

# 1.017/1.010 Class 21

## Multifactor Analysis of Variance

---

### Multifactor Models

We often wish to consider several factors contributing to variability rather than just one. Extend concepts of single factor ANOVA to multiple factors. Focus on the two-factor case.

Suppose there  $I$  treatments for Factor  $A$  and  $J$  treatments for factor  $B$ ., giving  $IJ$  random variables described by CDFs  $F_{x_{ij}}(x_{ij})$ . The different  $F_{x_{ij}}(x_{ij})$  are assumed **identical** (except for their means) and **normally distributed** (check this, as in single factor case).

A random sample  $[x_{ij1}, x_{ij2}, \dots, x_{ijk}]$  of size  $K$  is obtained for treatment combination  $(i, j)$ . Two-factor model describing  $x_{ijk}$ :

$$x_{ijk} = \mu_{ij} + e_{ijk} = \mu + a_i + b_j + c_{ij} + e_{ijk}$$

$\mu_{ij} = E[x_{ijk}] = \mu + a_i + b_j + c_{ij}$  = unknown mean of  $x_{ijk}$  (for all  $k$ )

$\mu$  = unknown **grand mean** (average of  $\mu_i$ 's).

$a_i$  = unknown **main effects of Factor A**

$b_j$  = unknown **main effects of Factor B**

$c_{ij}$  = unknown **interactions** between Factors  $A$  and  $B$

$e_{ijk}$  = **random residual** for treatment  $i$ , replicate  $j$

$E[e_{ijk}] = 0, Var[e_{ijk}] = \sigma^2$ , for all  $i, j, k$

Note that  $c_{ij}$  can only be distinguished from  $e_{ijk}$  if number of replicates  $K > 1$ . Constraints:

$$\sum_{i=1}^I a_i = \sum_{j=1}^J b_j = 0 \quad \sum_{i=1}^I c_{ij} = 0 \quad \forall j \quad \sum_{j=1}^J c_{ij} = 0 \quad \forall i$$

Objective is to estimate/test values of  $a_i$ 's,  $b_j$ 's, and  $c_{ij}$ 's, which are distributional parameters for the  $F_{x_{ij}}(x_{ij})$ 's.

### Formulating the Problem as a Hypothesis Test

Formulate three sum-of-squares hypotheses that insure that all  $a_i$ 's, all  $b_j$ 's, or all  $c_{ij}$ 's are zero:

$$\text{HOA} : \sum_{i=1}^I a_i^2 = 0$$

$$\text{HOB} : \sum_{j=1}^J b_j^2 = 0$$

$$\text{HOAB} : \sum_{j=1}^J \sum_{i=1}^I c_i^2 = 0$$

Derive test statistics based on sums-of-squares of data.

## Sums-of-Squares Computations

Define treatment and grand sample means:

$$m_{xi} = \frac{1}{JK} \sum_{j=1}^J \sum_{k=1}^K x_{ijk} = \bar{x}_{i..} \quad m_{xj} = \frac{1}{IK} \sum_{i=1}^I \sum_{k=1}^K x_{ijk} = \bar{x}_{.j.}$$

$$m_{xij} = \frac{1}{K} \sum_{k=1}^K x_{ijk} = \bar{x}_{ij.}$$

$$m_x = \frac{1}{IJK} \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K x_{ijk} = \bar{x}_{...}$$

Test statistics are computed from sums-of-squares:

$$SSE = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (x_{ijk} - m_{xij})^2$$

$$SSA = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (m_{xi} - m_x)^2 \quad SSB = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (m_{xj} - m_x)^2$$

$$SSAB = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (m_{xij} - m_{xi} - m_{xj} + m_x)^2$$

$$SST = \sum_{i=1}^I \sum_{j=1}^J \sum_{k=1}^K (x_{ijk} - m_x)^2 = SSE + SSA + SSB + SSAB$$

Corresponding mean-sums-of-squares are:

$$MSE = \frac{SSE}{IJ(K-1)}$$

$$MSA = \frac{SSA}{I-1} \quad MSB = \frac{SSB}{J-1}$$

$$MSAB = \frac{SSAB}{(I-1)(J-1)}$$

Expected values of these mean-sums-of-squares show depends on main effects and interactions:

$$E[MSE] = \sigma^2$$

$$E[MSA] = \sigma^2 + \frac{JK}{I-1} \sum_{i=1}^I a_i^2 \quad E[MSB] = \sigma^2 + \frac{IK}{I-1} \sum_{j=1}^J b_j^2$$

$$E[MSAB] = \sigma^2 + \frac{K}{(I-1)(J-1)} \sum_{i=1}^I \sum_{j=1}^J c_{ij}^2$$

## Test Statistic

Use ratios as test statistics for the three hypotheses:

$$F_A(MSA, MSE) = \frac{MSA}{MSE}$$

$$F_B(MSB, MSE) = \frac{MSB}{MSE}$$

$$F_{AB}(MSAB, MSE) = \frac{MSAB}{MSE}$$

When  $H_0$  is true each statistic follows **F distribution** with degree of freedom parameters  $v_A = I - 1$ ,  $v_B = J - 1$ ,  $v_{AB} = (I - 1)(J - 1)$ , and  $v_E = IJ(K-1)$ .

**One-sided** rejection regions

$$ROA: F(MSA, MSE) \geq F_{F, v_A, v_E}^{-1}[\alpha]$$

$$ROB: F(MSB, MSE) \geq F_{F, v_B, v_E}^{-1}[\alpha]$$

$$ROAB: F(MSAB, MSE) \geq F_{F, v_{AB}, v_E}^{-1}[\alpha]$$

**One-sided** p-values:

$$p_A = 1 - F_{\mathcal{F}, v_A, v_E}[\mathcal{F}(MSA, MSE)]$$

$$p_B = 1 - F_{\mathcal{F}, v_A, v_E}[\mathcal{F}(MSB, MSE)]$$

$$p_{AB} = 1 - F_{\mathcal{F}, v_{AB}, v_E}[\mathcal{F}(MSAB, MSE)]$$

**Unbalanced ANOVA** problems with **different sample sizes for different treatments** can be handled by modifying formulas slightly.

## Two-Factor ANOVA Tables

Source	SS	df	MS	$\mathcal{F}$	$p$
Factor $A$	$SSA$	$v_A = I - 1$	$MSA = SSA/v_A$	$\mathcal{F}_A = MSA/MSE$	$p = 1 - F_{\mathcal{F}, v_A, v_E}(\mathcal{F})$
Factor $B$	$SSB$	$v_B = J - 1$	$MSB = SSB/v_B$	$\mathcal{F}_B = MSB/MSE$	$p = 1 - F_{\mathcal{F}, v_B, v_E}(\mathcal{F})$
Interaction $AB$	$SSAB$	$v_{AB} = (I-1)(J-1)$	$MSAB = SSAB/v_{AB}$	$\mathcal{F}_{AB} = MSAB/MSE$	$p = 1 - F_{\mathcal{F}, v_{AB}, v_E}(\mathcal{F})$
Error	$SSE$	$v_E = IJ(K-1)$	$MSE = SSE/v_E$		
Total	$SST$	$v_T = IJK - 1$			

## Exercise: Two Factor ANOVA

Relevant MATLAB functions: `normplot`, `anova2`

## Concepts and Definitions

Objective: Identify factors responsible for variability in observed data

Specify one or more **factors** that could account for variability (e.g. location, time, etc.). Each factor is associated with a particular set of populations or **treatments** (e.g. particular sampling stations, sampling days, etc.). **One-way analysis of variance** (ANOVA) considers only a single factor.

Suppose a random sample  $[x_{i1}, x_{i2}, \dots, x_{ij}]$  is obtained for treatment  $i$ . There are  $i=1, \dots, I$  treatments (e.g. each treatment may correspond to a different sampling location).

Arrange data in a table/array -- rows are treatments, columns are replicates:

$$\begin{bmatrix} x_{11}, x_{12}, \dots, x_{1j} \\ x_{21}, x_{22}, \dots, x_{2j} \\ \vdots \\ x_{I1}, x_{I2}, \dots, x_{Ij} \end{bmatrix}$$

Here we assume each treatment has same number of replicates  $J$ . The ANOVA procedure may be generalized to allow different number of replicates for each treatment.

Each random sample has a CDF  $F_{xi}(x_i)$ . The different  $F_{xi}(x_i)$  are assumed **identical** except for their means, which may differ. Classical ANOVA also assumes that all data are **normally distributed**.

Each random variable  $x_{ij}$  is decomposed into several parts, as specified by the following **one-factor model**:

$$x_{ij} = \mu_i + e_{ij} = \mu + a_i + e_{ij}$$

$\mu_i = E[x_{ij}]$  is unknown mean of  $x_i$  (for all  $j$ ).

$\mu$  = unknown **grand mean** (average of  $\mu_i$ 's).

$a_i = \mu_i - \mu$  = unknown deviation of treatment mean from grand mean (often called an **effect**)

$e_{ij}$  = **random residual** for treatment  $i$ , replicate  $j$

$E[e_{ij}] = 0, Var[e_{ij}] = \sigma^2$ , for all  $i, j$

Objective is to estimate/test values of  $a_i$ 's, which are the unknown distributional parameters of the  $F_{xi}(x_i)$ 's.

## Formulating the Problem as a Hypothesis Test

If the factor does not affect variability in the data then all  $a_i$ 's = 0. Use hypothesis test:

$$H_0: a_1 = a_2 = \dots = a_I = 0$$

It is better to test all  $a_i$  simultaneously than individually or in pairs. Test that sum-of-squared  $a_i$ 's = 0.

$$H_0: \sum_{i=1}^I a_i^2 = 0$$

Derive a test statistic based on sums-of-squares of data.

## Sums-of-Squares Computations

Define the sample treatment and grand means:

$$m_{xi} = \frac{1}{J} \sum_{j=1}^J x_{ij} = \bar{x}_i.$$

$$m_x = \frac{1}{IJ} \sum_{i=1}^I \sum_{j=1}^J x_{ij} = \bar{x}_{..}$$

The **total sum-of-squares**  $SST$  measures variability of  $x_{ij}$  around  $m_x$ :

$$\begin{aligned} SST &= \sum_{i=1}^I \sum_{j=1}^J (x_{ij} - m_x)^2 \\ &= \sum_{i=1}^I \sum_{j=1}^J (x_{ij} - m_{xi})^2 + \sum_{i=1}^I \sum_{j=1}^J (m_{xi} - m_x)^2 \\ &= SSE + SSTr \end{aligned}$$

$SST$  can be divided into **error sum-of-squares**  $SSE$  and **treatment sum-of-squares**  $SSTr$ .

$SSE$  measures variability of  $x_{ij}$  around  $m_{xi}$ , within treatments:

$$SSE = \sum_{i=1}^I \sum_{j=1}^J (x_{ij} - m_{xi})^2$$

$SSTr$  measures variability of  $m_{xi}$  around  $m_x$ , across treatments:

$$SSTr = \sum_{i=1}^I \sum_{j=1}^J (m_{xi} - m_x)^2$$

Error and treatment **mean squared values**:

$$MSE = \frac{SSE}{I(J-1)}$$

$$MStr = \frac{SStr}{I-1}$$

$$E[MSE] = \sigma^2$$

$$E[MStr] = \sigma^2 + \frac{J}{I-1} \sum_{i=1}^I a_i^2$$

*MSE* is an unbiased estimate of  $\sigma^2$ , even if  $a_i$ 's are not zero.  
*MStr* is an unbiased estimate of  $\sigma^2$ , **only** if all  $a_i$ 's are zero.

## Test Statistic

Use ratio *MStr* / *MSE* as a test statistic:

$$F(MSE, MStr) = \frac{MStr}{MSE}$$

When  $H_0$  is true and  $x_{ij}$ 's are **normally distributed** this statistic follows *F* **distribution** with  $v_{Tr} = I - 1$  and  $v_E = I(J-1)$  degrees of freedom. Check normality by plotting  $(x_{ij} - m_{xi})$  with `normplot`.

**One-sided** rejection region (rejects only if *MStr* is large):

$$R_0 : F(MSE, MStr) \geq F_{F, v_{Tr}, v_E}^{-1}[\alpha]$$

**One-sided** *p*-value:

$$p = 1 - F_{F, v_{Tr}, v_E}[F(MSE, MStr)]$$

**Unbalanced** ANOVA problems with **different sample sizes for different treatments** can be handled by modifying formulas slightly (see Devore, Section 10.3).

## Single Factor ANOVA Tables

Above calculations are typically summarized in an **ANOVA** table:

Source	SS	df	MS	$\mathcal{F}$	$p$
Treatments	$SSTr$	$\nu_{Tr} = I-1$	$MSTr = SSTr/\nu_{Tr}$	$\mathcal{F} = MSTr/MSE$	$p = 1 - F_{\mathcal{F}, \nu_{Tr}, \nu_E}(\mathcal{F})$
Error	$SSE$	$\nu_E = I(J-1)$	$MSE = SSE/\nu_E$		
Total	$SST$	$\nu_T = IJ-1$	$MST = SST/\nu_T$		

## Example -- Effect of Season on Oxygen Level

Consider following set of dissolved oxygen concentration data ( $x_{ij}$ ) obtained in 4 different seasons/treatments (rows), 6 replicates per season (columns):

5.62	6.12	6.62	6.21	7.08	5.36
7.70	8.31	8.80	8.24	7.87	7.44
2.52	5.44	4.94	2.99	4.39	4.44
6.77	6.65	6.01	6.26	7.09	6.05

Use a single factor ANOVA to determine if season has a significant impact on oxygen variability.

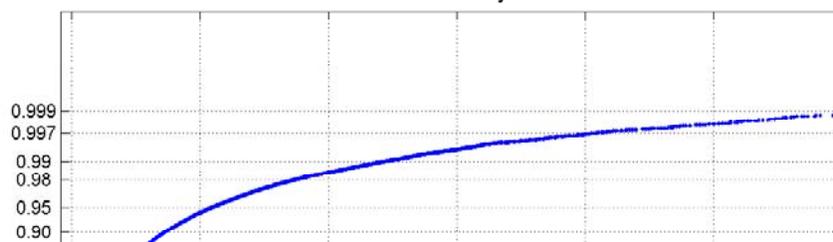
The MATLAB `anova1` function derives the error and treatment sums of squares and computes  $p$  value. **When using `anova1` be sure to transpose the data array** (MATLAB requires treatments in columns and replicates in rows).

Results are presented in this standard single factor **ANOVA table**:

Source	SS	df	MS=SS/df	$\mathcal{F}$	$p$
Treatments	47.1642	3	15.7214	29.8	1.4E-7
Error	10.5518	20	0.5276		
Total	57.716	23			

The very low  $p$  value indicates that seasonality is **highly significant** in this case. Note that  $MSTr$ , which depends on the  $a_i$ 's, is much larger than  $MSE$

F CDF,  $\nu_{Tr} = 3$ ,  $\nu_E = 5$





*Copyright 2003 Massachusetts Institute of Technology  
Last modified Oct. 8, 2003*