

Week 4 (LAB2): Analysis of Variance

PSTAT122: Design and Analysis of Experiments

Delaney Thau

💡 Submission Instructions

- This LAB must be completed and submitted **individually**. Collaboration is allowed for discussion, but each student must submit their own work.
- Ensure that all R code are presented clearly and appropriately.
- All figures should be numbered, and axes must be labeled.
- Please use the provided `LAB 2 M.qmd` file to type your solutions and submit the completed LAB as a PDF file.

You can utilize RStudio for this purpose. For guidance, refer to the [Tutorial: Hello, Quarto](#).

- Submit your solution via **Gradescope as a single PDF file**.

🔥 Due Date

Due Date: Wednesday, January 28, 2026, 11:59 PM

❗ Overview

In this assignment, you will implement the ANOVA (Analysis of Variance) procedure from scratch, using basic operations in R. You'll compare your custom function's output with the built-in `aov` function to ensure correctness, and apply it to real datasets from class examples and homework.

1 Objectives:

1. Develop a custom ANOVA function that works with unbalanced datasets.
2. Validate your function by comparing its output to R's built-in `aov` function.
3. Apply the function to multiple datasets to demonstrate its accuracy and versatility.

! Implementing ANOVA from Scratch (15 Points)

Write a function to perform ANOVA from scratch on any arbitrary dataframe.

- Your function needs to
- take a `dataframe` as its input.
- expect that the dataframe will have two columns:
- the outcome variable in the first column
- the factor variable in the second column
- It should do the version of the calculations from the slides that do not rely on having a balanced dataset (i.e consider different values of n_i as in slide 51 of Chapter 3). Include a code comment that indicates exactly where in your calculations you are doing this.
- It should output the same exact values that a `summary` of the `aov` function does, in a matrix with the proper column and row names (elements of the matrix that have nothing in it should be filled in with `NA`).

ANSWERS TO Implementing ANOVA from Scratch:

```
1 library("dplyr")
2 anova_from_scratch <- function(dataframe) {
3   #Rename columns for consistency
4   colnames(dataframe) <- c("y", "group")
5
6   # Group by the factor and summarize: mean, variance, and count per group
7   summary_stats <- dataframe %>%
8     group_by(group) %>%
9     summarize(
10       y_bar_i = mean(y),
11       S_i = var(y),
```

```

12     n_i = n(),
13     .groups = 'drop'
14 )
15
16 # Grand mean
17 grand_mean <- mean(dataframe$y)
18
19 # Number of groups (a) and total sample size (N)
20 a <- nrow(summary_stats) # number of groups
21 N <- nrow(dataframe) #total sample size
22
23 # Sum of Squares Between (SSB)
24 # Taking account for unequal group sizes using n_i (slide 51 from Ch 3)
25 SSB <- sum(summary_stats$n_i * (summary_stats$y_bar_i - grand_mean)^2)
26
27 # Sum of Squares Within/Error (SSE)
28 SSE <- sum(summary_stats$S_i * (summary_stats$n_i - 1))
29
30 # Degrees of Freedom
31 df_between <- a - 1
32 df_within <- N - a
33
34 # Mean Squares
35 MSB <- SSB / df_between
36 MSE <- SSE / df_within
37
38 # F Statistic
39 F_value <- MSB / MSE
40
41 # p-value
42 p_value <- pf(F_value, df_between, df_within, lower.tail = FALSE)
43
44 # Construct ANOVA table matrix
45 anova_matrix <- matrix(NA, nrow = 2, ncol = 5)
46 rownames(anova_matrix) <- c("group", "Residuals")
47 colnames(anova_matrix) <- c("DF", "Sum Sq", "Mean Sq", "F value", "Pr(>F)")
48
49 # Fill in values
50 anova_matrix[1, ] <- c(df_between, SSB, MSB, F_value, p_value)
51 anova_matrix[2, 1:3] <- c(df_within, SSE, MSE)
52

```

```
53     return(anova_matrix)
54 }
```

! Testing on Class Example Datasets (10 Points)

Test this function on the navigation data from Chapter 3 (Soft Drink) and (Paper Strength). Also re-run the `aov` function with its `summary` here to show that you get the same values either way.

ANSWERS TO Testing on Class Example Datasets:

```
1 softdrink<- c(5.43, 5.71, 6.22, 6.01, 5.29, 6.24, 6.71, 5.98,
2           5.66, 6.60, 8.79, 9.20, 7.90, 8.15, 7.55)
3 y <- as.vector(t(softdrink))
4 n = rep(5,3)
5 group <- rep(1:3, n)
6 data = data.frame(y = y, group = factor(group))
7
8 anova_from_scratch(data)
```

	DF	Sum Sq	Mean Sq	F value	Pr(>F)
group	2	18.78305	9.3915267	35.77225	8.781737e-06
Residuals	12	3.15044	0.2625367	NA	NA

```
1 anova_real <- aov(y ~ group, data = data)
2 summary(anova_real)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
group	2	18.78	9.392	35.77	8.78e-06 ***
Residuals	12	3.15	0.263		

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

```
1 paperstrength <- c(14.7, 15.1, 16.4, 14.1, 14.5, 12.7, 9.88, 10.7,
2                     13, 13.3, 15.9, 12.8, 15, 18.4, 18.7, 19.3, 20,
3                     18, 17.7, 16.3, 15.6, 18.6, 17.6, 15.1)
4 y <- as.vector(t(paperstrength))
5 n = rep(6,4)
6 group <- rep(1:4, n)
7 data = data.frame(y = y, group = factor(group))
8
```

```

9  anova_from_scratch(data)

      DF   Sum Sq  Mean Sq F value    Pr(>F)
group      3 110.76672 36.922239 13.61481 4.561864e-05
Residuals 20  54.23833  2.711917       NA          NA

1  anova_real <- aov(y ~ group, data = data)
2  summary(anova_real)

      Df Sum Sq Mean Sq F value    Pr(>F)
group      3 110.77   36.92   13.62 4.56e-05 ***
Residuals 20  54.24    2.71
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

```

! Testing on Homework Dataset (5 Points)

Test your function on the data from Question 5 in the Homework 2. Run the `aov` function with its `summary` here as well to show that you get the same values either way.

ANSWERS TO Testing on Homework Dataset:

```
1 curlingmethods <- c(42.5, 40.8, 39.6, 41.2, 45.1, 46.3, 44, 45.7, 38.9,
2                               40.1, 41.5, 42, 36.2, 37.4, 35.8, 38)
3 y <- as.vector(t(curlingmethods))
4 n = rep(4,4)
5 group <- rep(1:4, n)
6 data = data.frame(y = y, group = factor(group))
7
8 anova_from_scratch(data)
```

	DF	Sum Sq	Mean Sq	F value	Pr(>F)
group	3	142.5069	47.502292	35.11643	3.202802e-06
Residuals	12	16.2325	1.352708	NA	NA

```
1 anova_real <- aov(y ~ group, data = data)
2 summary(anova_real)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
group	3	142.51	47.50	35.12	3.2e-06 ***
Residuals	12	16.23	1.35		

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