

Simulation-Based Power-Analysis for Factorial ANOVA Designs

Daniël Lakens¹ & Aaron Caldwell²

¹ Eindhoven University of Technology, The Netherlands

² Department of Health, Human Performance and Recreation, University of Arkansas, USA

Author Note

All code used to create this manuscript is provided in an OSF repository at
<https://osf.io/xxxxx/>.

Correspondence concerning this article should be addressed to Daniël Lakens, ATLAS
9.402, 5600 MB, Eindhoven, The Netherlands. E-mail: D.Lakens@tue.nl

Abstract

Researchers need to design informative studies. When the goal of an experiment is to test a hypothesis based on a frequentist hypothesis test it is important to justify the sample size based on the statistical power of the study. However, researchers are faced with challenges when they try to calculate power for factorial ANOVA designs. First, current software solutions are limited as they do not allow power analyses for more complex designs involving multiple within factors. Second, they require partial eta-squared or Cohen's f as input, which are not the most intuitive way to specify predicted effects in an ANOVA, and do not generalize to different experimental designs. We have created R functions and an online Shiny app that performs simulations for ANOVA designs for up to three factors with an unlimited number of levels. Predicted effects are entered by specifying means, standard deviations, and correlations (for within factors). The simulation provides a-priori power analyses for all effects in the ANOVA, and all simple comparisons. No other software is currently available that allows researchers to so easily perform power analyses for a wide range of ANOVA designs. The app plots p -value distributions for all tests, and allows researchers to select a range of options to correct for multiple comparisons. This tutorial will teach researchers how to perform power analysis for ANOVA designs, and through simulations illustrate important factors that determine the statistical power of factorial ANOVA designs.

Keywords: power analysis, ANOVA, hypothesis test, sample size justification

Word count: 5000

Simulation-Based Power-Analysis for Factorial ANOVA Designs

Statistical power is the probability of observing a statistically significant result, given a specified effect size, and assuming there is a true effect. When a researcher aims to analyze a study based on an analysis of variance (ANOVA), the sample size of the study should be justified based on the statistical power of the test. A study has low power to detect effects the researcher is interested in has a high Type 2 error rate, and leads to a high probability of saying there is no effect, when there actually is a true effect to be found. Whereas power analyses for simple comparisons are relatively easy to perform, an a-priori power analysis for factorial ANOVA designs is a challenge. Available software solutions do not provide easy options for more complex designs (e.g., a 2x2x2 design, where the first factor is manipulated between participants, and the last two factors are manipulated within participants). Popular software solutions such as G*power require participants to enter their predictions as Cohen's f or partial eta squared, which are not the most intuitive ways to specify a hypothesized pattern of results, and which do not generalize to different experimental designs.

Here, we demonstrate how to perform power analyses for factorial ANOVA designs based on simulations. We provide R code and a Shiny app that can be used to calculate the statistical power based on a predicted pattern of means, standard deviations, and correlations (for within factors). Simulating studies, and calculating their p -values and effect sizes, are a useful way to gain a better understanding of the factors that determine the statistical power of hypothesis tests. After providing an introduction to statistical power in general, and for the F -test specifically, we will demonstrate how the power of factorial ANOVA designs depends on the pattern of means across conditions, the number of factors and levels, and the sample size. We will also illustrate the importance of controlling Type 1 error rates in exploratory ANOVA's, and the importance to design studies that have high power for main effects and interactions in an ANOVA, but also for follow-up tests for simple effects.

Factors That Determine Power in ANOVA Designs

You perform a study in which participants interact with an artificial voice assistant who sounds either cheerful or sad. You measure how much 70 participants in each condition enjoy interacting with the voice assistant on a line marking scale (coded continuously from -5 to 5) and observe a mean of 0 in the sad condition, and a means of 1 in the cheerful condition, with an estimated standard deviation of 2. After submitting your manuscript for publications, reviewers ask you to add a study with a neutral control condition to examine whether cheerful voices increase, or sad voices decrease enjoyment (or both). Depending on what the mean in the neutral condition is, which sample size would you need to have a high powered study for the expected pattern or means? A collaborator suggests to switch from a between design to a within design, to more efficiently collect the data. Which consequence will this switching to a within-subject design have on the sample size planning? Because the effect size in the first study could be considered “medium” based on the benchmarks by Cohen (1988), does it make sense to plan for a “medium” effect size in either the between of within ANOVA design? And if you justify the sample size based on the ANOVA, will the study also have sufficient statistical power for the simple effects (or vice versa)?

Performing simulation studies is an excellent way to develop intuitions about these questions. The power in ANOVA designs depends on the pattern of means, the number of groups, the standard deviation, the correlation between dependent measures, the alpha level, and the sample size. After these factors have been specified, a simulation study can be performed to run the planned statistical test on generated data many times, and summarize the results. Where analytic power solutions exist in software such as G*Power for ANOVA designs with up to one within-subject factor, simulation studies are more flexible (and can for example be used to see what happens in a 3x3x3 within-subject design).

Statistical Power When Comparing Differences Among Group Means

Let's consider the initial study we described above, where two group means are compared. We can test the difference between two means with a t -test of a one-way single factor ANOVA, and the two tests are mathematically equivalent. Figure 1 and Figure 2 visualize the distribution of Cohen's d and partial eta-squared that should be observed when there is no effect (grey curves) and when the observed difference between means equals the true effect. In both figures the light grey areas under the curve mark the observed results that would lead to a Type 1 error (observing a statistically significant result if the null-hypothesis is true) and the dark grey areas under the curve marks the observed effect sizes that would lead to a Type 2 error. An observed effect is statistically significant if the observed effect size is larger than the critical value. Critical values are often expressed as t -values or F -values, but can be expressed as effect sizes (Cohen's d and partial eta-squared), and any observed effect size larger than the critical effect size will be statistically significant. The goal of power analysis is to choose a sample size so that the probability of observing a statistically significant effect for a specified effect size reaches a desired probability. In both Figures, 83.5% of the expected effect sizes, if the true effect size is $d = 0.5$ or $\eta_p^2 = 0.0588$ and 70 participants in each condition are collected, will be more extreme than the critical effect size (which is $d = 0.334$ or $\eta_p^2 = 0.028$).

Two relationships between the t -test and the F -test when comparing two means are worth pointing out. First, $F = t^2$, or the F -value equals the t -value, squared. Whereas we typically think of a t -test as the difference between means, the relationship between the t -test and F -test reveals that we can also use the group means to calculate the variance of the difference between the two means $(m_1 - m_2)^2$. The F -test is used to compute the ratio of the between group variance and the within group variance, and the between group variance can be used to compare more than two means. An F -value of 1 means the two variances are equal (as is expected when the null-hypothesis is true) and Cohen's f and η_p^2

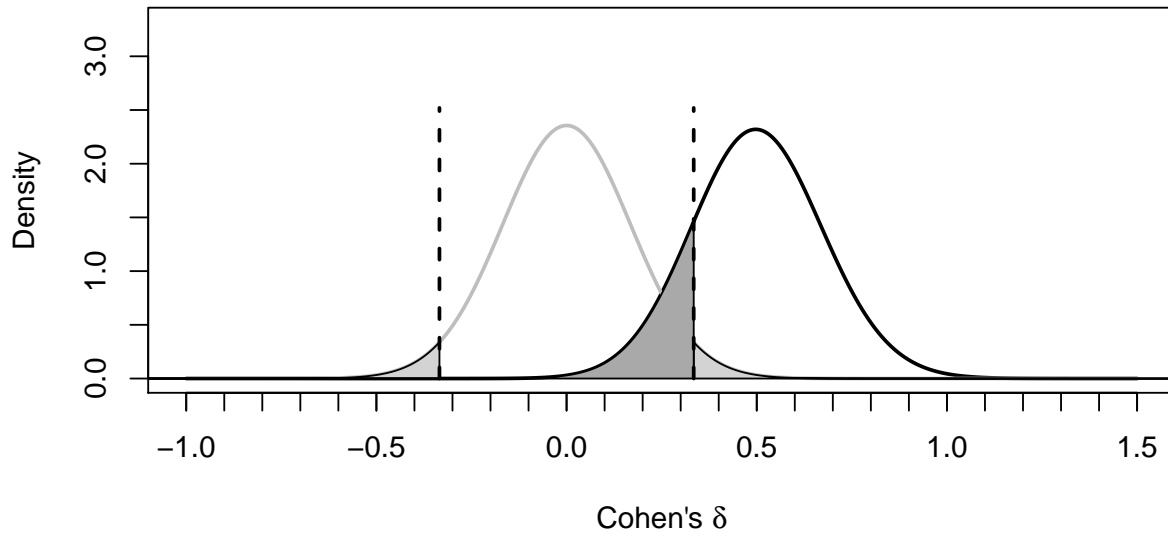


Figure 1. Distribution of Cohen's d under the null-hypothesis (grey curve) and alternative hypothesis assuming $d = 0.5$ (black curve). Note that

would be zero. Second, for two groups the effect size for an ANOVA, Cohen's f , is half the size of the effect size for standardized mean differences, Cohen's d , or $f = \frac{1}{2}d$. Cohen's d is calculated by dividing the difference between means by the standard deviation, or

$$d = \frac{m_1 - m_2}{\sigma}. \quad (1)$$

If we have two groups with means of 1 and 2, and the standard deviation is 2, Cohen's d is $(2-1)/2$, or 0.5. Cohen's f is the standard deviation of the population means divided by the population standard deviation (Cohen, 1988), or:

$$f = \frac{\sigma_m}{\sigma} \quad (2)$$

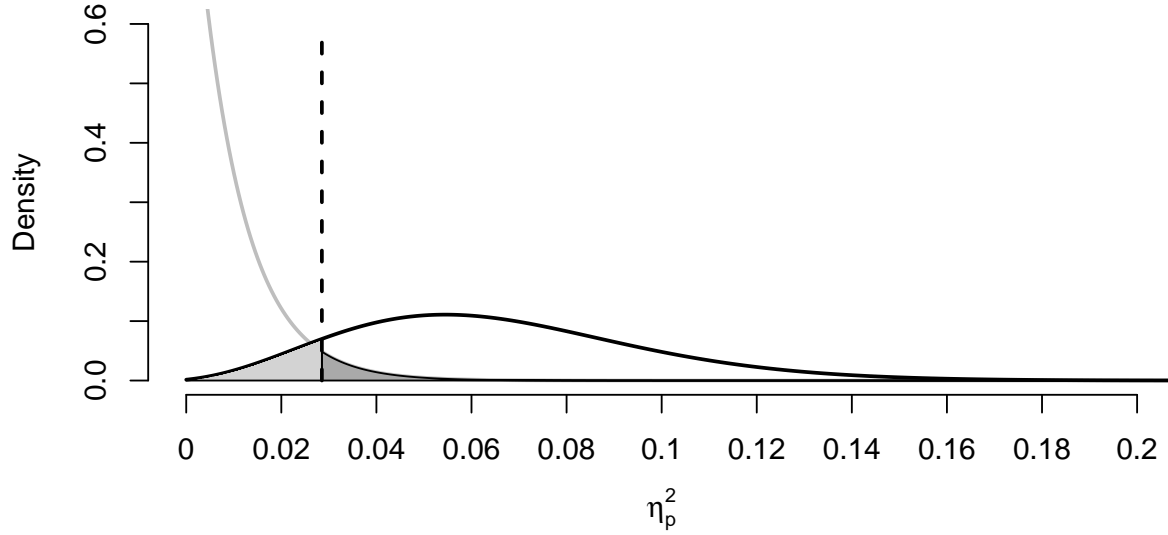


Figure 2. Distribution of eta-squared under the null-hypothesis (grey curve) and alternative hypothesis assuming partial eta-squared = 0.0588 (black curve).

111 where for equal sample sizes,

$$\sigma_m = \sqrt{\frac{\sum_{i=1}^k (m_i - m)^2}{k}}. \quad (3)$$

112 Because Cohen's f is an essential part of power power analyses for factorial ANOVA
 113 designs, it is worth illustrating how it is calculated in an example. If we again take two
 114 means of 1 and 2, and a standard deviation of 2, the grand mean is 1.5. We subtract each
 115 condition mean from the grand mean, take the square, calculate the sum of squares, divide it
 116 by two, and take the square root. $\sigma_m = \sqrt{\frac{(1-1.5)^2 + (2-1.5)^2}{2}} = \sqrt{\frac{0.25+0.25}{2}} = 0.5$, and
 117 $f = \frac{0.5}{2} = 0.25$. We see Cohen's f is half as large as Cohen's d . Power analyses for ANOVA
 118 are based on Cohen's f , but popular power anaysis software such as G*Power (Faul,
 119 Erdfelder, Lang, & Buchner, 2007) also allows researchers to specify the effect size as partial
 120 eta-squared (η_p^2). Partial eta-squared can be converted into Cohen's f :

$$f = \sqrt{\frac{\eta_p^2}{1 - \eta_p^2}} \quad (4)$$

and Cohen's f can be converted into partial eta-squared:

$$\eta_p^2 = \frac{f^2}{f^2 + 1} \quad (5)$$

In the example above, $\eta_p^2 = 0.25^2 / (0.25^2 + 1) = 0.0588$. To calculate the statistical power assuming a specific true effect size we need the noncentrality parameter of the distribution. In both Figure 1 and Figure 2 we see examples of the non-central t distribution and non-central F distribution, or the shape of the expected test statistics when there is a true effect. Power calculations rely on the noncentrality parameter (often referred to as λ). Based on λ (which specifies the shape of the expected distribution) and the critical test statistic (which specifies the part of the distribution that is more extreme than the test statistic needed for a statistically significant test result) we can calculate how often, in the long run, we can expect test results that will be statistically significant.

Simulating Statistical Power for Different Factorial Designs

We have created R code and a Shiny app that simulate factorial ANOVA designs. At the core, the code generates data for each condition in the design, performs the test, and calculates the percentage of significant results and average effect size of all main effects and interactions, as well as all simple effects. It requires specifying the design of the study in a string of numbers (for the levels in each factor) and letters (a "b" for between factors or "w" for within factors). Our initial study above is a "2b" or two level between subject design. For ease of interpreting the simulation results, the factors and levels can be named (for our example, Condition, cheerful, sad). The planned sample size should be specified, which is 70

140 participants in each between subjects condition. The means for each condition should be
 141 specified (0 and 1), as well as the standard deviation (2). For within designs the correlation
 142 between variables should be specified. With these variables specified, the design is set up for
 143 the simulation. For a visual confirmation of the input, a figure is created that displays the
 144 means.

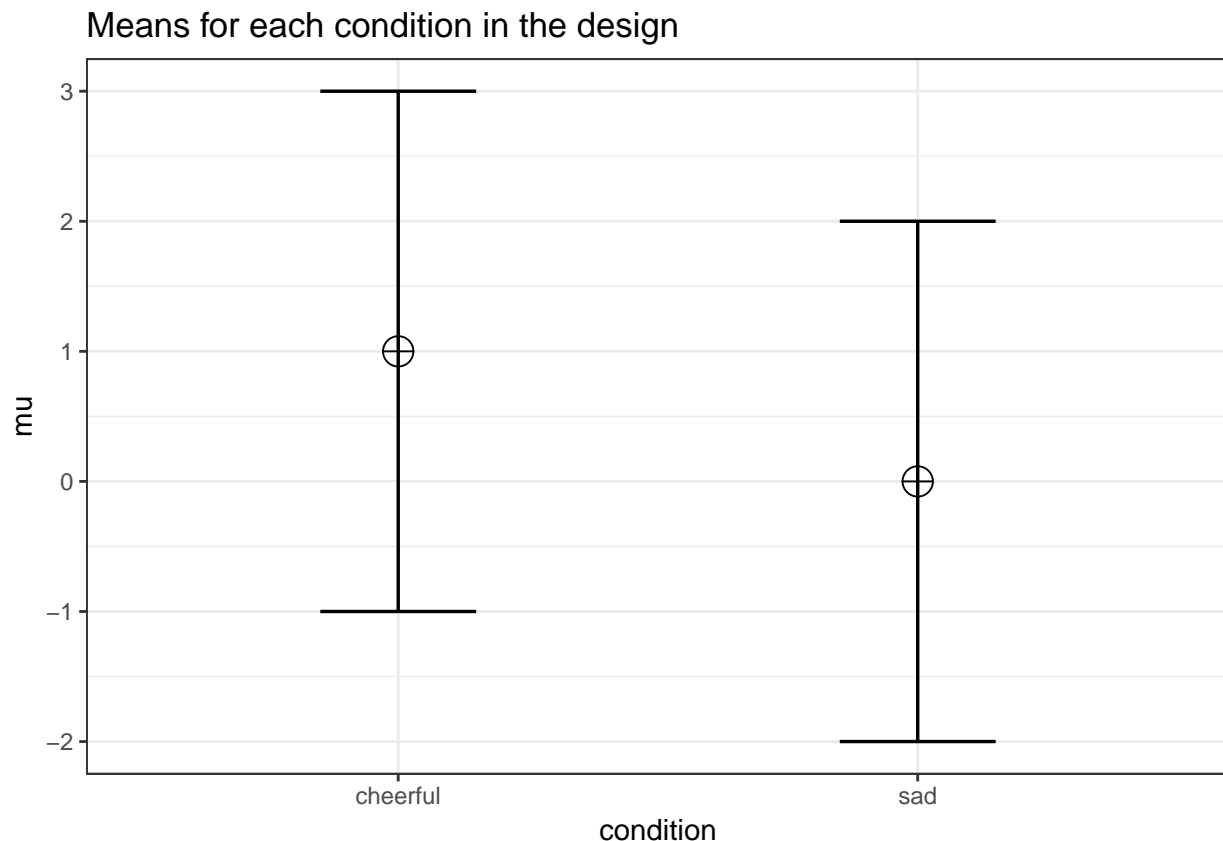


Figure 3. (#fig:sim_1_design)Means for each condition in the design.

145 The results of a simulation study will vary each time the simulation is performed (but
 146 can be made reproducible by specifying a “seed” number). Results become more stable, the
 147 more simulations are performed, but the simulations also take longer. The functions that
 148 performs the simulation requires specifying the number of simulations, the alpha level for the
 149 tests, and any adjustment for multiple comparisons for the ANOVA. If 100.000 simulations
 150 are performed in R with a seed set to 2019 (these settings will be used for all simulation

results in this manuscript), we see the statistical power (based on the percentage of $p < \alpha$ results) is 81.50 and the average effect size is 0.06. The simulation also provides the results for the simple effects, based on t -tests comparing each group. Since there are only two groups in this example, the results for the statistical power are identical, but the effect size is -0.50, which is the effect size in Cohen's d . Now the basic idea behind the simulation is clear, we can start exploring how changing the experimental design influences our power. We will first examine what happens if we add a third condition to the design. Let's assume we expect the neutral voice condition to fall either between the cheerful and sad conditions, or perhaps to be equal to the cheerful condition (based on the idea that the sad voice leads to less enjoyment, but the cheerful voice does not lead to more enjoyment). The design now has 3 between subject conditions, and we can explore what happens if we would collect 70 participants in each condition.

If we assume the mean falls exactly between between the cheerful and sad conditions the simulations show the statistical power for our design is reduced to 81.50, and the effect size is -0.50. If we assume the mean the mean is equal to the cheerful condition, the power increases. Compared to the two group design, three things have changed. First, an additional group was added to the design, which increases the numerator degrees of freedom, which makes the non-central F -distribution more similar to the central F -distribution, and reduces the statistical power. Second, the total sample size is 50% larger, which increases the statistical power. Third, the effect size has decreased, which reduces the statistical power. The exact effect of these three changes on the statistical power is difficult to predict. The most important conclusion based on these simulations is that when the design is changed, one can not assume the effect size remains unchanged.

What happens if we would perform the second study as a within-participants design? Instead of collecting three groups of participants, we only collect one group, and let this group evaluate the cheerful, neutral, and sad voice assistants. The sample size needed in a

177 within-design (NW), relative to the sample needed in between-design (NB), is (from Maxwell
178 & Delaney, 2004, p. 562, formula 47):

$$N_W = \frac{N_B(1 - \rho)}{a} \quad (6)$$

179 Here a is the number of within-subject levels, ρ is the correlation between the
180 measurements, and the formula assumes normal distributions and compound symmetry, and
181 ignores the difference in degrees of freedom between the two types of tests, so is a rough (but
182 useful) approximation.

183 We need to be able to estimate the correlation between the dependent variables we can
184 expect. Ideally, this estimate will be based on data from previous experiments. If we plan to
185 collect 70 participants, and assume the correlation between all dependent variables is 0.5, the
186 power for the repeated-measures ANOVA is . Note that simulation studies allow for great
187 flexibility, and the Shiny app also allows researchers to enter a correlation matrix specifying
188 the expected correlation between each individual pair of measurements, instead of assuming
189 the correlations between dependent variables are identical for all pairs.

190 `str(power_result)`

References

- Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed.). Hillsdale, N.J: L. Erlbaum Associates.
- Faul, F., Erdfelder, E., Lang, A.-G., & Buchner, A. (2007). GPower 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behavior Research Methods*, 39(2), 175–191. doi:10.3758/BF03193146