

1-Way Completely Randomized ANOVA

James H. Steiger

Department of Psychology and Human Development
Vanderbilt University

1-Way Completely Randomized Design

- 1 Introduction
- 2 An Introductory Example
- 3 The Basic Idea behind ANOVA
- 4 The ANOVA Structural Model
- 5 Computations
 - Computational Formulas
 - Calculation in R
- 6 Distribution of the F Statistic and Power Calculation

Introduction

- In this module, we introduce the single-factor completely randomized analysis of variance design.
- We discover that there are several ways to conceptualize the design.
- For example, we can see the design as a generalization of the 2-sample t -test on independent groups.
- Or, we can see the design as a special case of multiple regression with fixed regressors.

Introduction

- We will investigate several key aspects of any design we investigate:
 - ① Understanding the major hypotheses the design can test;
 - ② Measures of effect size we can extract from the statistical analysis;
 - ③ Relationships between power, precision, sample size, and effect size, and how we exploit them in planning our experiments;
 - ④ Statistical assumptions, restrictions, and robustness properties.

An Introductory Example

- MWL present data from a hypothetical study of memory in their Table 8.1.
- 40 participants were divided randomly into 4 groups.
- Each participant studied a list of 20 words and was tested for recall a day later.
- Each of the 3 experimental groups was instructed to use a different special memorization strategy.
- A 4th group was simply told to try to memorize the list.
- Data are shown in Table 8.1 (MWL, p. 171). (There is an error in the Table: $\bar{Y}_{\bullet j} = 9.95$ should be $\bar{Y}_{\bullet\bullet} = 9.95$.)

An Introductory Example

Table 8.1 Recall scores from a hypothetical memory study

	Control	Loci	Image	Rhyme	
	11	10	13	16	
	4	18	16	9	
	8	6	3	7	
	3	20	6	10	
	11	15	13	9	
	8	9	10	14	
	2	8	13	16	
	5	11	9	3	
	8	12	5	9	
	5	12	19	12	
$\bar{Y}_j =$	6.5	12.1	10.7	10.5	$\bar{Y}_j = 9.95$
$s_j^2 =$	10.056	19.433	25.567	16.722	

ANOVA: The Basic Idea

- At its foundation, ANOVA attempts to assess whether a group of means are all the same. If there are a groups, then the null hypothesis might be written

$$H_0 : \mu_1 = \mu_2 = \dots = \mu_a \quad (1)$$

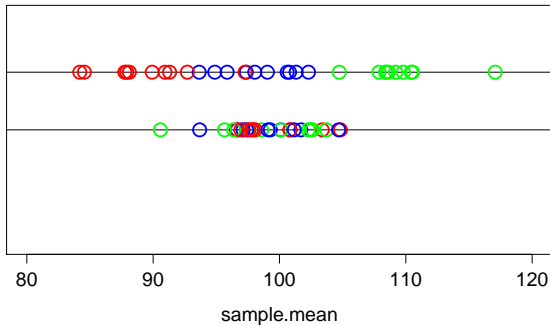
ANOVA: The Basic Idea

- ANOVA does this by comparing the variability of the sample means to what it *should be* if the population means are all the same.
- If the population means are all the same, the sample sizes are all the same (n per group), and the populations all have the same variance σ^2 , then the means represent repeated samples from the same population.
- In that case, variance of the sample means should be approximately σ^2/n . It is crucial to realize that, under all the assumptions, that is the smallest long run variance that these means can show for any set of population mean values.
- On the other hand, if the population means are not all the same, then the variance of the sample means should be larger than σ^2/n , because the spread of the means reflects more than sampling variability.

ANOVA: The Basic Idea

- The plot on the next slide shows 10 sample means taken from each of 3 groups. All populations had standard deviations of 15 and sample sizes of 25. The lower line shows 30 means from 3 populations with identical means of 100,100,100. Group 1 is red, Group 2 in blue, Group 3 in green. The dispersion that you see along the line is totally due to within group sampling variation.
- The upper line shows 30 means sampled from populations with means of 90,100,and 110. The dispersion of the means on the top line reflects both within group sampling variation and the spread of the means between groups.

ANOVA: The Basic Idea



ANOVA: The Basic Idea

- As a consequence of the facts from the preceding slides, we can conceptualize the null hypothesis in ANOVA in an entirely different (but equivalent) way from Equation 1
- To emphasize that this is a *one-tailed hypothesis*, I write both the null and alternative hypotheses.

$$H_0 : \sigma_{\bar{Y}}^2 = \frac{\sigma^2}{n} \quad H_1 : \sigma_{\bar{Y}}^2 > \frac{\sigma^2}{n} \quad (2)$$

ANOVA: The Basic Idea

- As a consequence of this, and a lot of statistical machinery, we arrive at an F statistic that may be written, for a groups,

$$F_{a-1, a(n-1)} = \frac{s_Y^2}{\hat{\sigma}^2/n} = \frac{ns_Y^2}{\hat{\sigma}^2} \quad (3)$$

- We simply compute the sample variance of the sample means, and divide by an estimate of σ^2/n . With equal sample sizes, $\hat{\sigma}^2$ is simply the mean of the group variances.

ANOVA: The Basic Idea

- Here are the calculations in R.

```
> memory.data <- read.csv("Table 8_1 Memory Data.csv")
> means <- with(memory.data,
+   aggregate(Score, by=list(Method), FUN=mean))$x
> var.xbars <- var(means)
> variances <- with(memory.data,
+   aggregate(Score, by=list(Method), FUN=var))$x
> var.e <- mean(variances)
> n <- 10
> MS.A <- n*var.xbars
> MS.e <-var.e
> F <- MS.A/MS.e
> c(MS.A,MS.e,F)
```

```
[1] 57.966667 17.944444 3.230341
```

The ANOVA Structural Model

- An alternative way of conceptualizing ANOVA begins with a *structural model* that includes extensive assumptions and restrictions.
- Following the notation in RDASA3, we shall assume a balanced design with n subjects per cell, and a cells.
- Let Y_{ij} be the score for the i th person in the j th cell. Then

$$Y_{ij} = \mu + \alpha_j + \varepsilon_{ij} \quad (4)$$

where

- 1 μ is the “grand mean” or overall population mean,
 - 2 $\alpha_j = \mu_j - \mu$ is the difference between the mean of population j and the overall mean,
 - 3 ε_{ij} is independent, normal error. The ε_{ij} are i.i.d. $N(0, \sigma_e^2)$.
- These assumptions imply that the observations in each group are normally distributed, and that group variances are equal.
 - MWL summarize the model assumptions in their Box 8.1.

The ANOVA Structural Model

Box 8.1 Parameter Definitions and Assumptions

1. *The parent population mean, μ .* This is the grand mean of the treatment populations selected for this study and is a constant component of all scores in the a populations. It is the average of the treatment population means:

$$\mu = \sum_{j=1}^a \mu_j / a$$

2. *The effect of treatment A_j , α_j .* This equals $\mu_j - \mu$ and is a constant component of all scores obtained under A_j but may vary over treatments (levels of j).

2.1 Because the deviation of all scores about their mean is zero, $\sum_j \alpha_j = 0$.

2.2 If the null hypothesis is true, all $\alpha_j = 0$.

2.3 The population variance of the treatment effects is $\sigma_A^2 = \sum_{j=1}^a \alpha_j^2 / a$.

3. *The error, ε_{ij} .* This is the deviation of the i^{th} score in group j from μ_j and reflects uncontrolled, or chance, variability. It is the only source of variation within the j^{th} group, and if the null hypothesis is true, the only source of variation within the data set. We assume that

3.1 The ε_{ij} are independently distributed; i.e., the probability of sampling some value of ε_{ij} does not depend on other values of ε_{ij} in the sample.

3.2 The ε_{ij} are normally distributed in each of the a treatment populations. Also, because $\varepsilon_{ij} = Y_{ij} - \mu_j$, the mean of each population of errors is zero; i.e., $E(\varepsilon_{ij}) = 0$.

3.3 The distribution of the ε_{ij} has variance σ_e^2 (error variance) in each of the a treatment populations; i.e., $\sigma_1^2 = \dots = \sigma_a^2 = \dots = \sigma_e^2$. This is the assumption of *homogeneity of variance*. The error variance is the average squared error; $\sigma_e^2 = E(\varepsilon_{ij}^2)$.

Sums of Squares and Source Table

Table 8.3 The analysis of variance for the one-factor between-subjects design

(a) General form of the ANOVA

Source	df	SS	MS	F
Total	$an - 1$	$\sum_{j=1}^a \sum_{i=1}^n (Y_{ij} - \bar{Y}_{..})^2$		
A	$a - 1$	$n \sum_{j=1}^a (\bar{Y}_{.j} - \bar{Y}_{..})^2$	SS_A/df_A	$MS_A/MS_{S/A}$
S/A	$a(n - 1)$	$\sum_{j=1}^a \sum_{i=1}^n (Y_{ij} - \bar{Y}_{.j})^2$	$SS_{S/A}/df_{S/A}$	

(b) ANOVA of the data of Table 8.1

Source	Sum of squares	df	Mean square	F	p -value
Method	173.90	3	57.967	3.230	.034
Error	646.00	36	17.944		
Total	819.90	39			

Calculation in R

- Let's load in the data from Table 8.1.

```
> memory.data <- read.csv("Table 8_1 Memory Data.csv")
> head(memory.data)
```

```
  Method Score
1 Control    11
2 Control     4
3 Control     8
4 Control     3
5 Control    11
6 Control     8
```

```
> str(memory.data)
```

```
'data.frame':      40 obs. of  2 variables:
 $ Method: Factor w/ 4 levels "Control","Image",...: 1 1 1 1 1 1 1 1 1 1
 $ Score : int  11 4 8 3 11 8 2 5 8 5 ...
```

- We see from the above that `memory.data` has two variables, one of which is an integer variable, the other a “factor.”

Calculation in R

- It is vitally important that the factor variables are typed as factors.
- In this case, we are ready to go.

```
> summary(aov(Score ~ Method, data=memory.data))
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Method	3	173.9	57.97	3.23	0.0336 *
Residuals	36	646.0	17.94		

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Distribution of the F Statistic

- The F statistic with $a - 1$ and $a(n - 1)$ degrees of freedom has a noncentral F distribution with noncentrality parameter

$$\lambda = n \sum_{j=1}^a \left(\frac{\alpha_j}{\sigma} \right)^2 \quad (5)$$