

## **The benefits of reporting critical effect size values**

Ambra Perugini<sup>1</sup>, Filippo Gambarota<sup>1</sup>, Enrico Toffalini<sup>2</sup>, Daniël Lakens<sup>3</sup>, Massimiliano Pastore<sup>1</sup>, Livio Finos<sup>4</sup>, Psicostat<sup>1</sup>, and & Gianmarco Altoè<sup>1</sup>

<sup>1</sup> Department of Developmental and Social Psychology

University of Padova

Italy

<sup>2</sup> Department of General Psychology

University of Padova

Italy

<sup>3</sup> Eindhoven University of Technology

Netherlands

<sup>4</sup> Department of Statistics

University of Padova

Italy

### **Author Note**

The authors made the following contributions. Ambra Perugini: Conceptualization, Writing - Original Draft Preparation, Writing - Review & Editing; Filippo Gambarota: Writing - Review & Editing, Supervision; Enrico Toffalini: Writing - Review & Editing, Supervision; Daniël Lakens: Writing - Review & Editing, Supervision; Massimiliano Pastore: Writing - Review & Editing, Supervision; Livio Finos: Writing - Review & Editing, Supervision; Psicostat: Writing - Review & Editing, Supervision; Gianmarco Altoè: Writing - Review & Editing, Supervision.

Correspondence concerning this article should be addressed to Gianmarco Altoè, Via Venezia 8, 35131, Padova, Italy. E-mail: gianmarco.altoe@unipd.it

## **Abstract**

Scrivere abstract

*Keywords:* keywords

Word count: X

### The benefits of reporting critical effect size values

In the present paper, we suggest that researchers systematically report the critical effect size values when they make inference using statistical significance. The “critical effect size values” refers to the smallest statistically significant effect sizes, which depends on the test performed and the characteristics of the sample (Lakens, 2022). To cut to the chase, let us consider a study that tests (using for example a two-tailed test) a bivariate correlation with  $n = 20$  participants and uses  $p < 0.05$  as the threshold for statistical significance: in this case, the smallest significant effects (critical values) are Pearson’s  $r = -0.44$  and  $r = 0.44$ . As we will argue, the critical values can serve as a useful complement to aid in the interpretation of study results, communicate their relevance, help assessing their replicability, particularly when statistical significance is relied upon for inference. This is especially important in cases where the sample size was not (or could not) be predetermined based on a theory-driven target effect size or when statistical power may be difficult to determine or deviates from an optimal level.

Despite longstanding criticisms, the Null Hypothesis Significance Testing (NHST) remains a prominent approach for statistical inference in science (Cohen, 1994; Gigerenzer et al., 2004; Krueger, 2001). Over the past decades, numerous proposals have emerged to enhance inference based on this approach. These include conducting power analyses to determine sample size ( $n$ ) a priori based on theory-informed target effects, and reporting of effect sizes and their confidence intervals, accompanied by pertinent related comments, rather than drawing dichotomous conclusions solely based on statistical significance (Transue, 2019). Even for the most skeptical, there are suggestions that NHST may possess merits of its own. For instance, Wilson et al. (2020) propose that it could serve as a useful preliminary filter utilized by “original science” for screening potentially interesting effects, that should subsequently be validated by “replication science”.

In psychology, power is often severely limited for medium to small effect sizes (Szucs

& Ioannidis, 2017), which are generally expected to represent true replicable findings. In fact, resource constraints have been authoritatively listed among acceptable justifications for sample size (Lakens, 2022), and this scenario is arguably widespread in psychology. However, limited power makes it challenging to distinguish signal from noise, and on average it results in an overestimation of effect sizes when relying on statistical significance for inference (Altoè et al., 2020). On the other hand, very large samples are occasionally available in psychological science, albeit less frequently. However, as some have argued, excessively large sample sizes risk creating an “everything-is-significant” scenario, where researchers report and discuss practically negligible effects solely because they are statistically significant, regardless of theoretical relevance or whether they just reflect minor procedural artifacts (Wilson et al., 2020). Thus, while statistical significance can serve as a potentially useful filter, relying solely on it for inference may introduce risks in interpretation. In such cases, we contend that reporting critical values can help put into perspective and convey the relevance of the results.

There are two clear benefits of reporting the critical effect size for a corresponding test. First, when sample sizes are smaller, the critical values inform readers about whether the effect sizes that could lead to a rejection of the null hypothesis are in line with realistic expectations. If the sample size is small, and only very large effects would yield a statistically significant result, and the underlying mechanism that is examined is unlikely to lead to such large effect sizes. In these cases, researchers will realized they are not able to collect sufficient data to perform a meaningful test. An a-priori power analysis would typically lead to a similar conclusion, but reporting the critical effect sizes focuses the attention more strongly on which effect sizes are reasonable to expect. Second, in large studies the critical effect sizes will make clear that trivially small effect sizes will be statistically significant. This will focus the attention of researchers on the difference between statistical significance and practical significance. This is especially important in correlational studies on large samples in psychology, where small uncontrolled sources of

statistical noise may simultaneously affect and create some shared variance across otherwise unrelated variables, a phenomenon controversially labelled as the ‘crud factor’ (Orben & Lakens, 2020).

Let us consider the following two scenarios as examples. First, imagine researchers conducting a study involving a between-group comparison. Due to severe resource constraints, they are only able to collect a limited sample size of  $n = 30$  ( $n = 15$  per group). No statistically significant effect is detected with  $p < 0.05$ . Given their prior expectation that the effect of interest may not be large, they acknowledge that their study was likely underpowered, though uncertain to what extent. Subsequently, they compute the critical values, revealing Cohen’s  $d = -0.75$  and  $d = 0.75$ . This indicates that any estimated  $|\text{Cohen’s } d| < 0.75$  will certainly fail to reach statistical significance. By reporting these critical values, the researchers transparently convey that estimated effects up to a medium-to-large magnitude will consistently fall short of significance. In essence, this scenario exemplifies a version of the “winner’s curse” (Hedges, 1984), made explicit. In this context, the winner’s curse indicates the tendency for a statistically significant initial finding (i.e., a “winner”) to be associated with an overestimated effect size which will probably be deflated in subsequent replications (i.e., the “curse”), especially if such initial finding was obtained on a small sample. Knowing from the critical values that a specific result must necessarily be associated with an estimated large effect size or be non-significant, readers are, in principle, prevented from both over-interpreting the magnitude of an effect size and from equating lack of statistical significance with lack of an effect of interest.

In a second scenario, imagine researchers who gain access to a very large archival dataset ( $n = 5,000$ ) and decide to explore some bivariate correlations. With such a large sample size, researchers are not constrained by sample size limitations and could instead just focus on effect sizes. However, they opt to legitimately use statistical significance of a

two-tailed test as their primary criterion for whether to comment on effects in their further discussions. In this context, the critical effect size values for significance are determined to be correlations of  $r = \pm 0.03$ . This implies that practically any effect differing from zero, regardless of its small magnitude, may easily attain statistical significance. It is worth noting that such small effects may potentially reflect just minor artifacts (Wilson et al., 2020), such as subtle experimenter effects or slight non-independence among observations, even in meticulously designed studies. While readers may be tempted to interpret any reported effects simply because they attain statistical significance, preliminarily signaling the critical value level can serve as a clear warning that the significance filter might be somehow permissive, thus urging interpretive caution.

A similar scenario may arise in meta-analysis. Despite potential loss of precision due to substantial heterogeneity across effect sizes in different studies, meta-analyses typically synthesize a large amount of evidence, leading even very small average effect sizes to reach significance. While the focus of meta-analysis is generally on estimating effect sizes with uncertainty, statistical significance is routinely reported and interpreted. Thus, in meta-analysis as well, reporting the critical values can caution readers that the mere attainment of statistical significance by the average effect does not automatically ensure scientific or real-life relevance.

To provide more tangible illustrations of the practical application of critical values, let us examine two real-world instances drawn from published research. The first one features a small sample, and has been targeted in the series of replications by the Open Science Collaboration (Collaboration, 2015), while the second one presents with a very large sample and has already attracted attention in the scientific community. The first case pertains to “Study 5” conducted by McCrea (2008). This experiment compared two groups based on their performance percentage in a math test following a preliminary practice session with feedback of failure. One group was exposed to failure-excusing,

self-handicapping thoughts, while both groups underwent the math test afterward. With a modest sample size of 28 participants split into two groups (13 and 15 participants, respectively), a directional one-tailed t-test was conducted on the percentage of correct answers to compare the groups. The resulting effect was statistically significant,  $t(26) = 1.87, p(\text{one-tailed}) < .05$ . Based on the reported data, Cohen's  $d = 0.736$ , indicating a considerable effect size. However, the critical value for Cohen's  $d$  in this context was 0.645, suggesting that, given this sample size, statistical significance is achieved only for estimated effects of substantial magnitude. Subsequent replication efforts (Collaboration, 2015) reaffirmed the significance of the effect, and used a somewhat larger sample ( $n = 61$ ),  $t(59) = 2.325$ . The estimated effect size was somehow lower than that originally reported, however,  $d = 0.605$  (here, critical  $d = 0.428$ ). The second scenario involves the study by Kramer et al. (2014), which explored the impact of emotional content on Facebook users' experiences. With a notably large sample size of  $n = 689,003$ , the researchers observed a significant effect worsened emotional states when positive posts were reduced,  $t(310,044) = -5.63, p < 0.001$  (non-directional test). Based on this and other results of a similar, the authors conclude that emotions expressed by other users on Facebook influence our own emotions. Despite the statistical significance, however, the critical values for Cohen's  $d$  in this case are  $-0.006$  and  $+0.006$ , emphasizing the importance of considering effect size beyond mere significance. This prompts reflection on whether an effect of such magnitude on such a large sample (actual  $d = 0.02$ ) is truly meaningful in individuals' lives.

### How to Compute Critical Values

In this section, we provide guidance and formulas for computing critical values in general, with examples for frequently encountered effect sizes including Standardized Mean Differences (Cohen's  $d$ ), correlations (Pearson's  $r$ ), and raw and standardized coefficients in linear models. These formulas have been incorporated into R functions of the package "criticalESvalues" accessible at: [xxxxx], and are elaborated upon in the subsequent section.



**t-test**

For the t-test we considered the two-sample, paired and one-sample tests. As a general approach the  $t$  statistics is computed as reported in Equation (1).

$$t = \frac{b}{SE_b} \quad (1)$$

Where  $b$  is the unstandardized effect size that depends on the type of test. For example, in the two sample t-test or the single mean for a one sample t-test. The denominator is the standard error of the numerator.

Similarly, Equation (2) formalize a general form of the effect size.

$$d = \frac{b}{s} \quad (2)$$

Where  $b$  is still the unstandardized effect size and  $s$  is the standardization term. For example, in the two-sample case,  $s$  is the pooled standard deviation between the two samples or the standard deviation of the differences for paired samples.

***two-sample t-test***

For the two sample t-test the critical value ( $b_c$ ) is calculated using Equation (3). Where  $t_c$  is the critical  $t$  value calculated using a certain  $\alpha$  value (e.g., 0.05) and  $n_1 + n_2 - 2$  degrees of freedom. Then the  $b_c$  is simply divided by  $s$  obtaining the standardized critical value. For the two-sample case ( $s = s_p$  assuming equal variances Equation (4) reported the pooled standard deviation.

$$b_c = t_c \times SE_b \quad (3)$$

$$s_p = \sqrt{\frac{s_1^2(n_1 - 1) + s_2^2(n_2 - 1)}{n_1 + n_2 - 2}} \quad (4)$$

When relaxing the assumption of equal variances (i.e., Welch's test)  $s$  is simply the square root of the average between the two variances.

### *one sample t-test*

For the one sample case, the equations are the same. The only difference is that  $b$  is the sample mean and  $s$  is the sample standard deviation.

### *paired sample t-test*

For the paired sample case the situation is less straightforward. The reason is that the  $t$  statistic is computed using  $b$  being the average of the paired differences between the two (paired) samples and  $s$  being the standard deviation/error of the differences. While for hypothesis testing and power analysis it is common to calculate the  $d$  dividing by the standard deviation of the differences (Lakens, 2013), this effect size cannot be directly compared to the effect size used in the independent samples t-test, which divides by the pooled standard deviation (Morris & DeShon, 2002). Practically, Equation (5) depict the relationship between the pooled standard deviation and the standard deviation of the differences.

$$\begin{aligned} s_p &= \frac{s_D}{\sqrt{2(1 - \rho)}} \\ s_D &= s_p \sqrt{2(1 - \rho)} \end{aligned} \quad (5)$$

Where  $\rho$  is the correlation between the two paired samples. Even if the hypothesis testing is computed using  $s_D$ , we reported the critical value both using the pooled standard deviation and the standard deviation of the differences. Beyond this difference, the equations are still the same as the previous example.

***Hedges's correction***

The effect size calculated as in the previous step is known to be inflated thus we calculated also the corrected version usually known as Hedges's  $g$  (Hedges, 1981).

**Correlation Test**

For the Pearson's correlation test, the critical value is computed using Equation (6) where  $t_c$  is still the critical value for the  $t$  distribution and  $n$  is the sample size.

$$r_c = \frac{t_c}{\sqrt{n - 2 + t_c^2}} \quad (6)$$

Another approach for hypothesis testing of the Pearson's correlation coefficient is using the Fisher's  $z$  transformation. Equation (7) shows how to calculate the critical  $F(r_c)$  value where  $z_c$  is the critical value of the standard normal distribution with a chosen  $\alpha$  level. Then the  $F(r_c)$  can be transformed back into a correlation coefficient using Equation (8).

$$F(r_c) = \frac{z_c}{\sqrt{n - 3}} \quad (7)$$

$$r_c = \frac{\exp(2z_c) - 1}{\exp(2z_c) + 1} = \tanh(z_c) \quad (8)$$

**Linear Regression**

Hypothesis testing on single coefficients in linear regression (using the `lm` function in R) is performed using the Wald test. This test is basically a  $t$ -test where the  $t$  value is calculated dividing the regression parameter  $\beta_0, \dots, \beta_j$  by the standard error. Then a  $t$  distribution using  $n - p - 1$  degrees of freedom where  $n$  is the number of observations and  $p$  is the number of coefficients beyond the intercept. The  $t$  value is calculated as reported in Equation (9) and the critical value as reported in Equation (10).

$$t = \frac{\beta_j}{SE_{\beta_j}} \quad (9)$$

$$\beta_{jc} = t_c SE_{\beta_j} \quad (10)$$

In this case  $\beta_{jc}$  is the unstandardized critical regression coefficient (as  $b$  in the t-test section). We implemented also the possibility to use standardized regression coefficients. We can identify two types of standardization namely full and partial. The full standardization involves both the response variable and the predictors. Regression coefficients are interpreted as the increase in standard deviations of the response variable for an increase of one standard deviation in the predictor, keeping fixed all other predictors. The partial standardization involves only the predictors and regression coefficients are interpreted as the increase in the (raw) response variable for a increase in one standard deviation of the predictor. By default, standardizing means dividing a quantity by the standard deviation and eventually centering (i.e., subtracting the mean). Gelman (Gelman, 2008) suggested to standardize using two standard deviations because when numeric and binary variables are included the coefficients are not on the same metric. Equation (10) can be used but to calculate the critical value for the standardized regression coefficients however the result need to be interpreted according to the chosen approach.

### ***Meta-analysis***

Meta-analysis allows to pool information from multiple studies related to specific research question. The main advantage of meta-analysis is pooling multiple studies to obtain a more precise and powerful estimation of the effect. From a statistical point of view, a meta-analysis can be considered as a weighted linear regression with heterogeneous variances. Similarly to standard linear regression, hypothesis testing is performed using Wald  $t$  or  $z$  tests. The calculation of critical values is the same as reported in Equations (9) and (10). The only difference is that the standard error of the meta-analysis parameters

included both the within-studies variances and eventually the between-studies heterogeneity.

### ***Other models***

Despite we discussed only linear models, the same approach could be theoretically applied to other types of models such as generalized linear models. In fact, we simply need to multiply the critical value of the chosen distribution (e.g., Student  $t$  or Standard Normal) by the standard error of the regression coefficient.

## **Examples in R**

In this section, we introduce a user-friendly implementation of the aforementioned mathematical computations as functions of the package “criticalESvalues” in R. The complete package can be accessed at: [xxxx]. Here, we demonstrate its application through two examples: one example of a t-test on real data and a computation of the critical value for a correlation from sample size. In the Supplementary online materials additional examples can be found for correlation, t-test, paired t-test, linear models and regression coefficients, both from data and from sample size.

First the package should be downloaded and opened with the library function:

```
# devtools::install_github("filippogambarota/criticalvalue")  
library(criticalvalue)
```

For our examples on real data we used from the package ‘psych’ the dataset “holzinger.swineford” which has a series of demographics and scores of different subtests measuring intelligence on 301 subjects. Once the package is retrieved with ‘library’, the dataset can be opened using ‘data(“name of the dataset”)’. For simplicity we decided to rename it with a shorter name.

```
library(psych)
library(psychTools)
data("holzinger.swineford")
Holz <- holzinger.swineford
```

We want to know the critical value for a t-test comparing boys and girls on a cognitive variable of visual perception. In this case, it can be easily done with the same procedure using the ‘t.test’ function:

```
tt <- t.test(Holz$t01_visperc[Holz$female == 1],
             Holz$t01_visperc[Holz$female == 2])
critical(tt)
```

```
#>
#> Welch Two Sample t-test
#>
#> data: Holz$t01_visperc[Holz$female == 1] and Holz$t01_visperc[Holz$female == 2]
#> t = 1.4095, df = 298.9, p-value = 0.1597
#> alternative hypothesis: true difference in means is not equal to 0
#> 95 percent confidence interval:
#> -0.06419229 0.38823118
#> sample estimates:
#> mean of x mean of y
#> 4.314090 4.152071
#>
#> |== Effect Size and Critical Value ==|
#> d = 0.1623665 dc = ± 0.2269609 bc = ± 0.2262117
#> g = 0.1619587 gc = ± 0.2263909
```

The output now gives a wider range of values: the Cohen's  $d$  calculated on the data ( $d$ ), the critical Cohen's  $d$  ( $dc$ ), the numerator of the formula for the critical Cohen's  $d$  ( $bc$ ), the Cohen's  $d$  adjusted for small samples ( $g$ ) and the critical Cohen's  $d$  adjusted for small samples ( $gc$ ).

In the next example we will show the use of the package's function `critical_cor` to calculate the critical value for a correlation in a prospective framework.

```
n <- 60
critical_cor(n = n, hypothesis = "two.sided", test = "z")

#> $rc
#> [1] 0.2539247
#>
#> $df
#> [1] 58
#>
#> $test
#> [1] "z"
```

The direction of the hypothesis and the test to apply, either  $t$ -test or  $z$ -test, should be specified. The output will return the critical correlation value, the degrees of freedom and the type of test used.

## Discussion

With the present article, we propose that researchers compute and report the “critical effect size value(s)” in their empirical articles. This is not intended to replace other strategies aimed at enhancing the NHST approach to inference. Such strategies, such as the emphasis on estimating effect sizes with confidence intervals (Transue, 2019) or the a priori planning for statistical power are valuable in their own right. Instead, our proposal

serves as a complementary tool, especially beneficial for facilitating the interpretation of results when statistical power deviates from an optimal level (typically falling below, but occasionally exceeding it). Interestingly, the critical value can be retrospectively applied even to already published studies. This potential facilitates potential reframing of the original interpretations. Serving as a tool for retrospective analysis, the critical value may enable a reconsideration of the relevance of previously reported findings.

An advantage of reporting the critical value is that it can be precisely computed in any scenario, without requiring assumptions about the expected effect size, as is the case with power calculations. The critical value represents a directly interpretable benchmark that is especially useful in situations where statistical power is below the desired level and researchers are left otherwise uncertain about how to proceed with the interpretation of a study findings. For example, let us say that we read a published article reporting some effects as statistically significant, while others as not: we suspect that the study may be underpowered, but we are widely uncertain about the magnitude of possible true effects. To what extent can the reported results be interpreted, precisely? Knowing the critical value provides us with a clear benchmark. Conversely, let us say that an effect achieves significance in a very large sample: researchers tend to draw substantive conclusions based on this. But is it of real theoretical relevance? If in comparing two groups, such as controls versus treatment, any Cohen's  $d > 0.07$  would reach significance, is statistical significance enough to signal a "successful" treatment? Maybe yes, even if effects are tiny (e.g., Funder & Ozer, 2019), but knowing the critical value certainly prompts some appropriate interpretive caution.

Reporting the critical value(s) can also be an efficient way to allow researchers to evaluate which findings are statistically significant. For example, in a correlation table researchers customarily add an asterisk to all statistically significant correlations. But as long as all correlations are based on the same sample size, researchers can simply remark



‘the critical effect size is  $r = 0.3$ ’ and readers will know all correlations larger than this value are statistically significant.

Beyond enhancing study design and statistical inferences based on hypothesis tests, reporting critical value(s) can also serve an educational purpose. It underscores how the distinction between a significant and non-significant result is not solely determined by the presence or absence of a true effect, but also by the sample size. By highlighting a critical value, researchers can become more aware of the possibility of Type 2 errors when results are non-significant. Conversely, in studies with exceedingly large samples and in many meta-analyses, the critical value(s) may serve as a reminder that any observed effect larger than a trivially small value will likely achieve significance. This emphasizes that the mere attainment of statistical significance in a test is not particularly surprising, especially in non-experimental studies.

Real-case scenarios may not always be that simple. Hence, we chose to expand the application of computing critical significance values beyond Cohen’s  $d$  and correlation to include linear regression with both raw and standardized coefficients. This serves as a first step in computing critical values for a wider array of effects encountered in practical scenarios, where linear models and their extensions are commonly utilized for modeling purposes. A prerequisite is that researchers must be able to identify what parameters in their statistical models reflect the effect sizes of interest, and that they are able to assess their relevance. Notably, however, this prerequisite aligns with the requirements of APA style guidelines concerning the reporting of effect sizes. For further illustration and application, additional examples are provided in the Supplementary online material.

We suggest that reporting critical values is particularly valuable when sample size planning was not feasible or did not occur *a priori*. In cases where optimal power can be attained with a sufficiently large sample size for an effect of a specific magnitude of interest, and this is truly determined *a priori*, the interpretation of both significance and

non-significance becomes straightforward. However, when power analysis did not inform the sample size or when power is likely but undeterminedly low, reporting critical values for the obtained sample can help provide context for interpretation. Critical values can be computed and interpreted even retrospectively or for studies that have already been published.

In conclusion, the reporting of the critical value in empirical articles serves as a valuable addition to researchers' toolkit, aimed at augmenting transparency and facilitating the interpretability of their findings. While not designed to supplant existing practices, it provides a useful aid in interpreting newly presented and previously published results, thus advancing the understanding of research outcomes.

## References

- Altoè, G., Bertoldo, G., Zandonella Callegher, C., Toffalini, E., Calcagnì, A., Finos, L., & Pastore, M. (2020). Enhancing statistical inference in psychological research via prospective and retrospective design analysis. *Frontiers in Psychology, 10*, 499756.
- Cohen, J. (1994). The earth is round ( $p < .05$ ). *American Psychologist, 49*(12), 997–1003.
- Collaboration, O. S. (2015). Estimating the reproducibility of psychological science. *Science, 349*(6251), aac4716.
- Funder, D. C., & Ozer, D. J. (2019). Evaluating effect size in psychological research: Sense and nonsense. *Advances in Methods and Practices in Psychological Science, 2*(2), 156–168.
- Gelman, A. (2008). Scaling regression inputs by dividing by two standard deviations. *Statistics in Medicine, 27*, 2865–2873. <https://doi.org/10.1002/sim.3107>
- Gigerenzer, G., Krauss, S., & Vitouch, O. (2004). The null ritual. *The Sage Handbook of Quantitative Methodology for the Social Sciences*, 391–408.
- Hedges, L. V. (1981). Distribution theory for glass's estimator of effect size and related estimators. *Journal of Educational and Behavioral Statistics: A Quarterly Publication Sponsored by the American Educational Research Association and the American Statistical Association, 6*, 107–128. <https://doi.org/10.3102/10769986006002107>
- Hedges, L. V. (1984). Estimation of effect size under nonrandom sampling: The effects of censoring studies yielding statistically insignificant mean differences. *Journal of Educational Statistics, 9*(1), 61–85.
- Kramer, A. D., Guillory, J. E., & Hancock, J. T. (2014). Experimental evidence of massive-scale emotional contagion through social networks. *Proceedings of the National Academy of Sciences of the United States of America, 111*(24), 8788.
- Krueger, J. (2001). Null hypothesis significance testing: On the survival of a flawed method. *American Psychologist, 56*(1), 16.
- Lakens, D. (2013). Calculating and reporting effect sizes to facilitate cumulative science: A

- practical primer for t-tests and ANOVAs. *Frontiers in Psychology*, 4, 62627.
- Lakens, D. (2022). Sample size justification. *Collabra: Psychology*, 8(1), 33267.
- McCrea, S. M. (2008). Self-handicapping, excuse making, and counterfactual thinking: Consequences for self-esteem and future motivation. *Journal of Personality and Social Psychology*, 95(2), 274.
- Morris, S. B., & DeShon, R. P. (2002). Combining effect size estimates in meta-analysis with repeated measures and independent-groups designs. *Psychological Methods*, 7(1), 105.
- Orben, A., & Lakens, D. (2020). Crud (Re)Defined. *Advances in Methods and Practices in Psychological Science*, 3(2), 238–247. <https://doi.org/10.1177/2515245920917961>
- Szucs, D., & Ioannidis, J. P. (2017). Empirical assessment of published effect sizes and power in the recent cognitive neuroscience and psychology literature. *PLoS Biology*, 15(3), e2000797.
- Transue, B. (2019). *APA style 7th edition*.
- Wilson, B. M., Harris, C. R., & Wixted, J. T. (2020). Science is not a signal detection problem. *Proceedings of the National Academy of Sciences*, 117(11), 5559–5567.