Justify Your Alpha: A Primer on Two Practical Approaches

Maximilian Maier[1] & Daniël Lakens[2]

[1] University of Amsterdam, The Netherlands

[2] Eindhoven University of Technology, The Netherlands

Author Note

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

16                                                Abstract

17   The default use of an alpha level of 0.05 is suboptimal for two reasons. First, decisions

18   based on data can be made more efficiently by choosing an alpha level that minimizes the

19   combined Type 1 and Type 2 error rate. Second, it is possible that in studies with very

20   high statistical power $p$-values lower than the alpha level can be more likely when the null

21   hypothesis is true than when the alternative hypothesis is true (i.e., Lindley's paradox).

22   This manuscript explains two approaches that can be used to justify a better choice of an

23   alpha level than relying on the default threshold of 0.05. The first approach is based on

24   the idea to either minimize or balance Type 1 and Type 2 error rates. The second

25   approach lowers the alpha level as a function of the sample size to prevent Lindley's

26   paradox. An R package and Shiny app are provided to perform the required calculations.

27   Both approaches have their limitations (e.g., the challenge of specifying relative costs and

28   priors), but can offer an improvement to current practices, especially when sample sizes

29   are large. The use of alpha levels that are better justified should improve statistical

30   inferences and can increase the efficiency and informativeness of scientific research.

31   *Keywords:* Hypothesis Testing, Type 1 Error, Type 2 Error, Statistical Power

32   Word count: 7413 words

Justify Your Alpha: A Primer on Two Practical Approaches

Scientists regularly need to make dichotomous decisions when they perform lines of research. Should a pilot study be performed, or not? When multiple possible manipulations or measures are available, which should be used for the next study? Should the design of a study include a possible moderator, or can it be ignored? Should a research line be continued, or abandoned? These decisions come with costs and benefits for the scientist, as well as for society, when bad decisions lead to research waste. In a Neyman-Pearson approach to hypothesis testing (J. Neyman & Pearson, 1933) studies are designed such that erroneous decisions that determine how we act are controlled in the long run at some desired maximum level. If resources were infinite we could collect enough data to make the chance of a wrong decision incredibly small by using an extremely low alpha level while still achieving very high statistical power. However, since resources are limited, researchers need to decide how to choose the rate at which they are willing to make errors (Wald, 1949). After data is collected researchers can incorrectly act as if there is an effect when there is no true effect (a Type 1 error) or incorrectly act as if there is no effect when there is a true effect (a Type 2 error). With the same number of observations, a reduction in the Type 1 error rate will increase the Type 2 error rate (and vice versa).

The question how error rates should be set in any study requires careful consideration of the relative costs of a Type 1 error or a Type 2 error. Regrettably, researchers rarely provide such a justification and predominantly use an alpha level of 5%. In the past, the strong convention to use a 5% alpha level might have functioned as a de facto prespecification of the alpha level, which was useful given that the alpha level needs to be decided upon before the data is analyzed (Uygun-Tunç, Tunç, & Lakens, 2021). Nowadays, researchers can transparently preregister a statistical analysis plan in an online repository, which makes it possible to specify more appropriate but less conventional alpha levels. Even though it is possible to preregister non-conventional alpha levels, there is relatively little practical guidance on how to choose an alpha level

for a study. This article explains why error rates need to be justified and provides two practical approaches that can be used to justify the alpha level. In the first approach the cost of Type I and Type II error rates are balanced or minimized and in the second approach the alpha level is lowered as a function of the sample size.

## Why Do We Use a 5% Alpha Level and 80% Power?

We might naively assume that when all researchers do something, there must be a good reason for such an established practice. An important step towards maturity as a scholar is the realization that this is not the case. Neither Fisher nor Neyman, two statistical giants largely responsible for the widespread reliance on hypothesis tests in the social sciences, recommended the universal use of any specific threshold. Ronald A. Fisher (1971) writes: "It is open to the experimenter to be more or less exacting in respect of the smallness of the probability he would require before he would be willing to admit that his observations have demonstrated a positive result." Similarly, J. Neyman and Pearson (1933) write: "From the point of view of mathematical theory all that we can do is to show how the risk of the errors may be controlled and minimized. The use of these statistical tools in any given case, in determining just how the balance should be struck, must be left to the investigator."

Even though in *theory* alpha levels should be justified, in *practice* researchers tend to imitate others. R. A. Fisher (1926) notes: "Personally, the writer prefers to set a low standard of significance at the 5 per cent point, and ignore entirely all results which fail to reach this level." This sentence is preceded by the statement "If one in twenty does not seem high enough odds, we may, if we prefer it, draw the line at one in fifty (the 2 percent point), or one in a hundred (the 1 percent point)." Indeed, in his examples Fisher often uses an alpha of 0.01. Nevertheless, researchers have copied the value Fisher preferred, instead of his more important take-home message that the significance level should be set by the experimenter. The default use of an alpha level of 0.05 can already be found in work of Gosset on the *t*-distribution (Cowles & Davis, 1982; Kennedy-Shaffer, 2019), who believed that a difference of two standard deviations (a z-score of 2) was sufficiently rare.

The default use of 80% power (or a 20% Type 2, or beta (b) error) is similarly based on personal preferences by Cohen (1988), who writes: "It is proposed here as a convention that, when the investigator has no other basis for setting the desired power value, the value .80 be used. This means that beta is set at .20. This value is offered for several reasons (Cohen, 1965, pp. 98-99). The chief among them takes into consideration the implicit convention for alpha of .05. The beta of .20 is chosen with the idea that the general relative seriousness of these two kinds of errors is of the order of .20/.05, i.e., that Type I errors are of the order of four times as serious as Type II errors. This .80 desired power convention is offered with the hope that it will be ignored whenever an investigator can find a basis in his substantive concerns in his specific research investigation to choose a value ad hoc."

We see that conventions are built on conventions: the norm to aim for 80% power is built on the norm to set the alpha level at 5%. This normative use of statistics was criticized in a statement by the American Statistical Association (Wasserstein & Lazar, 2016), who wrote: "We teach it because it's what we do; we do it because it's what we teach." The real lesson we should take away from Cohen is to determine the relative seriousness of Type 1 and Type 2 errors, and to balance both types of errors when a study is designed. If a Type 1 error is considered to be four times as serious as a Type 2 error, the *weighted* error rates in the study are balanced with a 5% Type 1 error rate and a 20% Type 2 error rate.

**Justifying the Alpha Level**

In 1957 Neyman wrote: "it appears desirable to determine the level of significance in accordance with quite a few circumstances that vary from one particular problem to the next" (Jerzy Neyman, 1957). Despite this advice, the mindless application of null hypothesis significance tests, including setting the alpha level at 5% for all tests, is so universal that it has been criticized for more than half a century (Bakan, 1966; Gigerenzer, 2018). The default use of a 5% alpha level might have been difficult to abandon, even if it was a mediocre research practice, without an alternative approach in

117 which alpha levels are better justified.

118      There are two main reasons to abandon the universal use of a 5% alpha level. The

119 first reason to carefully choose an alpha level is that decision-making becomes more

120 efficient (Mudge, Baker, Edge, & Houlahan, 2012). If researchers use hypothesis tests to

121 make dichotomous decisions from a methodological falsificationist approach to statistical

122 inferences (Uygun-Tunç et al., 2021), and have a certain maximum sample size they are

123 willing or able to collect, it is typically possible to make decisions more efficiently by

124 choosing error rates such that the combined cost of Type 1 and Type 2 errors is

125 minimized. If we aim to either minimize or balance Type 1 and Type 2 error rates for a

126 given sample size and effect size, the alpha level should be set not based on convention,

127 but by weighting the relative cost of both types of errors.

128      The second reason is most relevant for large data sets (Harford, 2014). As the

129 statistical power increases, some $p$-values below 0.05 (e.g., $p = 0.04$) can be more likely

130 when there is *no* effect than when there *is* an effect. This is known as Lindley's paradox

131 (Bartlett, Jordan, & Mcauliffe, 1957; Cousins, 2017; Jeffreys, 1935, 1936b, 1936a; Lin,

132 Lucas Jr, & Shmueli, 2013; Lindley, 1957), or sometimes the Jeffreys-Lindley paradox

133 (Spanos, 2013), as Harold Jeffreys discussed the paradox long before Lindley

134 (Wagenmakers & Ly, n.d.). The distribution of $p$-values is a function of the statistical

135 power (Cumming, 2008), and the higher the power, the more right-skewed the

136 distribution becomes (i.e., the more likely it becomes that small $p$-values are observed).

137 When there is no true effect $p$-values are uniformly distributed, and 1% of observed

138 $p$-values fall between 0.04 and 0.05. When the statistical power is extremely high, not

139 only will most $p$-values fall below 0.05, most will also fall below 0.01. In Figure 1 we see

140 that with high power very small $p$-values are more likely to be observed when there *is* an

141 effect than when there is *no* effect (e.g., the black curve representing $p$-values when the

142 alternative is true falls above the dashed horizontal line for a $p$-value of 0.01). But

143 observing a $p$-value of 0.04 is more likely when the null hypothesis (H0) is true than when

144 the alternative hypothesis (H1) is true and we have very high power, as illustrated by the

**P−value distribution for d = 0.5 and N = 150**



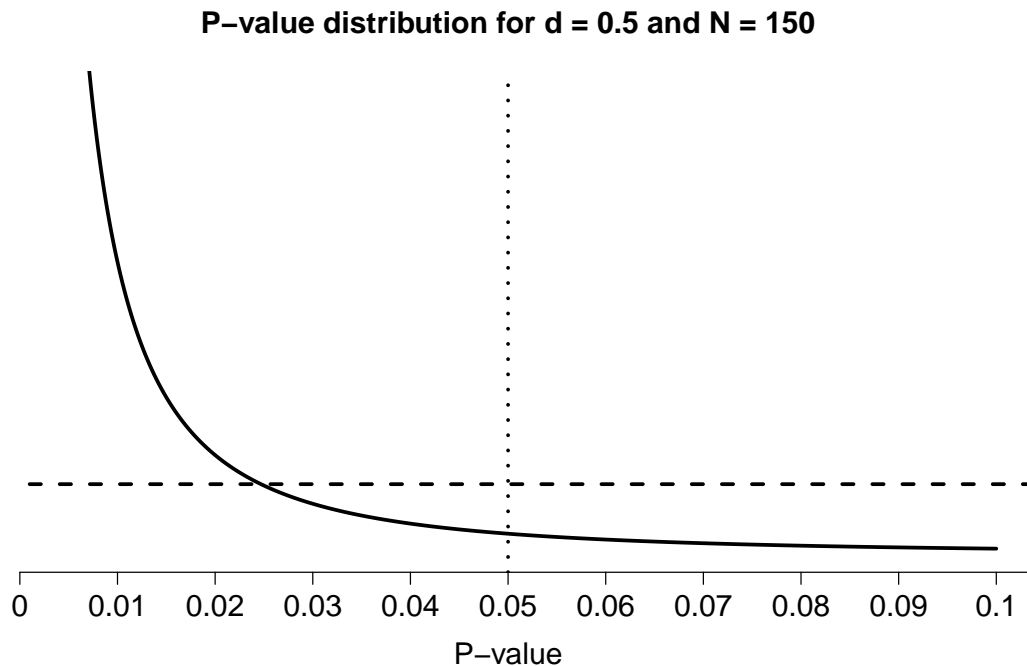*Figure 1*. *P*-value distributions for a two-sided independent *t*-test with N = 150 and d = 0.5 (black curve) or d = 0 (horizontal dashed line) which illustrates how *p*-values just below 0.05 can be more likely when there is no effect than when there is an effect.

147    Although it is not necessary from a Neyman-Pearson error-statistical perspective,

148 researchers often want to interpret a significant test result as evidence for the alternative

149 hypothesis. In other words, in addition to controlling the *error rate*, researchers might be

150 interested in interpreting the *relative evidence* in the data for the alternative hypothesis

151 over the null hypothesis. If so, it makes sense to choose the alpha level such that when a

152 significant *p*-value is observed, the *p*-value is actually more likely when the alternative

153 hypothesis is true than when the null hypothesis is true. This means that when statistical

154 power is very high (e.g., the sample size is very large), the alpha level should be reduced.

155 For example, if the alpha level in Figure 1 is lowered to 0.02 then the alternative

156 hypothesis is more likely than the null hypothesis for all significant *p*-values that would

157 be observed. This approach to justifying the alpha level can be seen as a

158 frequentist/Bayesian compromise (Good, 1992). The error rate is controlled, but at the

same time the alpha level is set to a value that guarantees that whenever we reject the null hypothesis, the data is more likely under the alternative hypothesis than under the null.

**Minimizing or Balancing Type 1 and Type 2 Error Rates**

If both Type 1 as Type 2 errors are costly, then it makes sense to optimally reduce both errors as you design studies. This idea is well established in applied statistics (Cornfield, 1969; DeGroot, 1975; Kim & Choi, 2021; Lindley, 1953; Mudge et al., 2012; Pericchi & Pereira, 2016) and leads to studies where you make decisions most efficiently. Researchers can choose to design a study with a statistical power and alpha level that minimizes the *weighted combined error rate*. For example, a researcher designs an experiment where they assume H0 and H1 are a-priori equally probable (the prior probability for both is 0.5). They set the Type 1 error rate to 0.05 and collect sufficient data such that the statistical power is 0.80. The weighted combined error rate is 0.5 (the probability H0 is true) × 0.05 (the probability of a Type 1 error) + 0.5 (the probability that H1 is true) × 0.20 (the probability of a Type 2 error) = 0.125. This weighted combined error rate might be lower if a different choice for Type 1 and Type 2 errors was made.

Assume that in the previous example data will be analyzed in an independent *t*-test and the researcher was willing to collect 64 participants in each condition to achieve the 0.05 Type 1 error rate and 0.8 power. The researcher could have chosen to set the alpha level in this study to 0.1 instead of 0.05. If the Type 1 error rate is 0.1, the statistical power (given the same sample size of 64 per group) would be 0.88. The weighted combined error rate is now $(0.5 \times 0.1 + 0.5 \times 0.12) = 0.11$. In other words, increasing the Type 1 error rate from 0.05 to 0.1 reduced the Type 2 error rate from 0.2 to 0.12 and the combined error rate from 0.125 to 0.11. In the latter scenario, our total probability of making an erroneous decision has become 0.015 smaller. As shown below, this approach can be extended to incorporate scenarios where the prior probability of H0 and H1 differ. Mudge et al. (2012) and Kim and Choi (2021) show that by choosing an alpha level

based on the relative weight of Type 1 errors and Type 2 errors and assuming beliefs about the prior probability that H0 and H1 are correct, decisions can be made more efficiently than when the default alpha level of 0.05 is used. Kim (2020) also provides an R-package to justify the alpha level based on decision-theoretic approaches, which provides solutions for a smaller set of power functions than the package accompanying this paper, and only allows users to minimize the costs of errors.

Winer (1962) writes: "The frequent use of the .05 and .01 levels of significance is a matter of convention having little scientific or logical basis. When the power of tests is likely to be low under these levels of significance, and when Type 1 and Type 2 errors are of approximately equal importance, the .30 and .20 levels of significance may be more appropriate than the .05 and .01 levels." The reasoning here is that a design that has 70% power for the smallest effect size of interest would not balance the Type 1 and Type 2 error rates in a sensible manner. Similarly, and perhaps more importantly, one should carefully reflect on the choice of the alpha level when an experiment achieves very high statistical power for all effect sizes that are considered meaningful. If a study has 99% power for effect sizes of interest, and thus a 1% Type 2 error rate, but uses the default 5% alpha level, it also suffers from a lack of balance. This latter scenario is quite common in meta-analyses, where researchers by default use a 0.05 alpha level, while the meta-analysis often has very high power for all effect sizes of interest. It is also increasingly common when analyzing large existing data sets or when collecting thousands of observations online. In such cases where power for all effects of interest is very high, it is sensible to lower the alpha level for statistical tests to reduce the weighted combined error rate and increase the severity of the test.

Researchers can decide to either balance Type 1 and Type 2 error rates (e.g., designing a study such that the Type 1 and Type 2 error rate are equal) or minimize the weighted combined error rate. For any given sample size and effect size of interest there is an alpha level that minimizes the weighted combined Type 1 and Type 2 error rates. Because the chosen alpha level also influences the statistical power, and the Type 2 error

215 rate is therefore dependent upon the Type 1 error rate, minimizing or balancing error

216 rates requires an iterative optimization procedure.

217     As an example, imagine a researcher who plans to perform a study which will be

218 analyzed with an independent two-sided $t$-test. They will collect 50 participants per

219 condition, and set their smallest effect size of interest to Cohen's d = 0.5. They think a

220 Type 1 error is just as costly as a Type 2 error, and believe H0 is just as likely to be true

221 as H1. The weighted combined error rate is minimized when they set alpha to 0.13 (see

222 Figure 2, dotted line), which will give the study a Type 2 error rate of beta = 0.166 to

223 detect effects of d = 0.5. The weighted combined error rate is 0.148, while it would have

224 been 0.177 if the alpha level was set at 5%[1].

225     We see that increasing the alpha level from the normative 5% level to 0.13 reduced

226 the weighted combined error rate - any larger or smaller alpha level would increase the

227 weighted combined error rate. The reduction in the weighted combined error rate is not

228 huge, but we have reduced the overall probability of making an error. More importantly,

229 we have chosen an alpha level based on a justifiable principle, and clearly articulated the

230 relative costs of a Type 1 and Type 2 error. Perhaps counter-intuitively, decision-making

231 is sometimes slightly more efficient after *increasing* the alpha level from the default of

232 0.05 because a small increase in the Type 1 error rate can lead to a larger decrease in the

233 Type 2 error rate. Had the sample size been much smaller, such as n = 10, the solid line

234 in Figure 2 shows that the weighted combined error rate will always be high, but it is

235 minimized if we increase the alpha level to alpha to 0.283. If the sample size had been n

236 = 100, the optimal alpha level to minimize the weighted combined error rate (still

237 assuming H0 and H1 have equal probabilities, and Type 1 and Type 2 errors are equally

238 costly) is 0.0509 (the long-dashed line in Figure 2).

───────

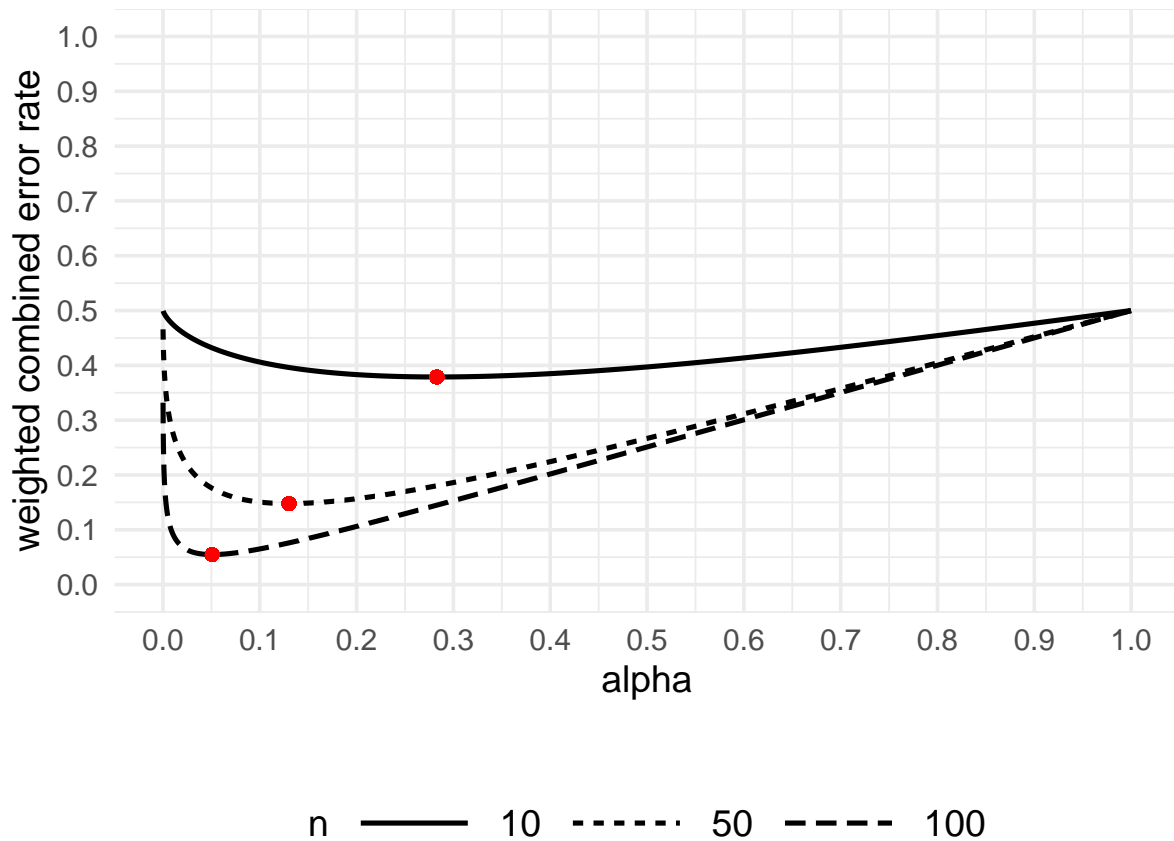[1] For the same scenario, balanced error rates are alpha = 0.149 and beta = 0.149.

*Figure 2*. Weighted combined error rate (y-axis) for an independent *t*-test with n = 10, n = 50, and n = 100 per group and a smallest effect of interest of d = 0.5, for all possible alpha levels (x-axis).

## Weighing the Relative Cost of Errors

Cohen (1988) recommended a study design with a 5% Type 1 error rate and a 20% Type 2 error rate. He personally felt "Type I errors are of the order of four times as serious as Type II errors." However, some researchers have pointed out, following Neyman (1933), that false negatives might be more severe than false positives (Fiedler, Kutzner, & Krueger, 2012). The best way to determine the relative costs of Type 1 and Type 2 errors is by performing a cost-benefit analysis. For example, Field, Tyre, Jonzén, Rhodes, and Possingham (2004) quantify the relative costs of Type 1 errors when testing whether native species in Australia are declining. In this example, the H1 is that the Koala population is declining and H0 the Koala population is not declining. The Type 1 error

would be to decide that the Koala population is declining, when in fact it is not; a Type 2 error would be to decide that the Koala population is not declining, when in fact it is. Field et al. (2004) conclude that when it comes to the Koala population, given its great economic value, a cost-benefit analysis indicates the alpha level should be set to 1. In other words, one should always act as if the population is declining because the relative cost of a Type 2 error compared to a Type 1 error is too high. Note that in this example, the decision not to collect data is deterministically dominant (Clemen, 1997). The alpha of 1 shows that the results of the data collection will not influence future decisions in any way - it is always beneficial to intervene. This is arguably rare, but not incredibly rare. If you are bitten by an animal, it is possible to observe the animal for 10 days to see if it has rabies, but given the costs and benefits, it is more cost-efficient to assume the animal has rabies and get a rabies shot. In psychology, it is possible that accurate pilot studies to determine which two possible manipulations has a larger effect size will require a larger sample than if one designs a study conservatively powered for the manipulation that based on a personal prior is believed to have the smallest effect size. There are similar situations where researchers might decide to skip a pilot study and immediately perform the main experiment because this is the most rational choice.

An applied example where the decision is not deterministically dominant can be found in Viamonte, Ball, and Kilgore (2006) who evaluate the benefits of a computerized intervention aimed at improving speed of processing to reduce car collisions in people aged 75 or older. They estimated that the risk of getting into an accident for these older drivers is 7.1%. The cost of a collision was estimated to be $22,000, or $22,000 * 0.071 = 1,562.84 per driver in the USA. Furthermore, they estimate that the intervention can prevent accidents for 86% of drivers. Therefore, the probability of a collision after intervention is now (1-0.86) * 0.071 = 0.00994. The total cost of completing the intervention was estimated to be $274.50. When the intervention is implemented, some drivers will still get into a collision, so the total cost of the intervention and collisions is $493.30 per driver ($274.50 + 0.00994 * $22,000).

277     We can implement the intervention when it does not actually work, making a Type
278 1 error. The waste is $274.50 per driver, as this is what the intervention costs even if it
279 offers no benefits. If the intervention works, but it is not implemented, we make a Type 2
280 error and the amount of money that is not saved is $1,562.84 (the cost of doing nothing) -
281 493.30 (the cost if the intervention was implemented), for a waste of 1.069,54 per driver.
282 This means that the relative cost of a Type 1 error compared to a Type 2 error is 274.50
283 /1.069,54 = 0.257, or the waste in money after a Type 1 error is 3.896 times
284 (1.069,54/274.50) worse than a Type 2 error. This ratio reflects that the intervention is
285 relatively cheap, and therefore a Type 1 error is not that costly, while the potential
286 savings if collisions are prevented is relatively large. Of course, quantifying costs and
287 benefits comes with uncertainties. The intervention might prevent more or less accidents,
288 the risks of an accident for drivers of 75 years or older might be greater or smaller,
289 etcetera. Sensitivity analyses can be used to compute a range of the ratio of the costs of
290 Type 1 and Type 2 errors (see Viamonte et al., 2006).

291     Although it can be difficult to formally quantify all relevant factors that influence
292 the costs of Type 1 and Type 2 errors, there is no reason to let the perfect be the enemy
293 of the good. In practice, even if researchers don't explicitly discuss their choice for the
294 relative weight of Type 1 versus Type 2 errors, they make a choice in every hypothesis
295 test they perform, even if they simply follow conventions (e.g., a 5% Type 1 error rate
296 and a 20% Type 2 error rate). It might be especially difficult to decide upon the relative
297 costs of Type 1 and Type 2 errors when there are no practical applications of the research
298 findings, but even in these circumstances, it is up to the researcher to make a decision
299 (Douglas, 2000). It is, therefore, worth reflecting on how researchers can start to think
300 about the relative weight of Type 1 and Type 2 errors.

301     First, if a researcher only cares about not making a decision error, but the
302 researcher does not care about whether this decision error is a false positive or a false
303 negative, Type 1 and Type 2 errors are weighed equally. Therefore, weighing Type 1 and
304 Type 2 errors equally is a defensible default, unless there are good arguments to weigh

false positives more strongly than false negatives (or vice versa). When deciding upon whether there is a reason to weigh Type 1 and Type 2 errors differently, researchers are in essence performing a multiple criterion decision analysis (Edwards, Miles Jr., & Winterfeldt, 2007), and it is likely that treating the justification of the relative weight of Type 1 and Type 2 errors as a formal decision analysis would be a massive improvement over current research practices. A first step is to determine the objectives of the decision that is made in the hypothesis test, assign attributes to measure the degree to which these objectives are achieved within a specific time-frame (Clemen, 1997), and finally to specify a value function. In a hypothesis test, we do not simply want to make accurate decisions, but we want to make accurate decisions given the resources we have available (e.g., time and money). Incorrect decisions have consequences, both for the researcher themselves, as for scientific peers, and sometimes for the general public. We know relatively little about the actual costs of publishing a Type 1 error for a researcher, but in many disciplines the costs of publishing a false claim are low, while the benefits of an additional publication on a resume are large. However, by publishing too many claims that do not replicate, a researcher risks gaining a reputation for publishing unreliable work. In addition, researcher might plan to build on work in the future, as might peers. The costs of experiments that follow up on a false lead might be much larger than the cost to reduce the possibility of a Type 1 error in an initial study, unless replication studies are cheap, will be performed anyway and will be shared with peers. However, it might also be true that the hypothesis has great potential for impact if true and the cost of a false negative might be substantial whenever it closes off a fruitful avenue for future research. A Type 2 error might be more costly than a Type 1 error, especially in a research field where all findings are published and people regularly perform replication studies to identify Type 1 errors in the literature (Fiedler et al., 2012).

Another objective might be to influence policy, in which case the consequences of a Type 1 and Type 2 error should be weighed by examining the relative costs of implementing a policy that does not work against not implementing a policy that works. The second author once attended a presentation by a policy advisor who decided whether

new therapies would be covered by the national healthcare system. She discussed Eye

Movement Desensitization and Reprocessing (EMDR) therapy. She said that, although

the evidence for EMDR was weak at best, the costs of the therapy (which can be done

behind a computer) are very low, it was applied in settings where no good alternative

therapies were available (e.g., inside prisons), and risk of negative side-effects was

basically zero. They were aware of the fact that there was a very high probability that

the claim that EMDR was beneficial might be a Type 1 error, but the cost of a Type 1

error was deemed much lower than the cost of a Type 2 error.

Imagine a researcher plans to collect 64 participants per condition to detect a d =

0.5 effect, and weighs the cost of Type 1 errors 4 times as much as Type 2 errors. To

minimize error rates, the Type 1 error rate should be set to 0.0327, which will make the

Type 2 error rate 0.252. If we would perform 20000 studies designed with these error

rates, and assume H0 and H1 are equally likely to be true, we would observe 0.5 (the

prior probability that H0 is true) $\times$ 0.0327 (the alpha level) $\times$ 20000 = 327 Type 1 errors,

and 0.5 (the prior probability that H1 is true) $\times$ 0.252 (the Type 2 error rate) $\times$ 20000 =

2524 Type 2 errors. Since we weigh Type 1 errors 4 times as much as Type 2 errors, we

multiple the cost of the 327 Type 1 errors by 4, which makes $4 \times 327 = 1308$, and to

keep the weighted error rate between 0 and 1, we also multiply the 10000 studies where

we expect H0 to be true by 4, such that the weighted combined error rate is (1308 +

2524)/(40000 + 10000) = 0.0766. Figure 3 visualizes the weighted combined error rate for

this study design across the all possible alpha levels, and illustrated the weighted error

rate is smallest when the alpha level is 0.0327.

If the researcher had decided to *balance* error rates instead of *minimizing* error

rates, we recognize that with 64 participants per condition, we are exactly in the scenario

Cohen (1988) described. When Type 1 errors are considered 4 times as costly as Type 2

errors, 64 participants per condition yield a 5% Type 1 error rate and a 20% Type 2 error

rate. If we would increase the sample size, The Type 1 and Type 2 error rates would

remain in a balanced 1:4 ratio, but both error rates would be smaller. With a smaller
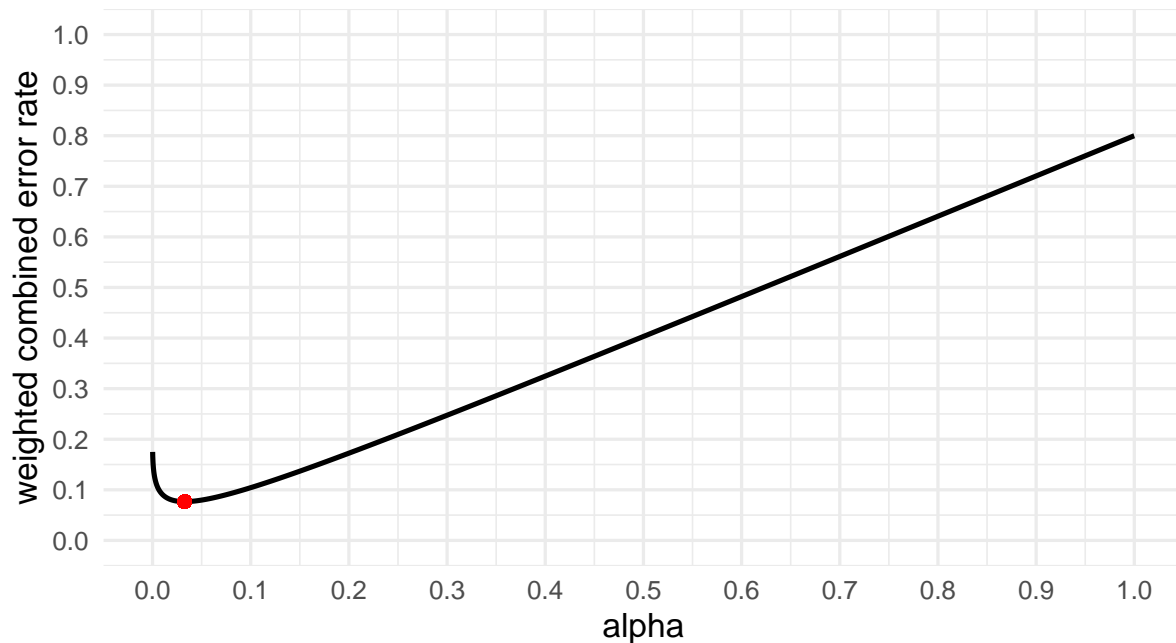
*Figure 3*. Weighted combined error rate (y-axis) for an independent *t*-test with n = 64 per group and a smallest effect of interest of d = 0.5, where Type 1 errors are weighed 4 times as much as Type 2 errors, for all possible alpha levels (x-axis).

362  sample size, both error rates would be larger.

363  **Incorporating Prior Probabilities**

364      The choice for an optimal alpha level depends not just on the relative costs of Type
365  1 and Type 2 errors, but also on the base rate of true effects (Miller & Ulrich, 2019). In
366  the extreme case, in all studies a researcher designs H1 is true. In this case, there is no
367  reason to worry about Type 1 errors, because a Type 1 error can only happen when the
368  null hypothesis is true. Therefore, you can set the alpha level to 1 without any negative
369  consequences. On the other hand, if the base rate of true H1s is very low, you are more
370  likely to test a hypothesis where H0 is true. Therefore, the probability of observing a
371  false positive becomes a more important consideration. Whatever the prior probabilities
372  are believed to be, researchers always need to specify the prior probabilities of H0 and H1.
373  Researchers should take their expectations about the probability that H0 and H1 are true
374  into account when evaluating costs and benefits.

375    For example, let's assume a researcher performs 1000 studies. The researcher

376  expects 100 studies to test a hypothesis where H1 is true, while the remaining 900 studies

377  test a hypothesis where H0 is true. This means H0 is believed to be 9 times more likely

378  than H1, or equivalently, that the relative probability of H1 versus H0 is 0.1111:1.

379  However, the researcher decides to ignore these prior probabilities and designs a study

380  that has the normative 5% Type 1 error rate and a 20% Type 2 error rate. The

381  researcher should expect to observe 0.9 (the prior probability that H0 is true) × 0.05 (the

382  alpha level) × 1000 = 45.00 Type 1 errors, and 0.1 (the prior probability that H1 is true)

383  × 0.2 (the Type 2 error rate) × 1000 = 20.00 Type 2 errors, for a total of 65.00 errors.

384    However, the total number of errors does not tell the whole story, as Type 1 errors

385  are weighed four times more than Type 2 errors. We therefore need to compute the

386  weighted combined error rates $w$ taking the relative cost of Type 1 and Type 2 errors into

387  account, and the prior probabilities of H0 and H1, which can be done with the following

388  formula from Mudge et al. (2012):

$$\frac{(cost_{T1T2} \times \alpha + prior_{H1H0} \times \beta)}{prior_{H1H0} + cost_{T1T2}} \tag{1}$$

389    For the previous example, the weighted combined error rate is (4 × 0.05 + 0.1111 ×

390  0.2) / (0.1111 + 4) = 0.054. If the researcher had taken the prior probabilities into

391  account when deciding upon the error rates, a lower combined error rate can be achieved.

392  With the same sample size (64 per condition) the combined weighted error rate was not

393  as small as possible, optimally balanced error rates (maintaining the 4:1 ratio of the

394  weight of Type 1 versus Type 2 errors) would require setting alpha to 0.011 and the Type

395  2 error rate to 0.402. The researcher should now expect to observe 0.9 (the prior

396  probability that H0 is true) × 0.011 (the alpha level) × 1000 = 9.89 Type 1 errors, and

397  0.1 (the prior probability that H1 is true) × 0.402 (the Type 2 error rate) × 1000 = 40.16

398  Type 2 errors. The weighted error rate is 0.0216.

399    Because the prior probability of H0 and H1 influence the expected number of Type

400 1 and Type 2 errors one will observe in the long run, the alpha level should be lowered as

401 the prior probability of H0 increases, or equivalently, the alpha level should be increased

402 as the prior probability of H1 increases. Because the base rate of true hypotheses is

403 unknown, this step requires a subjective judgment. This can not be avoided, because one

404 always makes assumptions about base rates, even if the assumption is that a hypothesis

405 is equally likely to be true as false (with both H1 and H0 having a 50% probability). In

406 the previous example, it would also have been possible minimize (instead of balance) the

407 error rates, which is achieved with an alpha of 0.00344 and a beta of 0.558, for a total of

408 58.86 errors, where the weighted error rate is 0.0184.

409        The two approaches (balancing error rates or minimizing error rates) typically yield

410 quite similar results. Where minimizing error rates might be slightly more efficient,

411 balancing error rates might be slightly more intuitive (especially when the prior

412 probability of H0 and H1 is equal). Note that although there is always an optimal choice

413 of the alpha level, there is always a range of values for the alpha level that yield quite

414 similar weighted error rates, as can be seen in Figure 3.

**Increasing the Alpha Level Above 0.05**

416        Many empirical sciences have recently been troubled by a replication crisis

417 (Camerer et al., 2016; Open Science Collaboration, 2015), which has in part been caused

418 by inflated alpha levels due to $p$-hacking (Simmons, Nelson, & Simonsohn, 2011),

419 publication bias, and low statistical power (Lindsay, 2015). In light of this low

420 replicability, a potential concern about allowing researchers to justify their alpha level is

421 that researchers can decide to increase the alpha level above the 0.05 threshold. This

422 could increase the rate of false positives published in the literature compared to when an

423 alpha level of 0.05 remains the norm. An increase of the alpha level should only be

424 deemed acceptable when authors can justify that the costs of the increase in the Type 1

425 error rate is sufficiently compensated by the benefit of decreased Type 2 error rate.

426 Furthermore, researchers should explicitly accompany claims by their error rates

427 throughout an article, especially when the alpha level is increased, and readers of claims

made with higher alpha level should understand such claims are made with greater uncertainty, and could very well be false.

There are circumstances under which optimal error rates will require an increase of the alpha level, which will also increase the number of false positives in the literature. Assuming the goal of scientists is to efficiently generate reliable knowledge, the proposal to increase the alpha level (and thus to increase the Type 1 error rate in the literature) should only be adopted if the cost of an increase in Type 1 errors is compensated in some way. So far we have focussed only on how the increase in the Type 1 error rate will lead to a greater reduction in the Type 2 error rate, which all else being equal, should improve decision making in hypothesis tests. In practice, it might be a challenge to reach agreement on the weight of Type 1 and Type 2 errors among different stakeholders. For example, where a team of researchers might believe a Type 1 and Type 2 error is equally costly, an editor of a journal might weigh Type 1 errors more than Type 2 errors. We should also consider the possibility that researchers try to opportunistically specify the relative cost of Type 1 and Type 2 error rates to increase their alpha level, and increase the probability of finding a 'significant' effect.

Nevertheless, in some cases, it can be justified to increase the alpha level above the 0.05 threshold. These will usually be cases where (1) the study will have directly decision-making relevant implications (as in the above EDM example), (2) a cost-benefit analysis is provided that gives a clear rationale for relatively high costs of a Type 2 error, (3) the probability of H1 being false is relatively low, and (4) it is not feasible to reduce overall error rates by collecting more data. In these cases, it will often be desirable to justify the alpha level during the first phase of a Registered Report so that the higher alpha level that will be used in a study can be discussed transparently during peer-review. At the same time, given the complexity of weighing the costs and benefits of research, it is understandable if some journals consider such discussions too great a burden for reviewers. If so, these journals could indicate that they limit deviations from an alpha level of 0.05 only where researchers increase the severity of their test by lowering the alpha level.

Journals might also prefer to use a default alpha level of 0.05 to reduce the burden on readers to examine at which alpha level claims in their journal are made. Especially if an increase in alpha levels was not evaluated by peers during the first phase of a Registered Report, the evaluation of whether this alpha level was appropriate is left to readers. In practice, the use of a higher alpha level will require readers to keep track of the fact that the claim of the presence of an effect was less severely tested than it would have been with a default alpha, instead of keeping track of the fact that claims of the absence of an effect were less severely tested than they would have been when the statistical power had been higher (i.e., by increasing the alpha level). In a science where people only focus on significant effects and treat all significant effects as equally well supported, increasing alpha levels could lead to a sense of false certainty about a body of work. If the practice to increase alpha levels becomes popular, it will be important to examine whether varying alpha levels are taken into account when interpreting and discussing research findings, and how negative side-effects can be mitigated.

Finally, the use of a high alpha level might be missed if readers skim an article. We believe this can be avoided by having each scientific claim accompanied by the alpha level under which it was made. Scientists should be required to report their alpha levels prominently, usually in the abstract of a paper alongside a summary of the main claim. The correct interpretation of a hypothesis test was never to label an effect as 'significant' or 'nonsignificant' but to reject effects implied by the null model with a specific error rate. Replacing 'the effect was significant' with 'we reject an effect size of 0 with a 10% error rate' might end up improving the interpretation of hypothesis tests. Note that by explicitly reporting the alpha level alongside a claim it will also become more visible when researchers lower their alpha level, and this practice will therefore clearly communicate whenever readers should be impressed by the fact that a claim passed an even more severe test than if a traditional alpha level of 0.05 would have been used.

**Sample Size Justification when Minimizing or Balancing Error Rates**

So far we have illustrated how to perform what is known as a *compromise power analysis* where the weighted combined error rate is computed as a function of the sample size, the effect size, and the desired ratio of Type 1 and Type 2 errors (Erdfelder, Faul, & Buchner, 1996). However, in practice researchers will often want to justify their sample size based on an *a-priori power analysis* where the required sample size is computed to achieve desired error rates, given an effect size of interest (Lakens, 2021). It is possible to determine the sample size at which we achieve a certain desired weighted combined error rate. This requires researchers to specify the effect size of interest, the relative cost of Type 1 and Type 2 errors, the prior probabilities of H0 and H1, whether error rates should be balanced or minimized, and the desired weighted combined error rate.

Imagine a researcher is interested in detecting an effect of Cohen's d = 0.5 with a two-sample *t*-test. The researcher believes Type 1 errors are equally costly as Type 2 errors and believes a H0 is equally likely to be true as H1. The researcher desires a minimized weighted combined error rate of 5%. Figure 4 shows the optimal alpha level, beta, and weighed combined error rate as a function of sample size for this situation. We can optimize the weighted combined error rate as a function of the alpha level and sample size through an iterative procedure, which reveals that a sample size of 105 participants in each independent condition is required to achieve the desired weighted combined error rate. In the specific cases where the prior probability of H0 and H1 are equal, this sample size can also be computed directly with common power analysis software by entering the desired alpha level and statistical power. In this example, where Type 1 and Type 2 error rates are weighted equally, and the prior probability of H0 and H1 is assumed to be 0.5, the sample size is identical to that required to achieve an alpha of 0.05 and a desired statistical power for d = 0.5 of 0.95. Note that it might be difficult to specify the desired *weighted* combined error rate for a power analysis when Type 1 and Type 2 errors are not weighted equally, and/or when H1 and H0 are not deemed equally probable.
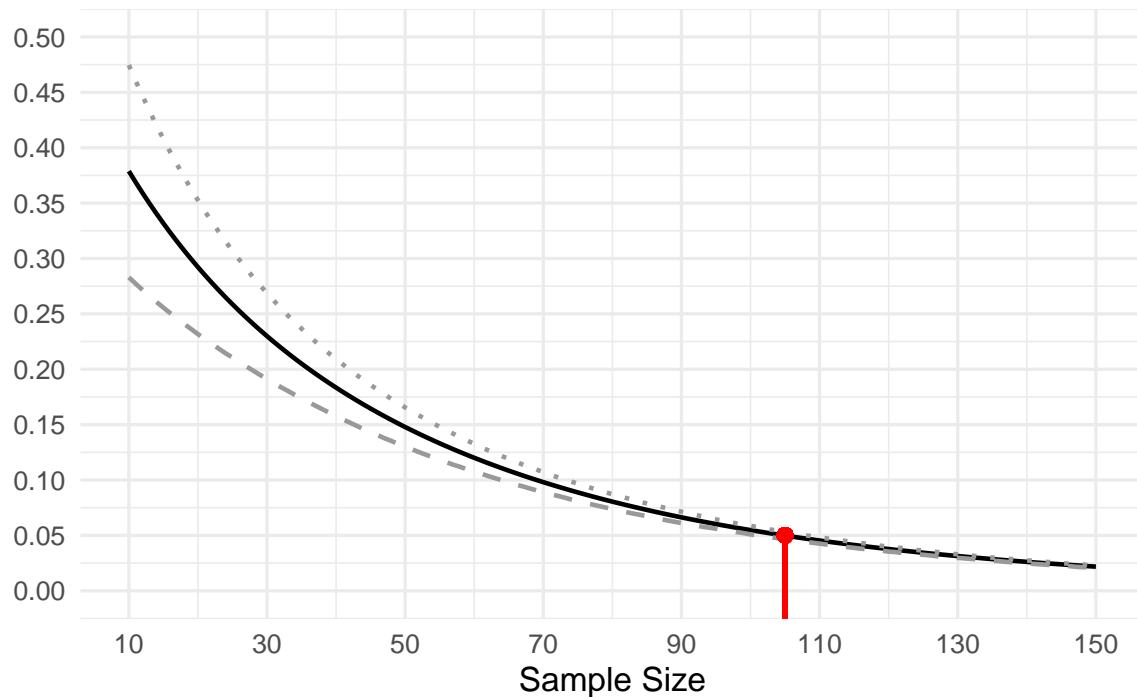
*Figure 4*. Weighted combined error rate (solid black line), alpha (lower grey dashed line), and beta (upper grey dotted line) for an independent *t*-test as a function of sample size when the alpha level is justified based on the goal to minimize the error rate at each sample size. The sample size corresponding to the red dot is the minimum required sample size to achieve a 5% weighted combined error rate.

## Lowering the Alpha Level to Avoid Lindley's Paradox

Formally controlling the costs of errors can be a challenge, as it requires researchers to specify the relative cost of Type 1 and Type 2 errors, prior probabilities, and the effect size of interest. Due to this complexity, researchers might be tempted to fall back on the heuristic use of an alpha level of 0.05. Fisher (1971) referred to the default alpha level of 0.05 as a "convenient convention" and one can argue suffices as a low enough threshold to make scientific claims in a scientific system where we have limited resources and value independent replications (Uygun-Tunç et al., 2021).

However, there is a well-known limitation of using a fixed alpha level that has lead statisticians to recommend choosing an alpha level as a function of the sample size. This was suggestion of a flexible decision criterion was already mentioned by the statistician

520 Harold Jeffreys in a letter he wrote to Fisher in 1934 (Wagenmakers & Ly, n.d.). Jeffreys

521 later stated more explicitly that the critical value should increase with the sample size:

522 "The results show that the probability that such a term is needed is increased or

523 decreased according as the coefficient is more or less than a certain multiple of its

524 standard error; the multiple needed, however, increases with the number of observations."

525 (Jeffreys, 1936b).

526      To understand the argument behind this recommendation, it is important to

527 distinguish between statistical inferences based on error control and inferences based on

528 likelihoods. An alpha level of 5% will limit incorrect decisions to a desired maximum (in

529 the long run, and when all test assumptions are met). However, from a likelihood

530 perspective it is possible that the observed data is much more likely when the null

531 hypothesis is true than when the alternative hypothesis is true, even when the observed

532 $p$-value is smaller than 0.05. This situation, known as Lindley's paradox, is visualized in

533 Figure 1.

534      To prevent situations where a frequentist rejects the null hypothesis based on $p <$

535 0.05, when the evidence in the test favors the null hypothesis over the alternative

536 hypothesis, it is recommended to lower the alpha level as a function of the sample size.

537 The need to do so is discussed extensively by Leamer (1978). He writes "The rule of

538 thumb quite popular now, that is, setting the significance level arbitrarily to .05, is shown

539 to be deficient in the sense that from every reasonable viewpoint the significance level

540 should be a decreasing function of sample size." The same point was already recognized

541 by Jeffreys (1939), who discusses ways to set the alpha level in the Neyman-Pearson

542 approach to statistics: "We should therefore get the best result, with any distribution of

543 alpha, by some form that makes the ratio of the critical value to the standard error

544 increase with n. It appears then that whatever the distribution may be, the use of a fixed

545 $P$ limit cannot be the one that will make the smallest number of mistakes." Similarly,

546 Good (1992) notes: "we have empirical evidence that sensible $P$ values are related to

547 weights of evidence and, therefore, that $P$ values are not entirely without merit. The real

objection to $P$ values is not that they usually are utter nonsense, but rather that they can be highly misleading, especially if the value of N is not also taken into account and is large."

Lindley's paradox emerges because in frequentist statistics the critical value of a test approaches a limit as the sample size increases (e.g., $t = 1.96$ for a two-sided $t$-test with an alpha level of 0.05). It does not emerge in Bayesian hypothesis tests because the inference criterium requires a larger test statistic as the sample size increases (Rouder, Speckman, Sun, Morey, & Iverson, 2009; Zellner, 1971). One possible inference criterium in Bayesian statistics is the Bayes factor (Kass & Raftery, 1995).

A Bayes factor contrasts the probability of the data under the competing hypotheses considered. When comparing H1 to H0 it is given by Equation 2.

$$\frac{p(data|H_1)}{p(data|H_0)} \tag{2}$$

Note that the equation shows a crucial difference between $p$-values and Bayes factors: A $p$-value depends only on the probability of the data or more extreme data under H0, whereas the Bayes factor takes both H0 and H1 into account.

A Bayes factor of 1 implies equal evidence for H0 and H1. Although any discretization inevitably results in loss of information, as a rule of thumb, Bayes factors between 3 and 10 imply moderate evidence for H1 and Bayes factors larger 10 strong evidence (Jeffreys, 1939; Lee & Wagenmakers, 2013). To prevent Lindley's paradox when using frequentist statistics one would need to adjust the alpha level in a way that the likelihood ratio (also called the Bayes factor) at the critical test statistic is not larger than 1. With such an alpha level, a significant $p$-value will always be at least as likely if H1 is true than if H0 is true, which avoids Lindley's paradox. Rouder et al. (2009) and Faulkenberry (2019) developed Bayes factors for $t$-tests and Analysis of Variance (ANOVA) which can calculate the Bayes factor from the test statistic and degrees of freedom. We developed a Shiny app that lowers the alpha level for a $t$-test or ANOVA, such that the critical value that leads researchers to reject H0 is also high enough to

574 guarantee (under the assumption of the priors) that the data provide relative evidence in

575 favor of H1.

576        There are two decisions that should be made when desiring to prevent Lindley's

577 paradox, the first about the prior, and the second about the threshold for the desired

578 evidence in favor of H1. Both Leamer (1978) and Good (1992) offer their own suggestions.

579 We rely on a unit information prior for the ANOVA and a Cauchy prior with scale 0.707

580 for $t$-tests (although the package allows users to adjust the r scale). Both of these priors

581 are relatively wide, which makes them a conservative choice when attempting to prevent

582 the Lindley's paradox. The choice for this prior is itself a 'convenient convention,' but the

583 approach extends to other priors researchers prefer, and researchers can write custom code

584 if they want to specify a different prior. A benefit of the chosen defaults for the priors is

585 that, in contrast to previous approaches that aimed to calculate a Bayes factor for every

586 $p$-value (Colquhoun, 2017, 2019), researchers do not need to specify the effect size under

587 the alternative hypothesis. This lowers the barrier of adopting this approach in situations

588 where it is difficult to specify a smallest effect size of interest or an expected effect size.

589        A second decision is the threshold of the Bayes factor used to lower the alpha level.

590 Using a Bayes factor of 1 formally prevents Lindley's paradox. It does mean that one

591 might reject the null hypothesis when the data provide just as much evidence for H1 as

592 for H0. Although it is important to note that researchers will often observe $p$-values well

593 below the critical value, and thus, in practice the evidence in the data will be in favor of

594 H1 when H0 is rejected, researchers might want to increase the threshold of the Bayes

595 factor that is used to lower the alpha level to prevent weak evidence (Jeffreys, 1939).

596 This can be achieved by setting the threshold to a larger value than 1 (e.g., BF > 3). The

597 Shiny app allows researchers to adjust the alpha level in a way that a significant $p$-value

598 will always provide moderate (BF > 3) or strong (BF > 10) evidence against the null

599 hypothesis.

600        To illustrate this approach to justifying the alpha level as a function of the sample

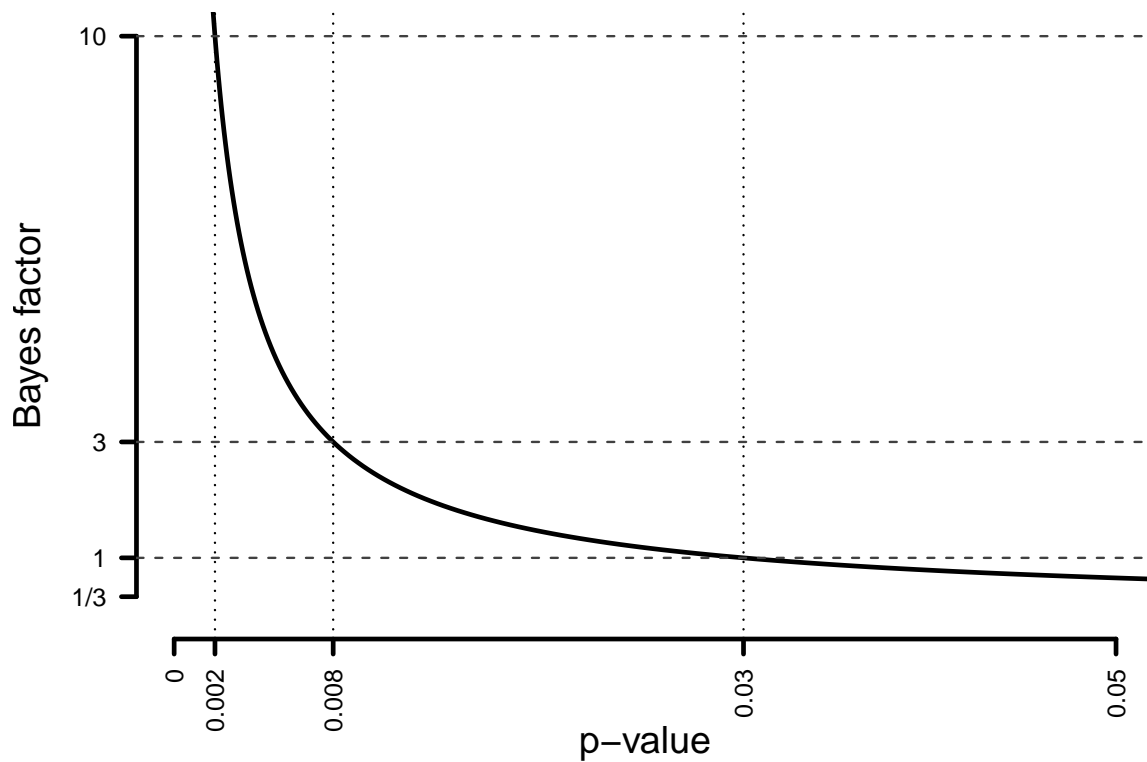601 size, imagine a researcher collected 150 observations in a within-subjects design where

*Figure 5*. Relationship between *p*-value and Bayes factor for a one-sample *t*-test with 150 participants using a Cauchy prior.

602  they aim to test a directional prediction in a dependent *t*-test. For any sample size and

603  choice of prior, a *p*-value is directly related to a Bayes factor. Figure 5 shows the

604  relationship of two-sided *p*-values and Bayes factors using a Cauchy prior with a r-scale of

605  0.707 given a sample size of 150 for a within-subjects *t*-test. To avoid Lindley's paradox,

606  the researcher would need to use an alpha level of 0.0302 for the one-sided *t*-test, given

607  the chosen prior, as this choice for an alpha level guarantees that a significant *p*-value will

608  correspond to evidence in favor of H1.

609       To give a practical example of how the alpha level can be justified to prevent

610  Lindley's paradox, we can re-examine a study by Pennycook, McPhetres, Zhang, Lu, and

611  Rand (2020) who investigated how an accuracy nudge can reduce the sharing of

612  COVID-19 misinformation on social media. In study 2 they report a main effect of

XXXXX - ADD MAX —- $F(1, 25623) = X.XX$, $p = .XXXX$. In addition, they report that in the treatment condition sharing intentions for true headlines were significantly higher than for false headlines, $F(1, 25623) = 8.89$, $p = .003$. However, given the large number of observations, which likely provide very high power for all effect sizes that would be considered large enough to be meaningful, one could have decided to reduce the alpha level so that any observed significant *p*-value can also be interpreted as evidence for the alternative hypothesis. If the authors had justified their alpha level as a function of their sample size described above, they would have set the alpha level to 0.0014. If this lower alpha level is used, the $p$-value reported by Pennycook et al. (2020) for the main effect of interest is statisticall significant at a level at which is can also be interpreted as evidence for the alternative hypothesis. However, the interaction effect is not statistically significant at this lower alpha level, indicating the interaction effect is more likely under the null hypothesis of no difference between the groups than under the assumption that the two groups differ. This is caused by the fact that the interaction effect is relatively small, but the data set is large. We can reject the null-hypothesis at the traditional alpha level, but the data do not provide evidence against the alternative hypothesis. Of course, such conclusions also depend on the prior adopted for the Bayes factor, and if the prior model had been specified such that small effects were more probable, the conclusions might have been different.

For small sample sizes it is possible to guarantee that a significant result is evidence for the alternative hypothesis using an alpha level that is higher than 0.05. It is not recommended to use the procedure outlined in this section to *increase* the alpha level above the conventional choice of an alpha level (e.g., 0.05). This approach to the justification of an alpha level assumes researchers first want to control the error rate, and as a secondary aim want to prevent Lindley's paradox by reducing the alpha level as a function of the sample size where needed. Figure 6 shows the alpha levels for different values of N for between and within subjects $t$-test. We can see that particularly for within-subjects $t$-tests the alpha level rapidly falls below 5% as the sample size increases.
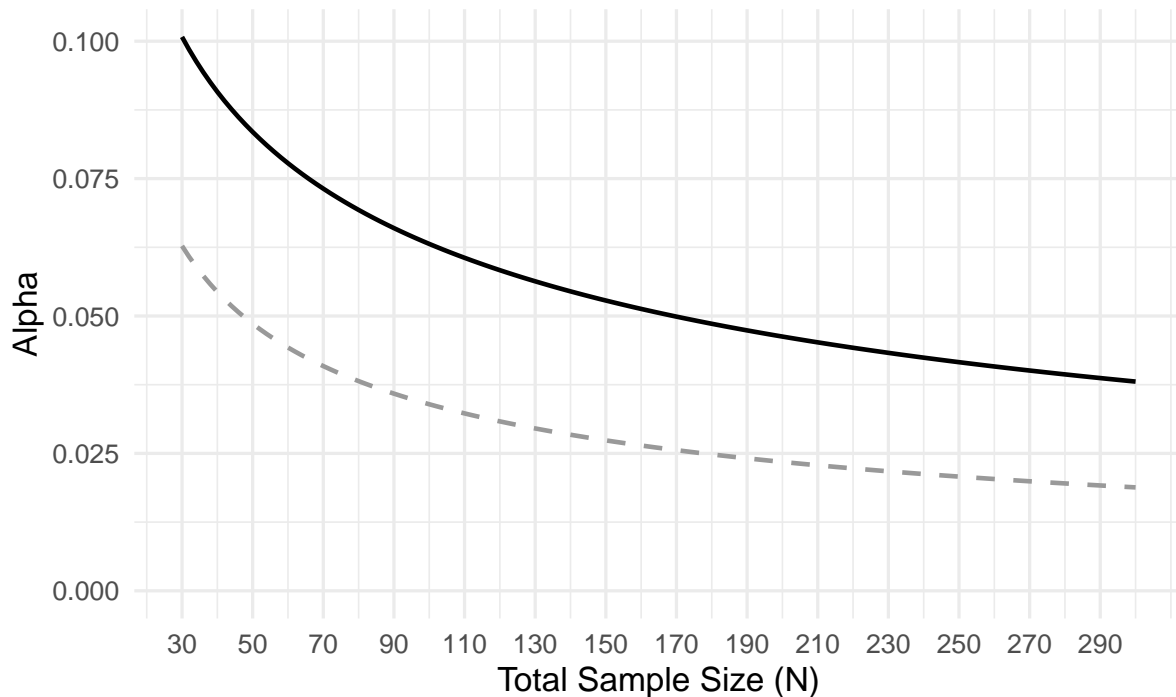
*Figure 6*. Optimal alpha level for within (grey dashed line) and between-sample (solid black line) two-sided *t*-tests.

## When to Minimize Alpha Levels and When to Avoid Lindley's Paradox

When should we minimize or balance error rates and when should we avoid Lindley's paradox? In practice, it might be most convenient to minimize or balance error rates whenever there is enough information to conduct a power analysis, and if researchers feel comfortable specifying the relative cost of Type 1 and Type 2 errors, and have a decent empirically justified estimate of prior probabilities of the null and alternative hypothesis. This is more likely for applied research, as in the case of the test of an intervention for older drivers discussed previously. When a study has direct policy implications the costs of Type 1 error (the policy being implemented although it does not work) in comparison to a Type 2 error (the policy is not implemented even though it does work) can often be assessed by means of cost-benefit analysis. It is important to note that the approach which tries to minimize or balance error rates will in practice also reduce the alpha level as a function of sample size and should therefore avoid Lindley's paradox in most applied cases (although it does not guarantee to do so). If researchers do

not feel they can specify these parameters, they can fall back on the approach to lower the alpha level as a function of the sample size to prevent Lindley's paradox. This might often be the more feasible approach in basic research.

In addition, the two approaches differ with regard to their underlying philosophy of science. The first is based on decision theoretical developments that build on a Neyman-Pearson approach and might, therefore, be more attractive to researchers whose inferential philosophy is based on statistical decision-theory. The second approach, on the other hand, offers a Bayes-Non-Bayes hybrid combining frequentist and Bayesian statistics, which might be more attractive to researchers who care about both statistical schools (Good, 1992).

## Discussion

As the choice of error rates is an important decision in any hypothesis test, authors should always be expected to justify their choice of error rates whenever they use data to make decisions about the presence or absence of an effect. As Skipper, Guenther, and Nass (1967) remark: "If, in contrast with present policy, it were conventional that editorial readers for professional journals routinely asked:"What justification is there for this level of significance? authors might be less likely to indiscriminately select an alpha level from the field of popular eligibles." It should especially become more common to lower the alpha level when analyzing large data sets or when performing meta-analyses, whenever each test has very high power to detect any effect of interest. Researchers should also consider increasing the alpha level when the combination of the effect size of interest, the sample size, the relative cost of Type 1 and Type 2 errors, and the prior probability of H1 and H0 mean this will improve the efficiency of decisions that are made.

Will departing from a fixed .05 alpha level make the scientific literature harder to interpret by having significant and non-significant mean different things in different papers? We believe this can be avoided by having each scientific claim accompanied by the alpha level under which it was made. Therefore, scientists should be required to

report their alpha levels prominently alongside any claim, usually in the abstract or even in the title of a paper.

A Shiny app is available that allows users to perform the calculations recommended in this article. It can be used to minimize or balance alpha and beta by specifying the effect size of interest and the sample size, as well as an analytic power function. The effect size should be determined as in a normal a-priori power analysis (preferably based on the smallest effect size of interest, for recommendations, see Lakens (2021)). Alternatively, researchers can lower the alpha level as a function of the sample size by specifying only their sample size. In a Neyman-Pearson approach to statistics the alpha level should be set before the data is collected. Whichever approach is used, it is strongly recommended to preregister the alpha level that researchers plan to use before the data is collected. In this preregistration, researchers should document and explain all assumptions underlying their decision for an alpha level, such as beliefs about prior probabilities or choices for the relative weight of Type 1 and Type 2 errors.

In this paper, we presented two ways of justifying alpha levels, the first based on minimizing or balancing the relative costs of errors, and the second based on avoiding Lindley's paradox. Additional approaches to justifying the alpha level have been presented, such as Bayarri, Benjamin, Berger, and Sellke (2016) , who propose to justify the alpha level based on the strength of evidence (1-beta)/alpha. We look forward to the development of additional approaches, and hope that in the future researchers will have multiple tools in their statistical toolbox to justify alpha levels.

Throughout this manuscript we have reported error rates rounded to three decimal places. Although we can compute error rates to many decimals, it is useful to remember that the error rate is a long run frequency, and in any finite number of tests (e.g., all the tests you will perform in your lifetime) the observed error rate varies somewhere around the long run error rate. The weighted combined error rate might be quite similar across a range of alpha levels, or when using different justifications (e.g., or balancing versus minimizing alpha levels in a cost-benefit approach) and small differences between alpha

levels might not be noticeable in a limited number of studies in practice. We recommend preregistering alpha levels up to three decimals, while keeping in mind there is some false precision in error rates with too many decimals.

Because of the strong norms to use a 5% error rate when designing studies, there are relatively few examples of researchers who attempt to justify the use of a different alpha level. Within specific research lines researchers will need to start to develop best practices to decide how to weigh the relative cost of Type 1 and Type 2 errors, or quantify beliefs about prior probabilities. It might be a challenge to get started, but the two approaches illustrated here provide one way to move beyond the mindless use of a 5% alpha level, and make more informative decisions when we test hypotheses.

## Funding

## Supplemental material

All code used to create this manuscript is provided at https://github.com/Lakens/justify_alpha_in_practice. Information about the JustifyAlpha R package and Shiny app is available at https://lakens.github.io/JustifyAlpha/index.html.

## Prior versions

A preprint of this article is available at https://doi.org/10.31234/osf.io/ts4r6.

## References

Bakan, D. (1966). The test of significance in psychological research. *Psychological Bulletin*, *66*(6), 423–437. https://doi.org/10.1037/h0020412

Bartlett, P. L., Jordan, M. I., & Mcauliffe, J. D. (1957). Comment on d. V. Lindley's statistical paradox. *Biometrika*, *44*, 533–534.

Bayarri, M., Benjamin, D. J., Berger, J. O., & Sellke, T. M. (2016). Rejection odds and rejection ratios: A proposal for statistical practice in testing hypotheses. *Journal of Mathematical Psychology*, *72*, 90–103. Retrieved from https://doi.org/10.1016/j.jmp.2015.12.007

Camerer, C. F., Dreber, A., Forsell, E., Ho, T.-H., Huber, J., Johannesson, M., ... Wu, H. (2016). Evaluating replicability of laboratory experiments in economics. *Science*, *351*(6280), 1433–1436. https://doi.org/10.1126/science.aaf0918

Clemen, R. T. (1997). *Making Hard Decisions: An Introduction to Decision Analysis* (2 edition). Belmont, Calif: Duxbury.

Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (2nd ed). Hillsdale, N.J: L. Erlbaum Associates.

Colquhoun, D. (2017). The reproducibility of research and the misinterpretation of p-values. *Royal Society Open Science*, *4*(12), 171085. https://doi.org/10.1098/rsos.171085

Colquhoun, D. (2019). The false positive risk: A proposal concerning what to do about p-values. *The American Statistician*, *73*(sup1), 192–201. https://doi.org/10.1080/00031305.2018.1529622

Cornfield, J. (1969). The bayesian outlook and its application. *Biometrics*, 617–657.

Cousins, R. D. (2017). The Jeffreys–Lindley paradox and discovery criteria in high energy physics. *Synthese*, *194*(2), 395–432. https://doi.org/10.1007/s11229-014-0525-z

Cowles, M., & Davis, C. (1982). On the origins of the. 05 level of statistical significance. *American Psychologist*, *37*(5), 553--558. https://doi.org/10.1037/0003-066X.37.5.553

Cumming, G. (2008). Replication and *p* Intervals: *P* Values Predict the Future Only
    Vaguely, but Confidence Intervals Do Much Better. *Perspectives on Psychological
    Science*, *3*(4), 286–300. https://doi.org/10.1111/j.1745-6924.2008.00079.x

DeGroot, M. (1975). *Probability and statistics. massachusetts.* Addison-Wesley
    Publishing Company, InC.

Douglas, H. E. (2000). Inductive risk and values in science. *Philosophy of Science*, *67*(4),
    559–579. https://doi.org/10.1086/392855

Edwards, W., Miles Jr., R. F., & Winterfeldt, D. von (Eds.). (2007). *Advances in
    decision analysis: From foundations to applications.* Cambridge: Cambridge
    University Press. https://doi.org/10.1017/CBO9780511611308

Erdfelder, E., Faul, F., & Buchner, A. (1996). GPOWER: A general power analysis
    program. *Behavior Research Methods, Instruments, & Computers*, *28*(1), 1–11.
    https://doi.org/10.3758/BF03203630

Faulkenberry, T. J. (2019). Estimating evidential value from analysis of variance
    summaries: A comment on ly et al.(2018). *Advances in Methods and Practices in
    Psychological Science*, *2*(4), 406–409. https://doi.org/10.1177/2515245919872960

Fiedler, K., Kutzner, F., & Krueger, J. I. (2012). The long way from $\alpha$-error control to
    validity proper: Problems with a short-sighted false-positive debate. *Perspectives on
    Psychological Science*, *7*(661–669), 661–669.
    https://doi.org/10.1177/1745691612462587

Field, S. A., Tyre, A. J., Jonzén, N., Rhodes, J. R., & Possingham, H. P. (2004).
    Minimizing the cost of environmental management decisions by optimizing statistical
    thresholds. *Ecology Letters*, *7*(8), 669–675. Retrieved from
    10.1111/j.1461-0248.2004.00625.x

Fisher, R. A. (1926). Introduction to "The arrangement of field experiments." *Journal of
    the Ministry of Agriculture*, *33*, 503–513.

Fisher, Ronald A. (1971). *The Design of Experiments* (9 edition). New York: Macmillan
    Pub Co.

Gigerenzer, G. (2018). Statistical Rituals: The Replication Delusion and How We Got

There. *Advances in Methods and Practices in Psychological Science*,

2515245918771329. https://doi.org/10.1177/2515245918771329

Good, I. J. (1992). The Bayes-Non-Bayes Compromise: A Brief Review. *Journal of the

American Statistical Association*, *87*(419), 597. https://doi.org/10.2307/2290192

Harford, T. (2014). Big data: A big mistake? *Significance*, *11*(5), 14–19. Retrieved from

https://doi.org/10.1111/j.1740-9713.2014.00778.x

Jeffreys, H. (1935). Some tests of significance, treated by the theory of probability.

*Mathematical proceedings of the cambridge philosophical society*, *31*, 203–222.

Cambridge University Press.

Jeffreys, H. (1936a). Further significance tests. *Mathematical proceedings of the

cambridge philosophical society*, *32*, 416–445. Cambridge University Press.

Jeffreys, H. (1936b). On some criticisms of the theory of probability. *The London,

Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, *22*(146),

337–359.

Jeffreys, H. (1939). *Theory of probability* (1st ed.). Oxford, UK: Oxford University Press.

Kass, R. E., & Raftery, A. E. (1995). Bayes factors. *Journal of the American Statistical

Association*, *90*(430), 773–795.

Kennedy-Shaffer, L. (2019). Before p < 0.05 to beyond p < 0.05: Using history to

contextualize p-values and significance testing. *The American Statistician*, *73*(sup1),

82–90. https://doi.org/10.1080/00031305.2018.1537891

Kim, J. H. (2020). Decision-theoretic hypothesis testing: A primer with r package

OptSig. *The American Statistician*, *74*(4), 370–379.

https://doi.org/https://doi.org/10.1080/00031305.2020.1750484

Kim, J. H., & Choi, I. (2021). Choosing the level of significance: A decision-theoretic

approach. *Abacus*, *57*(1), 27–71. https://doi.org/10.1111/abac.12172

Lakens, D. (2021). *Sample Size Justification*. PsyArXiv.

https://doi.org/10.31234/osf.io/9d3yf

Leamer, E. E. (1978). *Specification Searches: Ad Hoc Inference with Nonexperimental Data* (1 edition). New York: Wiley.

Lee, M. D., & Wagenmakers, E.-J. (2013). *Bayesian cognitive modeling: A practical course.* Cambridge University Press.

Lin, M., Lucas Jr, H. C., & Shmueli, G. (2013). Research commentary—too big to fail: Large samples and the p-value problem. *Information Systems Research, 24*(4), 906–917.

Lindley, D. V. (1953). Statistical inference. *Journal of the Royal Statistical Society: Series B (Methodological), 15*(1), 30–65.

Lindley, D. V. (1957). A statistical paradox. *Biometrika, 44*(1/2), 187–192.

Lindsay, D. S. (2015). Replication in Psychological Science. *Psychological Science, 26*(12), 1827–1832. https://doi.org/10.1177/0956797615616374

Miller, J., & Ulrich, R. (2019). The quest for an optimal alpha. *PLOS ONE, 14*(1), e0208631. https://doi.org/10.1371/journal.pone.0208631

Mudge, J. F., Baker, L. F., Edge, C. B., & Houlahan, J. E. (2012). Setting an optimal $\alpha$ That minimizes errors in null hypothesis significance tests. *PLOS ONE, 7*(2), e32734. https://doi.org/10.1371/journal.pone.0032734

Neyman, Jerzy. (1957). "Inductive Behavior" as a Basic Concept of Philosophy of Science. *Revue de l'Institut International de Statistique / Review of the International Statistical Institute, 25*(1/3), 7–22. https://doi.org/10.2307/1401671

Neyman, J., & Pearson, E. S. (1933). On the problem of the most efficient tests of statistical hypotheses. *Philosophical Transactions of the Royal Society of London A: Mathematical, Physical and Engineering Sciences, 231*(694-706), 289–337. https://doi.org/10.1098/rsta.1933.0009

Open Science Collaboration. (2015). Estimating the reproducibility of psychological science. *Science, 349*(6251), aac4716–aac4716. https://doi.org/10.1126/science.aac4716

Pennycook, G., McPhetres, J., Zhang, Y., Lu, J. G., & Rand, D. G. (2020). Fighting COVID-19 misinformation on social media: Experimental evidence for a scalable accuracy-nudge intervention. *Psychological Science*, *31*(7), 770–780. Retrieved from https://doi.org/10.1177/0956797620939054

Pericchi, L., & Pereira, C. (2016). Adaptative significance levels using optimal decision rules: Balancing by weighting the error probabilities. *Brazilian Journal of Probability and Statistics*, *30*(1), 70–90.

Rouder, J. N., Speckman, P. L., Sun, D., Morey, R. D., & Iverson, G. (2009). Bayesian t tests for accepting and rejecting the null hypothesis. *Psychonomic Bulletin & Review*, *16*(2), 225–237. https://doi.org/10.3758/PBR.16.2.225

Simmons, J. P., Nelson, L. D., & Simonsohn, U. (2011). False-positive psychology: Undisclosed flexibility in data collection and analysis allows presenting anything as significant. *Psychological Science*, *22*(11), 1359–1366. https://doi.org/https://doi.org/10.1177%2F0956797611417632

Skipper, J. K., Guenther, A. L., & Nass, G. (1967). The sacredness of .05: A note concerning the uses of statistical levels of significance in social science. *The American Sociologist*, *2*(1), 16–18.

Spanos, A. (2013). Who should be afraid of the Jeffreys-Lindley paradox? *Philosophy of Science*, *80*(1), 73–93. https://doi.org/10.1086/668875

Uygun-Tunç, D., Tunç, M. N., & Lakens, D. (2021). *The epistemic and pragmatic function of dichotomous Claims based on statistical hypothesis tests.* https://doi.org/10.31234/osf.io/af9by

Viamonte, S. M., Ball, K. K., & Kilgore, M. (2006). A cost-benefit analysis of risk-reduction strategies targeted at older drivers. *Traffic Injury Prevention*, *7*(4), 352–359. Retrieved from https://doi.org/10.1080/15389580600791362

Wagenmakers, E. J., & Ly, A. (n.d.). History and nature of the jeffreys-lindley paradox. *Manuscript in Preparation.*

Wald, A. (1949). Statistical decision functions. *The Annals of Mathematical Statistics*, 165–205. Retrieved from https://doi.org/10.1214/aoms/1177730030

869 Wasserstein, R. L., & Lazar, N. A. (2016). The ASA statement on p-values: Context,

870     process, and purpose. *The American Statistician, 70*(2), 129–133.

871     https://doi.org/10.1080/00031305.2016.1154108

872 Winer, B. J. (1962). *Statistical principles in experimental design.* New York :

873     McGraw-Hill.

874 Zellner, A. (1971). *An introduction to Bayesian inference in econometrics.* New York:

875     Wiley.