

Statistical Analysis



JOHNS HOPKINS
CAREY BUSINESS SCHOOL

Analysis of Variance

1. History and Uses
2. Meaning of Analysis of Variance
3. Terminology
4. One Way ANOVA
5. Examples
6. Tukey Test
7. Least Significance Difference Test (LSD Test)
8. Other Tests
 - a. Two Way ANOVA
 - b. Factorial Analysis
 - c. Latin Square Design

History and Uses

History

- ANOVA originated in agricultural research
- Was used to test a large number of agrochemicals at one time under the same environmental conditions. Each chemical was a treatment

Uses in Business

- Comparison of advertising methods
- Comparison of sales territories
- Comparison of outputs from manufacturing sites
- Advantages
 - ▶ Small sample sizes can be used
 - ▶ Allows for each treatment to have different numbers of experimental units

Meaning of Analysis of Variance

- The method is comparing the variance between treatments to the variance within treatments.
- It is hypothesis testing for 3 or more populations
- ANOVA follows the same rules as covered in hypothesis testing
 - ▶ The null hypothesis has the equal statements
 - ▶ Uses a significance level (usually $\alpha = .05$ or $.01$)

Terminology

- *Experimental Units*- the objects receiving the treatments
- *Treatments*- exposure to an action
- *Completely Randomized Design*- the experimental units are assigned randomly to the various treatments. Each unit randomly chosen for the study has an equal chance of being assigned to any treatment.
- *Models*- the manner in which treatments are selected:
 - ▶ *Fixed-effects model*- specific treatments are chosen or fixed in advance of the study.
 - ▶ *Random-effects*- the levels (treatments) used in the study are chosen randomly from a population of possible levels.

One Way ANOVA- Completely Randomized Design *example*

Management Director of a large logistics company wants to determine if three different training programs have a different effect on passing a test for shipping hazardous materials

Three Treatments- (training programs) fixed effects model

- In-class
- Online
- Tutor

Experimental Units

- 14 employees
- Unbalanced problem
- Randomly assigned to treatments

Example Data - Test Scores

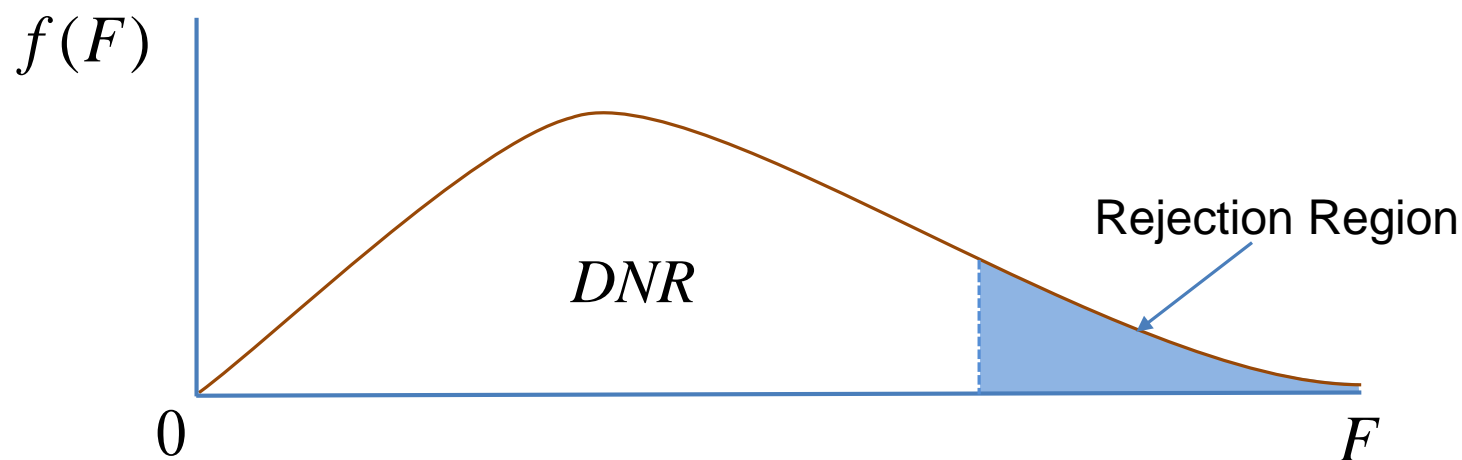
Treatments		
<i>In-class</i>	<i>Online</i>	<i>Tutor</i>
85	80	82
72	84	80
83	81	85
80	78	90
**	82	88
$\bar{x}_1 = 80$	$\bar{x}_2 = 81$	$\bar{x}_3 = 85$

Problem Set-up

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

H_a : not all means are equal

$$\alpha = .05, n = 14, c = 3$$



Equations

$$\overline{\overline{X}} = \frac{\sum_{j=1}^c r_j \overline{X_j}}{n}$$

$$SSTo = \sum_{i=1}^r \sum_{j=1}^c (X_{ij} - \overline{\overline{X}})^2$$

$$SSTr = \sum r_j (\overline{X_j} - \overline{\overline{X}})^2$$

$$SSE = \sum \sum (X_{ij} - \overline{X_j})^2$$

$$SSTo = SSTr + SSE$$

Equations

Treatments		
<i>In-class</i>	<i>Online</i>	<i>Tutor</i>
85	80	82
72	84	80
83	81	85
80	78	90
**	82	88
$\bar{x}_1 = 80$	$\bar{x}_2 = 81$	$\bar{x}_3 = 85$

$$\bar{\bar{X}} = \frac{\sum_{j=1}^c r_j \bar{X}_j}{n} = \frac{4(80) + 5(81) + 5(85)}{14} = 82.14$$

Equations

Treatments		
<i>In-class</i>	<i>Online</i>	<i>Tutor</i>
85	80	82
72	84	80
83	81	85
80	78	90
**	82	88
$\bar{x}_1 = 80$	$\bar{x}_2 = 81$	$\bar{x}_3 = 85$

$$\begin{aligned}\bar{\bar{X}} &= \frac{\sum_{j=1}^c r_j \bar{X}_j}{n} \\ &= \frac{4(80) + 5(81) + 5(85)}{14} \\ &= 82.14\end{aligned}$$

$$\begin{aligned}SSTo &= \sum_{i=1}^r \sum_{j=1}^c (X_{ij} - \bar{\bar{X}})^2 = (85 - 82.14)^2 + (72 - 82.14)^2 + \\ &\quad (83 - 82.14)^2 + \dots + (90 - 82.14)^2 + (88 - 82.14)^2 \\ &= 251.7\end{aligned}$$

Equations

Treatments		
<i>In-class</i>	<i>Online</i>	<i>Tutor</i>
85	80	82
72	84	80
83	81	85
80	78	90
**	82	88
$\bar{x}_1 = 80$	$\bar{x}_2 = 81$	$\bar{x}_3 = 85$

$$\begin{aligned}
 \bar{X} &= \frac{\sum_{j=1}^c r_j \bar{X}_j}{n} \\
 &= \frac{4(80) + 5(81) + 5(85)}{14} \\
 &= 82.14
 \end{aligned}$$

$$\begin{aligned}
 SSTr &= \sum r_j (\bar{X}_j - \bar{X})^2 \\
 &= 4(80 - 82.14)^2 + 5(81 - 82.14)^2 + 5(85 - 82.14)^2 \\
 &= \mathbf{65.7}
 \end{aligned}$$

Equations

Treatments		
<i>In-class</i>	<i>Online</i>	<i>Tutor</i>
85	80	82
72	84	80
83	81	85
80	78	90
**	82	88
$\bar{x}_1 = 80$	$\bar{x}_2 = 81$	$\bar{x}_3 = 85$

$$\begin{aligned}
 SSE &= \sum \sum (X_{ij} - \bar{X}_j)^2 = (85 - 80)^2 + (72 - 80)^2 + \dots + (80 - 80)^2 \\
 &\quad + (80 - 81)^2 + (84 - 81)^2 + \dots + (82 - 81)^2 \\
 &\quad + (82 - 85)^2 + (80 - 85)^2 + \dots + (88 - 85)^2 \\
 &= \mathbf{186}
 \end{aligned}$$

ANOVA Table

A. Generalized ANOVA Table

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Squares	F-value
Between Samples (treatment)	$SSTr$	$c-1$	$SSTr/(c-1)$	$MSTr/MSE$
Within Samples (error)	SSE	$n-c$	$SSE/(n-c)$	
Total Variation	$SSTo$	$n-1$		

$$SSTo = SSTr + SSE \longrightarrow 251.7 = 65.7 + 186.0$$

B. ANOVA Table for Employee Training Programs

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Squares	F-value
Between Samples (treatment)	65.7	2	32.9	1.94
Within Samples (error)	186.0	11	16.9	
Total Variation	251.7	13		

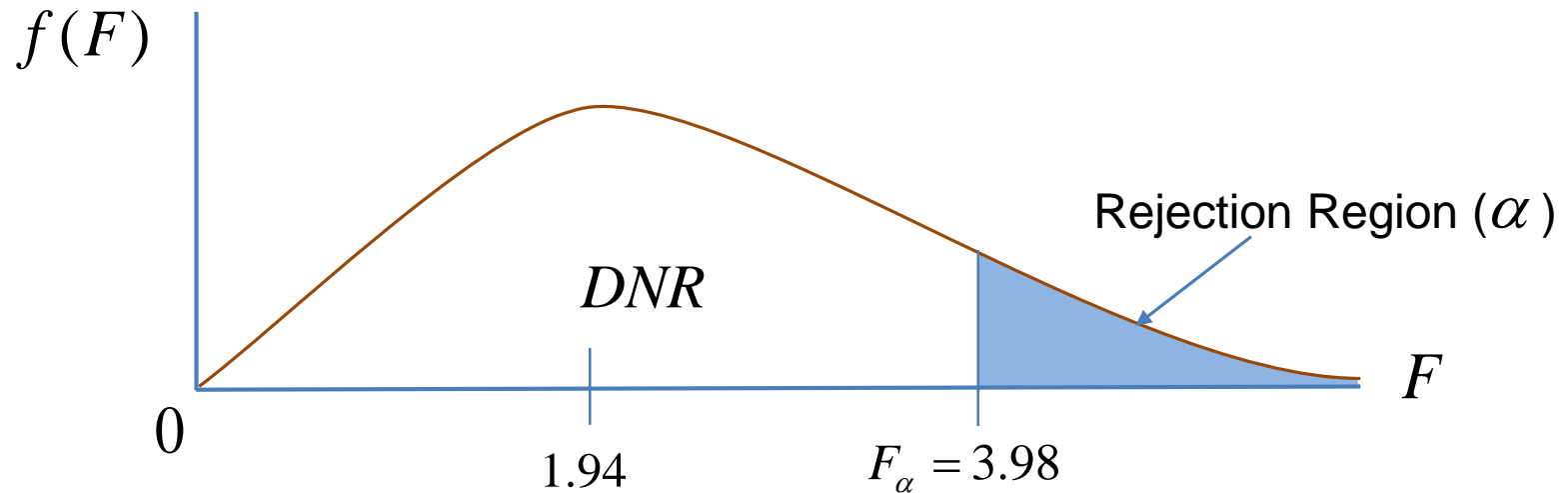
$$H_0 : \mu_1 = \mu_2 = \mu_3$$

H_a : Not all means are equal

Decision Rule : Do not reject if $F < 3.98$. Reject if $F > 3.98$

[F distribution table link](#)

Interpretation



Interpretation:

We do not reject the null because our calculated value 1.94 is less than the table value 3.98. Therefore, there are no significant differences in test scores for the three training methods.

One Way ANOVA- Balanced Problem *example*

Survey results on spending at health food stores were published in a medical trade journal. Four regions of the U.S. were surveyed on amount spent per month for six families in each region.

Are there significant differences in spending in the regions at the 5% level?

Region	Sample Mean
North	\$258.30
South	\$173.00
West	\$167.50
Midwest	\$138.30

$$H_o: \mu_1 = \mu_2 = \mu_3 = \mu_4$$

H_a : not all means are equal

In this example:

$$H_o: 258.30 = 173.00 = 167.50 = 138.30$$

H_a : not all means are equal

ANOVA Table for Regions

$$n = 24 \quad c = 4$$

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Squares	F-value
Between Samples (treatment)	48,023	3	16,008	9.02
Within Samples (error)	35,492	20	1,775	
Total Variation	83,515	23		

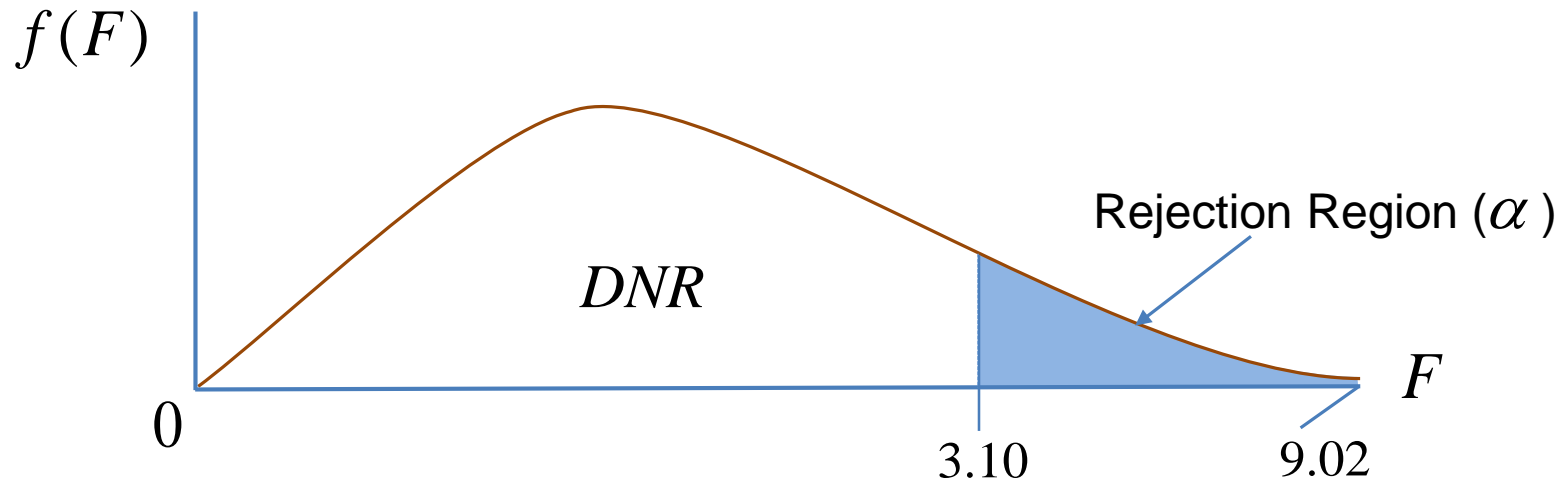
$$H_0 : \mu_1 = \mu_2 = \mu_3 = \mu_4$$

H_a : Not all means are equal

Decision Rule : Do not reject if $F < 3.10$. Reject if $F > 3.10$

Conclusion : Reject the null.

Interpretation



Interpretation:

We reject the null because our calculated value 9.02 is greater than 3.10. Therefore, not all regions spend the same amount on health foods.

Now our question is, which regions are significantly different from each other?

Tests for Differences between Individual Pairs

Tukey Test

- Can only be used for balanced problems
- Uses studentized range distribution

Least Significant Difference Test (LSD Test)

- Can be used for either balanced or unbalanced problems
- More conservative test than Tukey Test

Tukey Test

$$n = 24 \quad c = 4$$

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Squares	F-value
Between Samples (treatment)	48,023	3	16,008	9.02
Within Samples (error)	35,492	20	1,775	
Total Variation	83,515	23		

$$T = q_{\alpha, c, n-c} \sqrt{\frac{MSE}{r}}$$

If $\alpha = 0.05$, then $q_{0.05, 4, 20} = 3.96$,

$$T = 3.96 \sqrt{\frac{1.775}{6}} = 68.11$$

Tukey Test

$$T = 68.11$$

Region	Sample Mean
North	\$258.30
South	\$173.00
West	\$167.50
Midwest	\$138.30

$$|258.3 - 173.0| = 85.3 > 68.11$$

$$|258.3 - 167.5| = 90.8 > 68.11$$

$$|258.3 - 138.3| = 120.0 > 68.11$$

$$|173.0 - 167.5| = 5.5 < 68.11$$

$$|173.0 - 138.3| = 34.7 < 68.11$$

$$|167.5 - 138.3| = 29.2 < 68.11$$

Interpretation:

\$258.30 is significantly different from the other three numbers. Therefore, the North region spends significantly more money at health food stores than the other three regions. There are no significant differences between the other three regions.

LSD Test

$$n = 24 \quad c = 4$$

Source of Variation	Sum of Squares	Degrees of Freedom	Mean Squares	F-value
Between Samples (treatment)	48,023	3	16,008	9.02
Within Samples (error)	35,492	20	1,775	
Total Variation	83,515	23		

$$LSD = \sqrt{\frac{2(\text{MSE})F_{\alpha,1, n-c}}{r}}$$

$F_{0.05,1,20} = 4.35$. Then

$$LSD = \sqrt{\frac{2(1,775)(4.35)}{6}} = 50.73$$

LSD Test

$$LSD = 50.73$$

Region	Sample Mean
North	\$258.30
South	\$173.00
West	\$167.50
Midwest	\$138.30

$$|258.3 - 173.0| = 85.3 > 50.73$$

$$|258.3 - 167.5| = 90.8 > 50.73$$

$$|258.3 - 138.3| = 120.0 > 50.73$$

$$|173.0 - 167.5| = 5.5 < 50.73$$

$$|173.0 - 138.3| = 34.7 < 50.73$$

$$|167.5 - 138.3| = 29.2 < 50.73$$

Interpretation:

258.30, 173.00, 167.50, 138.30

173.00, 167.50, 138.30 are not significantly different from each other. However, 258.30 is significantly different from all the other three numbers.

Other Tests

Two way ANOVA- the randomized block design

- A second factor is considered in the model and is blocked to remove its effect. Such as operator skill, age, or gender

Factorial Analysis

- When two factors of interest are examined at the same time.
- Used to detect interactions between two factors

Latin Square Design

- Used to block out extraneous effect of two variables that might cause equivocal results
- Allows the researcher to gain more information with a smaller sample size
- Contains only one element per treatment per cell