	C	D	E	F	G	H	I
52	Comparison table								
53	Pair of qualification levels	n_i	n_j	\bar{T}_i	\bar{T}_j	$ \bar{T}_i - \bar{T}_j $	$t_{(27), 0.025}$	CD	Inference
54	G, PG	12	12	77.5	84.5	7	2.0518	5.70124	Significant
55	G, D	12	12	77.5	76.4167	1.0833	2.0518	5.70124	Insignificant
56	PG, D	12	12	84.5	76.4167	8.0833	2.0518	5.70124	Significant

Fig. 11.18

From the comparison table shown in Fig. 11.18, we observe that the effect of qualification level Graduate (G) with respect to the training is significantly different from that of Post Graduate (PG) level. However, Graduate (G) level is not significantly different from Doctorate (D) level. Also, Post Graduate (PG) level is significantly different from Doctorate (D) level. Similarly, the final output for the pair-wise comparison for designation level is shown in Fig. 11.19.

	A	B	C	D	E	F	G	H	I
58	Comparison table								
59	Pair of Designations	n_i	n_j	\bar{T}_i	\bar{T}_j	$ \bar{T}_i - \bar{T}_j $	$t_{(27), 0.025}$	CD	Inference
60	SM, ZM	12	12	78	84.5	6.5	2.0518	5.70124	Significant
61	SM, RM	12	12	78	76.4167	1.5833	2.0518	5.70124	Insignificant
62	ZM, RM	12	12	84.5	76.4167	8.0833	2.0518	5.70124	Significant

Fig. 11.19

From the comparison table shown in Fig. 11.19, we observe that the effect of the designation level Sales manager (SM) with respect to the

training is significantly different from that of Zonal manager (ZM). However, the effect of the designation Sales manager (SM) is not significantly different from that of Regional manager (RM). Also, the effect of designation Zonal manager (ZM) is significantly different from that of Regional manager (RM).

You can now try the following exercises for practice.



Activity

Apply the suitable ANOVA test with the help of MS Excel 2007 for the following exercises and interpret the results:

- A1)** Examples 1 and 2 given in Unit 6 of MST-005.
- A2)** Exercises E1, E2 and E3 given in Unit 6 of MST-005.
- A3)** Example 1 given in Unit 7 of MST-005.
- A4)** Exercises E1 and E2 given in Unit 7 of MST-005.
- A5)** Examples 1 given in Unit 8 of MST-005.
- A6)** Exercises E1 given in Unit 8 of MST-005.

Match the results with the manual computation of data carried out in Units 6, 7 and 8 of MST-005.



Continuous Assessment 11

1. A researcher wishes to develop criteria for talent identification in different sports. In order to find the effectiveness of height in different sports, the heights of sportsmen in three different sports, namely cricket, hockey and volley ball were recorded. The data are shown in Table 4.

Table 4: Heights (in feet) of different sportsmen

Cricket	Hockey	Volley Ball
5.9	6.0	6.3
6.1	5.6	6.2
5.8	5.8	6.0
5.7	6.2	5.9
5.8	5.7	6.2
6.0	5.8	6.3
6.2	5.9	6.2
5.7	5.7	6.0
5.8	5.7	6.2
5.9	5.9	6.2
5.8	6.3	
5.7		

Assuming that the heights of the sportsmen are normally distributed in each sport with approximately equal variances, apply the one-way ANOVA to test the hypothesis that the average height of the sportsmen in all three sports is the same at 1% level of significance.

2. A leading shirt manufacturer has 1000 showrooms across the country. The company wishes to know the average difference in sales of these

showrooms. It also wishes to know the average sales difference between salesmen. For ascertaining the productivity of different salesmen, the company has adopted the practice of retaining one salesman for three months in one showroom. The company randomly selected four showrooms and four salesmen from each of these showrooms. Table 5 exhibits the average sales (in thousand rupees) from the four showrooms and the individual contribution of the five salesmen in different showrooms.

Table 5: Sales generated by four salesmen from different showrooms

Salesmen	Showrooms			
	Showroom 1	Showroom 2	Showroom 3	Showroom 4
Salesman 1	85	92	80	74
Salesman 2	86	90	75	86
Salesman 3	90	100	92	80
Salesman 4	80	95	82	75

Apply a two-way ANOVA to examine:

- i) whether the salesmen significantly differ in productivity by assuming that the effect of each salesman is normally distributed with approximately equal variances.
 - ii) whether there is a significant difference between the average sales of showrooms by assuming that the sale of each showroom is normally distributed with approximately equal variances at 5% level of significance.
3. A Steel and Iron company produces 8-metre long steel rods, which are used in the construction of buildings. The company has four machines that manufacture steel rods in three shifts. The company's quality control officer wishes to test whether there is a significant difference in the average length of the iron rods by shifts or by machines. The data produced by machines and shifts through a random sampling process were collected by him and are shown in Table 6.

Table 6: Data for length of steel rods for different machines in different shifts

Machine	Shift I	Shift II	Shift III
1	8.05	8.11	8.06
	8.01	8.10	8.04
	8.10	8.06	8.10
2	7.80	7.77	7.90
	7.90	7.90	7.88
	7.95	7.80	7.95
3	8.20	8.22	8.12
	8.15	8.25	8.10
	8.22	8.20	8.16
4	7.80	7.85	7.73
	7.90	7.94	7.80
	7.80	7.96	7.90

Employ a two-way analysis of variance and determine whether there is a significant difference in effects. Take $\alpha = 0.05$.



Home Work: Do It Yourself

- 1) Follow the steps explained in Secs. 11.3, 11.4 and 11.5 to apply the tests on the data of Tables 1, 2 and 3. Take the final screenshots and keep them in your record book.
- 2) Develop the spreadsheets for the exercises of “Continuous Assessment 11” as explained in this lab session. Take screenshots of the final spreadsheets.
- 3) **Do not forget** to keep all screenshots in your record book as these will contribute to your continuous assessment in the Laboratory.

