

How to Perform Welch's ANOVA in R: A Step-by-Step Guide

Authored by
stats writer

December 6, 2025

RECOMMENDED CITATION

stats writer (2025). *How to Perform Welch's ANOVA in R: A Step-by-Step Guide*.
PSYCHOLOGICAL SCALES. Retrieved from <https://scales.arabpsychology.com/?p=106085>

Understanding the Need for Welch's ANOVA

The standard Analysis of variance (ANOVA) is a fundamental statistical procedure used widely across experimental sciences to determine whether there are any statistically significant differences between the means of three or more independent groups. While incredibly efficient under ideal conditions, traditional ANOVA is classified as a parametric test, meaning its statistical validity rests upon several critical assumptions regarding the characteristics of the data. One of the most important of these is the assumption of the **homogeneity of variances**, also known as homoscedasticity. This assumption dictates that the variance, or spread, of the dependent variable must be approximately equal across all the groups being compared.

When real-world data fails to satisfy this condition--resulting in a state known as heteroscedasticity, where the group variances are significantly disparate--the results generated by the traditional ANOVA F-test can become severely compromised. Specifically, violating this assumption tends to inflate the Type I error rate, leading researchers to incorrectly conclude that a significant difference exists when in reality, the observed disparity is merely a function of unequal variability. Consequently, relying on standard ANOVA when this assumption is violated introduces a substantial risk of producing invalid statistical conclusions, necessitating a more robust methodology designed to handle variance heterogeneity.

This situation is precisely why **Welch's ANOVA** (or the Welch test) was developed. Welch's approach is a statistically robust alternative that does not assume equal variances. It achieves its robustness by adjusting the degrees of freedom utilized in the F-test calculation based on the observed differences in group variances, effectively weighting the data by the reciprocal of the sample variance. By implementing this corrective measure, Welch's ANOVA ensures that the comparisons between group means remain statistically accurate even when the assumption of homoscedasticity has been emphatically rejected. Utilizing this test is a core requirement for sound statistical practice whenever preliminary diagnostics reveal heterogeneity in variances across the comparison groups.

Step 1: Preparing the Data in R

The foundation of any successful statistical analysis is the meticulous preparation and structuring of the data. For this detailed demonstration of **Welch's ANOVA** in R, we are utilizing a dataset designed to investigate the effects of different studying techniques on student performance. Imagine a study where a professor seeks to determine if three distinct pedagogical methods--Technique A, Technique B, and Technique C--lead to significantly different outcomes in terms of final exam scores. To maintain experimental control, 10 students are randomly allocated to each of the three techniques, resulting in 30 total observations, and all students take an exam of identical difficulty.

In R, this data must be structured into a data frame containing at least two essential columns: the factor variable (`group`), which specifies the study technique used (A, B, or C), and the continuous dependent variable (`score`), which records the student's numerical exam performance. This arrangement is standard for analysis of variance procedures, requiring the grouping variable to be stored as a factor for accurate categorization by R's statistical functions. The following block of R code explicitly defines this data frame, ensuring the data is correctly mapped and ready for the forthcoming diagnostic and inferential tests.

The code below uses the `rep()` function to generate the grouping structure efficiently, ensuring that ten scores are associated with each of the three techniques. Following the creation of the `df` data frame, the `head(df)` command is executed to provide a quick visual verification of the dataset's structure. This check confirms that the grouping variable and the scores are correctly aligned, validating the integrity of the data structure before we proceed to test the necessary statistical assumptions.

```
#create data frame
```

```
df <- data.frame(group = rep(c('A','B', 'C'), each=10),  
score = c(64, 66, 68, 75, 78, 94, 98, 79, 71, 80,  
91, 92, 93, 85, 87, 84, 82, 88, 95, 96,  
79, 78, 88, 94, 92, 85, 83, 85, 82, 81))
```

```
#view first six rows of data frame
```

```
head(df)
```

```
group score
```

```
1 A 64
```

```
2 A 66
```

```
3 A 68
```

```
4 A 75
```

```
5 A 78
```

```
6 A 94
```

Step 2: Testing the Assumption of Equal Variances

The subsequent crucial step involves conducting a formal diagnostic test to rigorously evaluate whether the assumption of **homogeneity of variances** holds true for our dataset. If this assumption is upheld, we would proceed with standard ANOVA. If it is rejected, we transition immediately to Welch's alternative. Although other tests like Levene's test are often cited for their robustness against non-normality, **Bartlett's test** is readily available in the base R environment and serves as an effective initial check, provided the data within each group approximates a

normal distribution.

Bartlett's test is formulated to test the null hypothesis (H_0) that all population variances are equal ($\sigma^2_A = \sigma^2_B = \sigma^2_C$). The alternative hypothesis (H_A) states that at least two of these variances are significantly different. The statistical decision rests on the derived p-value. If this p-value falls below the predefined significance level, typically $\alpha = 0.05$, we possess compelling evidence to reject the null hypothesis, thereby confirming the existence of heteroscedasticity and mandating the use of a variance-adjusting test.

To perform this diagnostic check in R, we utilize the `bartlett.test()` function. This function requires the standard R formula syntax: `score ~ group`, indicating that the test should compare the variability of the `score` variable across the levels defined by the `group` factor. The following code execution demonstrates how to apply this function to our prepared data frame, yielding the crucial statistics needed to make an informed decision regarding the choice of the subsequent mean comparison test.

#perform Bartlett's test

```
bartlett.test(score ~ group, data = df)
```

Bartlett test of homogeneity of variances

data: score by group

Bartlett's K-squared = 8.1066, df = 2, p-value = 0.01737

Interpretation of Bartlett's Test Results

Upon examining the output generated by **Bartlett's test**, we find the critical statistics necessary for evaluating the core assumption of equal variances. The output provides a Bartlett's K-squared statistic of 8.1066, with 2 degrees of freedom, culminating in a critical p-value of **0.01737**. The decision rule requires comparing this p-value against our chosen significance level, $\alpha = 0.05$.

Since the calculated p-value (0.01737) is strictly less than the threshold of 0.05 , we are obligated to reject the null hypothesis that the population variances are equal across the three studying techniques. The data provides strong statistical evidence supporting the alternative hypothesis: that there is a significant difference in the variability of exam scores among the students assigned to different groups. This finding officially establishes that the assumption of **homogeneity of variances** has been violated.

This conclusion is critically important as it formally validates the decision to abandon the traditional ANOVA approach. Proceeding with standard ANOVA in the face of established heteroscedasticity

would jeopardize the reliability of the mean comparison. Therefore, the necessity of employing a method specifically designed to accommodate unequal variances is confirmed. Our next step must be the execution of the **Welch's ANOVA**, which provides the appropriate statistical corrections required for a valid comparison of the group means under these non-ideal conditions.

Step 3: Executing Welch's ANOVA in R

With the violation of the homogeneity assumption firmly established, we can now proceed to the main analysis: executing **Welch's ANOVA** to determine if the mean exam scores differ significantly across the studying techniques. In the R environment, this test is implemented using the versatile `oneway.test()` function, which is designed for one-way analysis of means. This function is preferred over the standard `aov()` function when variance equality cannot be assumed, as it includes the necessary parameters for the Welch correction.

To activate the Welch modification within R, we must explicitly include the argument `var.equal = FALSE` in the function call. This crucial parameter instructs R to perform the analysis without assuming equal variances, thereby calculating the adjusted F-statistic and the fractional degrees of freedom characteristic of the Welch test. The rest of the syntax remains consistent with the formula notation used previously: `score ~ group`, linking the dependent score variable to the independent grouping factor.

The execution of the following command initiates the robust mean comparison. The function calculates the F-statistic based on the separate variances of each group, providing a test that is far more reliable than the standard F-test when heterogeneity is present. This test will ultimately provide the overall p-value that determines the presence or absence of a significant effect of the studying technique on student scores.

#perform Welch's ANOVA

```
oneway.test(score ~ group, data = df, var.equal = FALSE)
```

One-way analysis of means (not assuming equal variances)

data: score and group

F = 5.3492, num df = 2.00, denom df = 16.83, p-value = 0.01591

Step 4: Interpreting the Welch's ANOVA Results

The output of the `oneway.test()` function provides the analytical results from the robust mean comparison. We observe the computed F-statistic ($F = 5.3492$), the numerator degrees of freedom ($\text{num df} = 2.00$), the calculated denominator degrees of freedom ($\text{denom df} = 16.83$), and the overall test p-value (0.01591). It is essential to note the denominator degrees of freedom: the

fractional value of 16.83 is a direct consequence of the Welch correction, which adjusts for the unequal variances found earlier, distinguishing this result from a standard ANOVA output.

The core of our inference lies in the overall p-value. This value addresses the primary research question by testing the null hypothesis that all population means across the three study techniques are equivalent ($\mu_A = \mu_B = \mu_C$). Comparing the p-value of **0.01591** against our standard significance level ($\alpha = 0.05$), we find that 0.01591 is decisively less than 0.05 . Based on this compelling statistical evidence, we formally reject the null hypothesis.

The ultimate conclusion from the **Welch's ANOVA** is that the three studying techniques do not result in statistically equivalent mean exam scores. A significant effect of the grouping factor exists on the outcome variable. However, this omnibus test only confirms that differences exist somewhere among the groups; it does not specify which particular pairs of techniques are significantly different from one another (e.g., A vs. B, A vs. C, or B vs. C). To isolate these specific group-wise differences, an additional analytical procedure is required, known as post-hoc testing.

Step 5: Conducting Post-Hoc Analysis

The confirmation of a significant overall effect through **Welch's ANOVA** naturally leads to the necessity of performing post-hoc tests. These follow-up tests are essential for conducting pair-wise comparisons between all combinations of the group means, allowing us to pinpoint precisely which study techniques yield significantly better or worse scores than others. This process is crucial for generating actionable insights from the research findings. However, conducting multiple comparisons significantly increases the probability of committing a Type I error (a false positive), which requires the use of specialized procedures to control the family-wise error rate.

Crucially, because our initial diagnostic confirmed a severe violation of the **homogeneity of variances**, standard post-hoc tests designed for equal variance conditions (such as Tukey's HSD) are unsuitable and would produce unreliable results. When the variances are unequal, the statistical community overwhelmingly recommends utilizing the **Games-Howell post-hoc test**. This test is specifically engineered to handle unequal sample sizes and unequal variances, making it the most appropriate companion to Welch's ANOVA. It operates by performing pair-wise t-tests using the modified degrees of freedom and variance estimates, ensuring robust control over the error rate.

To finalize the analysis and translate the statistical findings into meaningful conclusions for the professor, the Games-Howell test must be executed using an appropriate R package (as it is not included in base R). The resulting output will provide adjusted p-values for each pair-wise comparison, indicating exactly which techniques (A, B, or C) differ significantly from each other. Detailed tutorials and statistical references are highly recommended to ensure the correct package installation, execution, and accurate interpretation of the post-hoc results, confirming the precise

nature of the significant differences found in the study.

Selecting appropriate robust post-hoc tests (like Games-Howell) after Welch's ANOVA.

Implementing the necessary statistical packages in R for pair-wise comparisons under heteroscedasticity.

Interpreting the adjusted p-values from post-hoc comparisons to localize effects.

ARABPSYCHOLOGY.COM