Probability and Mathematical Statistics 503 Stat

Lecture 12

Analysis of Variance

Chapter Goals

After completing this chapter, you should be able to:

- Recognize situations in which to use analysis of variance
- Understand different analysis of variance designs
- Perform a one-way and two-way analysis of variance and interpret the results



 Evaluate the difference among the means of three or more groups

Examples: Average production for 1st, 2nd, and 3rd shifts Expected mileage for five brands of tires

- Assumptions
 - Populations are normally distributed
 - Populations have equal variances
 - Samples are randomly and independently drawn

Hypotheses of One-Way ANOVA

$$H_0: \mu_1 = \mu_2 = \mu_3 = \dots = \mu_K$$

All population means are equal

i.e., no variation in means between groups

• $H_1: \mu_i \neq \mu_i$ for at least one i, j pair

- At least one population mean is different
- i.e., there is variation between groups
- Does not mean that all population means are different (some pairs may be the same)





All Means are the same: The Null Hypothesis is True (No variation between groups)







Total variation can be split into two parts:

SST = Total Sum of Squares Total Variation = the aggregate dispersion of the individual data values across the various groupsSSW = Sum of Squares Within Groups Within-Group Variation = dispersion that exists among the data values within a particular groupSSG = Sum of Squares Between Groups Between-Group Variation = dispersion between the group sample means



$$SST = \sum_{i=1}^{N} \sum_{j=1}^{N} (x_{ij} - \overline{x})^2$$

Where:

SST = Total sum of squares

K = number of groups (levels or treatments)

- n_i = number of observations in group i
- $x_{ij} = j^{th}$ observation from group i
- \overline{x} = overall sample mean



$$SST = (x_{11} - \bar{x})^2 + (X_{12} - \bar{x})^2 + ... + (x_{Kn_{K}} - \bar{x})^2$$



Within-Group Variation

$$SST = SSW + SSG$$
$$SSW = \sum_{i=1}^{K} \sum_{j=1}^{n_i} (x_{ij} - \overline{x}_i)^2$$

Where:

SSW = Sum of squares within groups

K = number of groups

- n_i = sample size from group i
- \overline{x}_i = sample mean from group i
- $x_{ij} = j^{th}$ observation in group i

Within-Group Variation

(continued)

$$SSW = \sum_{i=1}^{K} \sum_{j=1}^{n_i} (\mathbf{x}_{ij} - \overline{\mathbf{x}}_i)^2$$

Summing the variation within each group and then adding over all groups





Mean Square Within = SSW/degrees of freedom



$$SSW = (x_{11} - \bar{x}_1)^2 + (x_{12} - \bar{x}_1)^2 + ... + (x_{Kn_K} - \bar{x}_K)^2$$





Where:

SSG = Sum of squares between groups

K = number of groups

- n_i = sample size from group i
- \overline{x}_i = sample mean from group i
- \overline{x} = grand mean (mean of all data values)

Between-Group Variation

(continued)

$$|SSG = \sum_{i=1}^{K} n_i (\overline{x}_i - \overline{x})^2|$$

Variation Due to Differences Between Groups

μ

 μ_{i}

$$MSG = \frac{SSG}{K-1}$$

Mean Square Between Groups = SSG/degrees of freedom



$$SSG = n_1(\overline{x}_1 - \overline{x})^2 + n_2(\overline{x}_2 - \overline{x})^2 + \ldots + n_K(\overline{x}_K - \overline{x})^2$$





$$MST = \frac{SST}{n-1}$$

$$MSW = \frac{SSW}{n-K}$$

$$MSG = \frac{SSG}{K-1}$$

One-Way ANOVA Table

Source of Variation	SS	df	MS (Variance)	F ratio
Between Groups	SSG	K - 1	$MSG = \frac{SSG}{K - 1}$	F = MSG MSW
Within Groups	SSW	n - K	MSW = $\frac{SSW}{n - K}$	
Total	SST = SSG+SSW	n - 1		

- K = number of groups
- n = sum of the sample sizes from all groups
- df = degrees of freedom

One-Factor ANOVA F Test Statistic

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

H₁: At least two population means are different

Test statistic

$$F = \frac{MSG}{MSW}$$

MSG is mean squares between variances MSW is mean squares within variances

- Degrees of freedom
 - $df_1 = K 1$ (K = number of groups)
 - $df_2 = n K$ (n = sum of sample sizes from all groups)



- The F statistic is the ratio of the between estimate of variance and the within estimate of variance
 - The ratio must always be positive
 - df₁ = K -1 will typically be small
 - df₂ = n K will typically be large





One-Factor ANOVA F Test Example

You want to see if three different golf clubs yield different distances. You randomly select five measurements from trials on an automated driving machine for each club. At the .05 significance level, is there a difference in mean distance?

Club 1	<u>Club 2</u>	<u>Club 3</u>
254	234	200
263	218	222
241	235	197
237	227	206
251	216	204



One-Factor ANOVA Example: Scatter Diagram

	<u>Club 1</u>	Club 2	<u>Club 3</u>					
	254	234	200					
	263	218	222					
	241	235	197					
	237	227	206					
	251	216	204					
X	$_{1} = 249.2$	$\bar{x}_2 = 226.0$	$\overline{X}_3 = 208$	5.8				
		x = 227.0						



One-Factor ANOVA Example Computations



 $SSG = 5 (249.2 - 227)^2 + 5 (226 - 227)^2 + 5 (205.8 - 227)^2 = 4716.4$ $SSW = (254 - 249.2)^2 + (263 - 249.2)^2 + ... + (204 - 205.8)^2 = 1119.6$ MSG = 4716.4 / (3-1) = 2358.2 2358.2 = 25.27593.3

MSW = 1119.6 / (15-3) = 93.3

One-Factor ANOVA Example Solution





ANOVA -- Single Factor: Excel Output

EXCEL: data | data analysis | ANOVA: single factor

SUMMARY						
Groups	Count	Sum	Average	Variance		
Club 1	5	1246	249.2	108.2		
Club 2	5	1130	226	77.5		
Club 3	5	1029	205.8	94.2		
ANOVA						
Source of Variation	SS	df	MS	F	P-value	F crit
Between Groups	4716.4	2	2358.2	25.275	4.99E-05	3.89
Within Groups	1119.6	12	93.3			
Total	5836.0	14				



Two-Way Analysis of Variance

Examines the effect of

- Two factors of interest on the dependent variable
 - e.g., Percent carbonation and line speed on soft drink bottling process
- Interaction between the different levels of these two factors
 - e.g., Does the effect of one particular carbonation level depend on which level the line speed is set?



- Assumptions
 - Populations are normally distributed
 - Populations have equal variances
 - Independent random samples are drawn



Two Factors of interest: A and B

- K = number of groups of factor A
- H = number of levels of factor B

(sometimes called a blocking variable)

	Group					
Block	1	2		K		
1	Х ₁₁	Х ₂₁		x _{K1}		
2	x ₁₂	X ₂₂		x _{K2}		
			-			
Н	x _{1H}	X _{2H}	•••	x _{KH}		



- Let x_{ji} denote the observation in the jth group and ith block
- Suppose that there are K groups and H blocks, for a total of n = KH observations
- Let the overall mean be $\overline{\mathbf{x}}$
- Denote the group sample means by

$$\overline{x}_{j\bullet}$$
 (j=1,2,...,K)

Denote the block sample means by

$$\overline{x}_{\bullet i}$$
 (i = 1,2,...,H)



Two-Way Sums of Squares

The sums of squares are

otal:
$$SST = \sum_{j=1}^{K} \sum_{i=1}^{H} (\mathbf{x}_{ji} - \overline{\mathbf{x}})^2$$

Between - Groups :

$$SSG = H \sum_{j=1}^{K} (\overline{x}_{j \bullet} - \overline{x})^2$$

H – 1

Between - Blocks :

Error :

$$SSB = K \sum_{i=1}^{H} (X_{\bullet i} - \overline{X})^2$$

i=1 i=1

 $SSE = \sum_{i=1}^{K} \sum_{j=1}^{H} (X_{ji} - \overline{X}_{j\bullet} - \overline{X}_{\bullet i} + \overline{X})^{2}$

$$(K - 1)(K - 1)$$



The mean squares are

$$MST = \frac{SST}{n-1}$$
$$MSG = \frac{SST}{K-1}$$
$$MSB = \frac{SST}{H-1}$$
$$MSE = \frac{SSE}{(K-1)(H-1)}$$

Two-Way ANOVA: The F Test Statistic

H₀: The K population group means are all the same





General Two-Way Table Format

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Squares	F Ratio
Between groups	SSG	K – 1	$MSG = \frac{SSG}{K-1}$	$\frac{MSG}{MSE}$
Between blocks	SSB	H – 1	$MSB = \frac{SSB}{H-1}$	MSB MSE
Error	SSE	(K – 1)(H – 1)	$MSE = \frac{SSE}{(K-1)(H-1)}$	
Total	SST	n - 1		





Example:

Four methods of manufacturing penicillin were compare in a randomized block design. The block are blends of the raw material, corn steep liquor, known to be quite variable. The yield of each method for five blends is given below.

	blend (block)				
Method	1	2	3	4	5
А	89	84	81	87	79
B	88	77	87	92	81
С	97	92	87	89	80
D	94	79	85	84	88

At level $\alpha = 0.01$,

- 1 are there significant differences between the blends (blocks)?
- 2 are there significant differences between the Methods (treatments)?

Example

		blend (block)				
Method	1	2	3	4	5	Method mean
A B C D	89 88 97 94	84 77 92 79	81 87 87 85	87 92 89 84	79 81 80 88	$ar{y}_{1.} = 84$ $ar{y}_{2.} = 85$ $ar{y}_{3.} = 89$ $ar{y}_{4.} = 86$
blend mean	$ \bar{y}_{.1} = 92$	$\bar{y}_{.2} = 83$	$\bar{y}_{.3} = 85$	$\bar{y}_{.4} = 88$	$\bar{y}_{.5} = 82$	$\bar{y}_{} = 86$

We have $k = 4, b = 5, \bar{y}_{..} = \frac{1}{4 \times 5} \sum_{i=1}^{4} \sum_{j=1}^{5} y_{ij} = \frac{89+84+81+87+79+...+94+79+85+84+88}{20} = \frac{1720}{20}$. $SSA = b \sum_{i=1}^{k} (\bar{y}_{i.} - \bar{y}_{..})^{2} = 5 \left[(84-86)^{2} + (85-86)^{2} + (89-86)^{2} + (86-86)^{2} \right] = 5 \left[4+1+9+0 \right] = 70$. $SSB = k \sum_{j=1}^{b} (\bar{y}_{.j} - \bar{y}_{..})^{2} = 4 \left[(92-86)^{2} + (83-86)^{2} + (85-86)^{2} + (88-86)^{2} + (82-86)^{2} \right] = 4 \left[36+9+1+4+16 \right] = 264$. $SST = \sum_{i=1}^{k} \sum_{j=1}^{b} (y_{ij} - \bar{y}_{..})^{2} = \sum_{i=1}^{4} \sum_{j=1}^{5} (y_{ij} - \bar{y}_{..})^{2} = (89-86)^{2} + (84-86)^{2} + \ldots + (84-86)^{2} + (88-86)^{2} = 9+4+\ldots + 4+4 = 560$ SSE = SST - SSA - SSB = 560 - 70 - 264 = 226

Example

ANOVA of CRD table

(Source)	(SS)	(df)	(MS)	F
Methods (A)	70	3	23.333	1.239
blends (B)	264	4	66.000	3.504
Error	226	12	18.833	
Total	560	19		

Test the blend (blocks) effects

 $H_0: \mu_{.1} = \mu_{.2} = \ldots = \mu_{.5}$

 H_1 :at least one blend mean is different.

 $F = 3.504 > F_{0.05,4,12} = 3.26$, then we reject H_0 , which means there do appear to be significant differences between blends.

2 Test the method effects

 $H_0: \mu_{1.} = \mu_{2.} = \ldots = \mu_{.4}$

 H_1 :at least one method mean is different.

 $F = 1.239 < F_{0.05,3,12} = 3.49$, then we fail to reject H_0 , which means there is no significant differences between methods.



Chapter Summary

Described one-way analysis of variance

- The logic of Analysis of Variance
- Analysis of Variance assumptions
- F test for difference in K means
- Described two-way analysis of variance
 - Examined effects of multiple factors
 - Examined interaction between factors