

## Design of Experiments – Completely randomized designs

In an experiment designed to study the effect of temperature on the mean yield of a chemical synthesis process, five batches were produced at each of three temperature levels selected. The replicated batch results obtained in kg are as below:

Temperature		
50°C	60°C	70°C
34	30	23
24	31	28
36	34	28
39	23	30
32	27	31

In experimental design terminology, the production temperature level is the independent variable or **factor**. Because three temperatures were selected in correspondence to this factor, we say there are three **treatments** associated with this experiment. This chemical synthesis process therefore is an example of a **single-factor experiment** involving a quantifiable factor (temperature level).

This chemical yield experiment is called a **completely randomized design** as this type of design requires that each of the temperature level (or treatment) be assigned randomly to each batch of the production. For example, 50°C temperature level might be assigned on Monday, 70°C temperature level to the Tuesday production and 60°C on Friday, with the other temperature levels randomly selected in the other days and weeks that followed.

The hypotheses we want to test in this completely randomized design are:

$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu$$

$H_1$  : Not all 3 population means (with temperature factor) were equal

where

$$\mu_j = \text{mean of the } j^{\text{th}} \text{ population}$$

We use analysis of variance (ANOVA) to test for the equality of means in this situation. In here, we need to consider two independent estimates of the population variance  $\sigma^2$ :

- **Within-treatments** estimate of population variance
- **Between-treatments** estimate of population variance

Before we go further, it is important to know the following important assumptions made for ANOVA:

1. **For each population, the response variable is normally distributed.** That means the daily production yield must be normally distributed at each temperature level;
2. **The variance of the response variable, denoted  $\sigma^2$ , is the same for all the populations.** That means the variance of yields must be the same for all three temperature levels;
3. **The observations must be independent.** In this example, the daily yield production must be independent of each other.

The formulae for the sample mean and sample variance for treatment  $j$  are as follow:

$$\bar{x}_j = \frac{\sum_{i=1}^{n_j} x_{ij}}{n_j} \quad [1]$$

$$s_{ij}^2 = \frac{\sum_{i=1}^{n_j} (x_{ij} - \bar{x}_j)^2}{n_j - 1} \quad [2]$$

The overall sample mean, denoted  $\bar{x}$ , is the sum of all the observations divided by the total number of observations  $n_T$ . That is:

$$\bar{x} = \frac{\sum_{j=1}^k \sum_{i=1}^{n_j} x_{ij}}{n_T} \quad [3]$$

where

$$n_T = n_1 + n_2 + \dots + n_k \quad [4]$$

If the size of repeats for each sample (treatment) is  $n$ , then  $n_T = kn$ ; in this case, equation [3] is simplified and reduces to:

$$\bar{x} = \frac{\sum_{j=1}^k \bar{x}_j}{k} \quad [5]$$

In other words, whenever the sample sizes are the same, the overall sample mean is just the average of the  $k$  sample means.

In this chemical yield example which consisted of  $n=5$  observations, we have:

Mean  $\bar{x}$  (50°C) = 33; variance  $s^2$  (50°C) = 32.0

Mean  $\bar{x}$  (60°C) = 29; variance  $s^2$  (60°C) = 17.5

Mean  $\bar{x}$  (70°C) = 28; variance  $s^2$  (70°C) = 9.5

and, the overall mean  $\bar{x} = \frac{33+29+28}{3} = 30$

#### **Within-treatments estimate of population variance**

Any experiment constitutes a population. The individual sample means obtained through the treatments are assumed to follow sampling distribution of population variance  $\sigma^2$  unless proven otherwise. Hence, we are going to use ANOVA to estimate its population variance from within-treatments and between-treatments perspectives.

When the sample sizes are equal, the within-treatments estimate of population variance  $\sigma^2$  can be obtained by computing the average of the individual sample variances. Hence, for this chemical yield experiment, we obtain:

$$\text{Within-treatments estimate of } \sigma^2 = \frac{32.0+17.5+9.5}{3} = 19.7$$

This estimate of  $\sigma^2$  is called the **mean square due to error** and is denoted  $MSE$ . The general formula for computing  $MSE$  is:

$$MSE = \frac{\sum_{j=1}^k (n_j - 1)s_j^2}{n_T - k} \quad [6]$$

The numerator in equation [6] is called the **sum of squares due to error** and is denoted  $SSE$ . The denominator of  $MSE$  is referred to as the degrees of freedom associated with  $SSE$ . Hence the formula for  $MSE$  can also be stated

as follows:

MEAN SQUARE DUE TO ERROR

$$MSE = \frac{SSE}{n_T - k} \quad [7]$$

where

$$SSE = \sum_{j=1}^k (n_j - 1) s_j^2 \quad [8]$$

Hence,  $MSE = 19.7$  and  $SSE = 236$ . It may be noted that  $MSE$  is based on the variation within each of the treatments; it is therefore not influenced by whether the null hypothesis  $H_0$  is true. Thus,  $MSE$  always provides an unbiased estimate of the population variance  $\sigma^2$ .

#### Between-treatments estimate of population variance

If the null hypothesis is true, we can think of each of the three means at different temperature levels (33, 29, 28) as values drawn at random from a population (sampling distribution). Hence, we can use the mean and variance of the three mean values to estimate its overall mean and population variance.

When the sample sizes are equal as in this experiment, the best estimate of the mean of the population is the mean or average of the sample means. Thus, an estimate of the overall population mean  $\bar{x}$  is  $(33+29+28)/3 = 30$  and its population variance of this mean,  $s_{\bar{x}}^2$  is 7.0.

Now, because  $\sigma_{\bar{x}}^2 = \frac{\sigma^2}{n}$ , solving for  $\sigma^2$  gives

$$\sigma^2 = n \sigma_{\bar{x}}^2$$

Hence,

$$\text{Estimate of } \sigma^2 = n (\text{Estimate of } \sigma_{\bar{x}}^2) = n s_{\bar{x}}^2 = 5(7.0) = 35.0$$

Note that the between-treatments estimate of  $\sigma^2$  is based on the assumption that the null hypothesis  $H_0$  is true. In other words, each sample (treatment) comes from the same population and hence there is only one sampling distribution of  $\bar{x}$ .

When the sample sizes are equal, this estimate of  $\sigma^2$  is called the mean square due to treatments and is denoted  $MSTr$ . The general formula for computing  $MSTr$  is :

$$MSTr = \frac{\sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2}{k-1} \quad [9]$$

The numerator in equation [9] is called the sum of squares due to treatments and is denoted  $SSTr$ . The denominator,  $k-1$ , represents the degrees of freedom associated with  $SSTr$ . Hence, the mean square due to treatments can be computed by the following formula:

#### MEAN SQUARE DUE TO TREATMENTS

$$MSTr = \frac{SSTr}{k-1} \quad [10]$$

where

$$SSTr = \sum_{j=1}^k n_j (\bar{x}_j - \bar{x})^2 \quad [11]$$

Hence, in this example,  $MSTr = 35$  and  $SSTr = 70$ . If  $H_0$  is true,  $MSTr$  provides an unbiased estimate of  $\sigma^2$  but if the means of the  $k$  populations are not equal (i.e.  $H_0$  is false),  $MSTr$  is not an unbiased estimate of  $\sigma^2$ ; in fact, in that case,  $MSTr$  tends to overestimate  $\sigma^2$ .

#### Comparing the variance estimates by the $F$ test

If the null hypothesis is true and the ANOVA assumptions are valid, the sampling distribution of  $MSTr/MSE$  is an  $F$  distribution with numerator degrees of freedom equal to  $k-1$  and denominator degrees of freedom equal to  $n_T - k$ .

Recall also that if the means of the  $k$  populations are not equal, the value of  $MSTr/MSE$  will be inflated because  $MSTr$  overestimates  $\sigma^2$ . Hence, we will reject  $H_0$  if:

$$\frac{MSTr}{MSE} \geq F_{critical}(v_1 = k-1, v_2 = n_T - k)$$

Let us return to the chemical yield experiment and use a level of significance  $\alpha = 0.05$  to conduct the hypothesis test. The value of the test statistic is:

$$F = \frac{MSTr}{MSE} = \frac{35}{19.7} = 1.78 < 3.89 (F_{critical}, v_1 = 2, v_2 = 12)$$

As the estimated  $F$  value of 1.78 is less than the critical  $F$  value of 3.89, we conclude that  $H_0$  is true, that means there were no significant differences amongst the three population means.

We can now write the result that shows how the total sum of squares  $SST$  is partitioned:

$$SST = SSTr + SSE$$

This result also holds true for the degrees of freedom associated with each of these sums of squares, i.e. the total degrees of freedom is the sum of the degrees of freedom associated with  $SSTr$  and  $SSE$ . We may present the general form of the ANOVA table for a completely randomized design as below:

**ANOVA TABLE FOR A COMPLETELY RANDOMIZED DESIGN**

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Square	$F$
Treatments	$SSTr$	$k - 1$	$MSTr = \frac{SSTr}{k - 1}$	$\frac{MSTr}{MSE}$
Error	$SSE$	$n_T - k$	$MSE = \frac{SSE}{n_T - k}$	
Total	$SST$	$n_T - 1$		

### Analysis of variance and Experimental Design with Excel

Microsoft Excel<sup>®</sup> – Data Analysis Tool can be used to test for the equality of  $k$  population means as shown by the table below:

Batch #	Temperature		
	50 Deg C	60 Deg C	70 Deg C
1	34	30	23
2	24	31	28
3	36	34	28
4	39	23	30
5	32	27	31

Anova: Single Factor				
SUMMARY				
Treatments	Count	Sum	Average	Variance
50 Deg C	5	165	33	32
60 Deg C	5	145	29	17.5
70 Deg C	5	140	28	9.5

ANOVA						
Source of Variation	SS	df	MS	F	P-value	F crit
Between Treatments	70	2	35	1.78	0.210	3.89
Within Treatments	236	12	19.67			
Total	306	14				