2

3

4

**No Estimation without Inference:**

**A Response to the International Society of Physiotherapy Journal Editors**

Keith Lohse, PhD[1]

[1] Physical Therapy and Neurology, Washington University School of Medicine, Saint Louis, MO

35    Recently, Elkins et al.[1] (hereafter referred to as "the Editorial") published an editorial on behalf of

36    the International Society of Physiotherapy Journal Editors (ISPJE), recommending that researchers stop

37    using null hypothesis significance tests and adopt "estimation methods". Further, the editorial warns that

38    this is not merely an idea to consider, but a coming policy of journals: "the [ISPJE] will be expecting

39    manuscripts to use estimation methods *instead* of null hypothesis statistical tests" (emphasis added).

40    However, the Editorial is deeply flawed in its statistical reasoning. If the proposed policies were adopted,

41    they could damage the statistical literacy and scientific integrity of the field.

42    I detail each of my critiques below, but in short the Editorial: (1) fails to adequately grapple with

43    the inherent connection between hypothesis testing and estimation as methods of statistical inference, (2)

44    presents several misleading arguments about the flaws of statistical significance tests, and (3) presents an

45    alternative that is, in itself, a form of significance testing – the minimal effects test[2] (but the alternative

46    does this implicitly and muddles two-sided and one-sided hypothesis testing). Finally, I end with a short

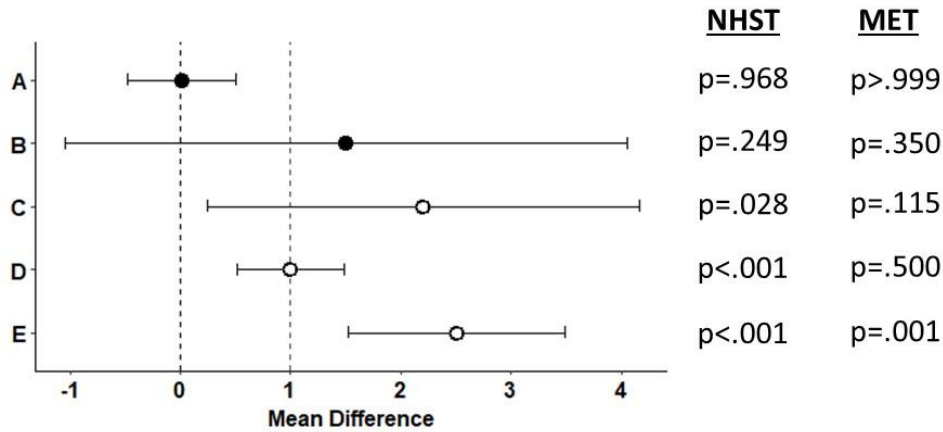47    list of more urgent problems that the ISPJE could work to address.

48    I commend the Editorial for encouraging researchers to think deeply about the statistical tools

49    available to them, to consider "practical significance" as well as "statistical significance", and for

50    bringing important methodological discussions to the forefront of physical therapy research. However, the

51    central argument of the Editorial is illogical and I worry what coming policy changes might mean for how

52    authors interpret their data. I think the antidote to researchers making faulty decisions is not to ban p-

53    values, but to improve education. A rising tide lifts all boats, and if the baseline statistical literacy in our

54    field were higher, authors would make fewer mistakes, reviewers would be more apt to catch remaining

55    mistakes, and readers would be better equipped to make their own conclusions given the available data.

56    Editors then need to hold the line and ensure rigorous review, not ban valid statistical tools.

**Hypothesis Testing and Estimation are Inescapably Intertwined**

57

58        The Editorial presents hypothesis testing and estimation as two distinct methodological

59   approaches. However, these approaches are two sides of the same coin, as illustrated by a simple example

60   in Figure 1. When a 95% confidence interval excludes the null value, then one can reject the null

61   hypothesis at p<.05. This is because hypothesis tests and confidence intervals are based on the same

62   underlying mathematics: e.g., how big is the observed effect relative to the variability we would expect

63   due to sampling? Although typically we think of the null-hypothesis as an assumption of "no effect", the

64   null hypothesis can assume zero or non-zero effects. So, as shown in the figure, we can ascertain the

65   probability of observing the data we did, assuming a null value of 0 or a null value of 1.

66        Hypothesis testing and estimation cannot be fully disentangled: estimation (frequentist or

67   Bayesian) asks about *plausible values* of the parameter in the population, hypothesis testing asks about

68   the plausibility of *a specific parameter value*. These are both inferences, because we are inferring

69   something about the population based on the data in our sample. In the frequentist paradigm, uncertainty

70   in the inference is accounted for with long-run error control; e.g., setting the Type 1 Error rate, α=0.05.

71   We can see this when running simulations as shown in Figure 1A-E: any confidence interval that does not

72   contain zero also has p<0.05, for the null hypothesis significance test (NHST).

73        The 95% confidence interval shows values in the population that are *compatible* with what we

74   observed in the sample.[3] That is, if you move outside of the confidence interval, any of those parameter

75   values (the "true" mean differences; Δ's) would be statistically different from the mean difference

76   observed in the sample ($\overline{x_d}$) at the *p*<0.05 level. Inside of the confidence interval, none of those parameter

77   values would be statistically different (p>0.05) from the observed mean difference. Recall that the p-value

78   is the probability of observing data as extreme or more extreme, assuming that the null hypothesis is true,

79   formally written as $p(\geq \overline{x_d}|H_0)$.

|  | NHST | MET |
|---|---|---|
| A | p=.968 | p>.999 |
| B | p=.249 | p=.350 |
| C | p=.028 | p=.115 |
| D | p<.001 | p=.500 |
| E | p<.001 | p=.001 |

80

**Figure 1.** 95% confidence intervals and corresponding p-values for testing $H_0: \Delta = 0$ (NHST, null hypothesis significance testing) and $H_0: \Delta \leq 1$ (MET, a one-sided minimal effects test).

Typically, the null hypothesis significance test (NHST) assumes that the true value in the population is 0 (i.e., $H_0: \Delta = 0$). The further the sample mean difference is away from 0, the lower the probability of observing that sample mean, if the null hypothesis were true. Importantly, the Editorial does not address the fact that we can set $H_0$ to be any value. For instance, rather than setting $H_0: \Delta = 0$ (sometimes referred to as the "nil-hypothesis")[4], we can set $H_0$ equal to any clinically meaningful value of interest. This is referred to as a minimal effects test (or minimum effect test, MET[2,5]). For the sake of argument, let's say this value is 1 in Figure 1. Comparing the confidence intervals to the new null value, you can see that any confidence intervals that only contain values larger than 1 also have a p<0.05 for the minimal effects test (i.e., Figure 1E).[A] Thus, we have both an inference about a specific hypothesis and an estimate in both the NHST and the MET[B], but the hypothesis test and the estimate are complementary and connected.

[A] METs are typically directional, using one-sided hypothesis tests (e.g., $H_0: \leq 1$) whereas NHSTs are often non-directional, using two-sided hypothesis tests (e.g., $H_0: = 0$). Thus, although the confidence interval for Figure 1A does not contain the null value of 1, the whole of the confidence interval is below 1, thus yielding a non-significant minimal effects test.

[B] For convenience, I am referring to NHST and MET as separate tests. However, it is more accurate to think of the MET as type of NHST where you have a one-sided test of a non-zero null value. I use the different terms because readers are likely more familiar with the term NHST when referring to the specific case of $H_0 = 0$.[4]

**Misleading Arguments about flaws with Significance Tests**

The Editorial bases many arguments on a previous list of perceived problems from Herbert (2019).[6] The Herbert paper is in itself an editorial that presents informed arguments, but is not an objective demonstration of any mathematical facts. So, reinforcing the Editorial's list through a citation to Herbert does not provide an evidentiary foundation: it is layering opinion on top of opinion. Second, each of the five "problems" outlined by the Editorial is either not really a problem inherent to p-values or the problem is a true but misleading statement. I address each problem from the Editorial (in quotes) below:

**1. *"A p-value is not the probability that a hypothesis is (or is not) true."*** – This is correct, but it does not follow that this makes p-values, or even statistical significance tests, unhelpful or uninformative. Knowing that the observed data are incompatible with some null value is a crucial step for many research questions. For instance, hypothesis testing in early phase research can help us make decisions about where to direct our resources, starting us down the road of replication and ultimately determining the efficacy and effectiveness of an intervention.

**2. *"A p-value does not constitute evidence"*** – This is an oversimplification and misleading. The Editorial is correct that a single p-value is not strictly speaking "evidence" and cannot tell us about the probability of the null hypothesis being true. However, p-values are still useful tools for making decisions.

Technical definitions of evidence can get a bit complicated and are debated.[7–9] However, I would invite readers to consider a simple example of absolute probability versus relative probability. If I find that eating green jelly beans reduces post-surgical recovery time for the ACL by 10% relative to controls with $p<0.05$, then the most likely explanation is still that jelly beans have no effect on recovery and what I observed was chance fluctuation. That is, the null hypothesis is still the most likely explanation even though $p$ was <0.05, because the baseline probability of "jelly bean efficacy" is very low and false positives occur 5% of the time when $\alpha = 0.05$. Thus, the p-value is not in itself a measure of evidence, because I would need additional *outside information* in order to change (or not change) my beliefs. As
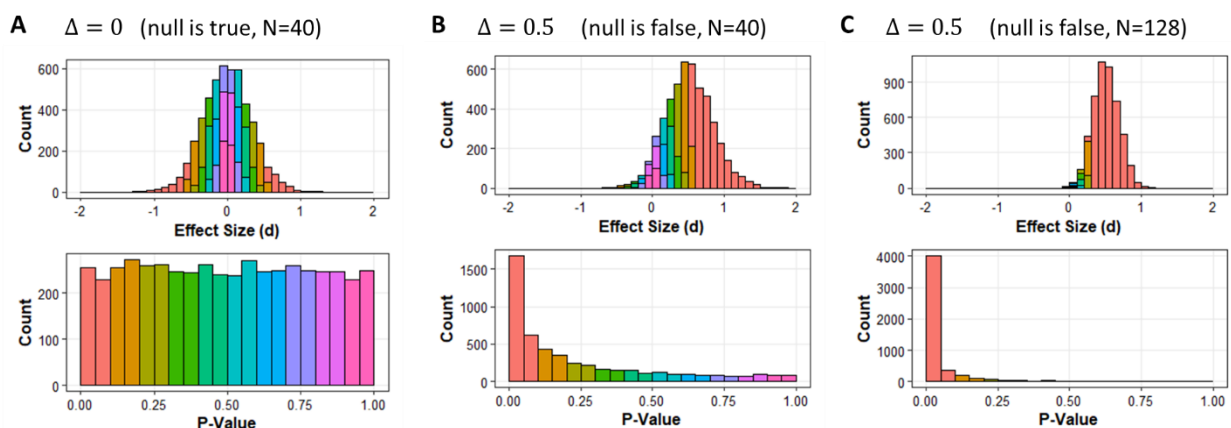
119     Goodman and Royall[9] write "The p-value is not adequate for inference because ***the measurement of***

120     ***evidence*** requires are least three components: the observations, and two competing explanations for how

121     they were produced" (p. 1569; emphasis added).

122             Some researchers might think of the p-value as evidence against the null specifically, without the

123     need for comparison to a given alternative. But the p-value is calculated assuming that the null is true, so

124     again the Editorial is correct that we cannot simply flip the question around, assume the data, and get the

125     likelihood of the null being true, i.e., $p(\bar{x}_d|H_0) \neq p(H_0|\bar{x}_d)$. To estimate the likelihood of the null

126     hypothesis being true, we would need Bayesian statistics in which we formalize some *prior* probability

127     about the null hypothesis.[9] If we have a strong enough prior probability that the null is true, then the

128     current data in the sample may not lead us to change our beliefs based on the *posterior* distribution, no

129     matter how small the p-value. This was the case in my jelly bean example, where $p<0.05$ still did not

130     shake my belief in the null hypothesis. For any given prior distribution, however, there is a smaller

131     *likelihood* of observing highly discrepant effects (e.g., $|\bar{x}_d|>>0$), leading to a smaller relative probability

132     of 0 in the posterior distribution compared to the prior distribution.[C] Updating the probability of 0 in the

133     posterior distribution reflects rational decision making in daily life. For instance, the first time I find jelly

134     beans reduce recovery time with $p<0.05$, I might rightly ignore that as a false positive. The fifth time I

135     find jelly beans reduce recovery time with $p<0.05$, I should take a long hard look at the ingredients and

136     maybe my study procedures; as $p<0.05$ is not always a sign that the null is wrong, but that some other

137     assumption has been violated.

138             Still, the p-value does not need to be a measure of evidence for it to be useful. Critically, small p-

139     values are *relatively* less likely to be observed when the null hypothesis is true compared to when an

140     alternative hypothesis is true. Thus, in a practical sense, a p-values can help us make decisions about what

---

[C] For a humorous demonstration see: https://xkcd.com/1132/ ; for a more quantitative visualization of the relationship between priors, p-values, and posteriors see: https://rpsychologist.com/d3/bayes/. More technically, the posterior (the updated probability density function after we've seen the evidence) is proportional to the prior (our expectation before we saw the evidence) multiplied by the likelihood (which is the probability of the current evidence given the hypothesis): $posterior \propto likelihood \cdot prior$.

141     effects to study, assuming that we are testing at least some real effects. As shown in Figure 2A, p-values

142     have a uniform distribution under the null hypothesis, with 5% of p-values necessarily below 0.05.

143     However, if the null is not true, then we will see a shift in the distribution of p-values, with small p-values

144     becoming more common. An example of this is shown in Figure 2B, where the null is false and 34% of p-

145     values are below 0.05. However, correctly rejecting the null hypothesis only 34% of the time is not ideal,

146     so consider Figure 2C, where I have now tripled the sample size and 80% of p-values are below 0.05.

147     That is, with 64 people per group, we now have 80% statistical power to detect a $\Delta = 0.5$.



148

**Figure 2.** P-values <0.05 are more likely to occur when the null is false, and critically will only occur 5%
of the time when the null is true. Plots show simulated experiments (k=5,000, σ=1 for all populations) in
which the means of two independent groups are compared using a t-test. In Panel A, the null hypothesis is
true and the true difference between population means is 0. In Panel B, the null hypothesis is false and the
true difference between population means is 0.5. In Panel C, the null-hypothesis is still false, but I have
increased the sample size from 40 to 128, yielding 80% of p-values <0.05 (i.e., 80% statistical power).
Quantiles are color coded with respect to their p-values and effects sizes are given as Cohen's d.

156

157         This is where the concept of a decision is important to distinguish from the term "evidence".[9]

158     Without knowing the actual *evidence* against the null-hypothesis, if I *decide* to reject the null when

159     p<0.05, then I will only be wrong 5% of the time (i.e., the Type 1 error rate). Similarly, if I have 80%

160     statistical power and a reasonable estimate for the smallest effect size of interest, then I only have a 20%

161     chance of missing an effect of that size (i.e., the Type 2 error rate). Mathematically, these probabilities are

162     robust if we accept the null-hypothesis as true and make minimal other assumptions, which is very helpful

163  when limited outside information is available. See Goodman quoting Neyman and Pearson about

164  hypothesis testing, "Without hoping to know whether each separate hypothesis is true or false, we may

165  search for rules to govern our behaviour with regard to them, in following which we insure that, in the

166  long run of experience, we shall not often be wrong."[10]

167      So, p-values are not a measure of evidence, but they are useful tools for helping us make the

168  correct decision. If we want a proper measure of evidence for one hypothesis versus another, then we can

169  do more work, but we also need to make more assumptions and/or bring in outside information. This can

170  be both a feature and bug of *using* hypothesis tests. We can control long run error rates with minimal

171  information, but if we do that so habitually that we forget other information is available, then that is on us

172  not the p-value.

173  **3. *"Statistically significant findings are not very replicable."*** – This is misleading. First, it is difficult to

174  precisely define replication,[11,12] but if we think about "being replicable" as the probability that a

175  statistically significant result represents a real, non-zero effect then we would expect more statistically

176  significant findings to "replicate" provided that hypothesis tests have adequate statistical power,

177  researchers have not engaged in p-hacking, there is not selective reporting of results, etc. Thus, not all

178  statistically significant findings will replicate,[13] but statistically significant findings in well-designed

179  studies are more likely to replicate.[14–16] Second and by any definition, threats to replicability are also

180  going to affect confidence intervals (the Editorial's proposed solution) as much as they affect p-values,

181  because, again, the p-value is intrinsically linked to the confidence interval. Thus, the Editorial is correct

182  in a practical sense: many statistically significant findings in the current literature do not replicate.

183  However, a lack of replication is the fault of poor study design and questionable research practices, not

184  the use of hypothesis tests as a method of inference.

185  **4. *"In most clinical trials, the null hypothesis must be false."*** – This is arguably true but very

186  misleading. It is true that real treatment effects are unlikely to be precisely 0 (e.g., they might be +0.001),

187  but it begs the question: do we really care if the true effect is 0 or 0.001? And will we ever have the

188  statistical precision to discern that difference? All measurement has some error, so I would argue that

189  many effects are functionally 0 even if the (unknowable) true value is not actually zero. But, in a strict

190  mathematical sense I will concede the Editorial is correct, if we accept a hyper-precise definition, the

191  null-hypothesis of $H_0: \Delta = 0.\overline{00}$ will usually be false. However, if we accept that definition, then all

192  point-estimates are false and no value will ever be precisely the minimum clinically important difference

193  either, which is the Editorial's proposed point-estimate in their alternative.

194         In response[17] to an independent critique by Lakens[18], this hyper-precise definition does seem to

195  be the argument that the editorial is making.[D] They claim, "The assertion that the null hypothesis is false

196  in most clinical trials does not require empirical evidence, because it is self-evidently true" and "The null

197  hypothesis may often be approximately true, but it is rarely if ever exactly true". The Editorial seems to

198  miss the point that the null is a useful *model*: testing against 0 is still useful for things that are

199  approximately 0. As an analogy, I have successfully gotten many places using maps, but none of those

200  maps was a photo-realistic version of reality.

201         Scientists are often working on the frontiers of human knowledge; this is costly work where we

202  need to explore a lot of different ideas and many them do not pan out. That is, many tested "effects" are

203  functionally zero.[14] So, simply because a point estimate of precisely 0 is unlikely to be true does not mean

204  that it is unhelpful to ask. It should be a very low bar to show that your clinical treatment has a non-zero

205  effect! Further, the Editorial is specifically critiquing this "nil" hypothesis (i.e., $H_0 = 0$), when we could

206  hypothesize any value, or avoid the point-null entirely with a one-sided test (i.e., $H_0 \leq 0$).[2,5] So, if

207  assuming $H_0 = 0$ is not desirable, we can set that null value to be anything we want (i.e., $H_0: \Delta \leq 0.4$

---

[D] I was very excited to see the Lakens commentary[19] and others[20], and even more excited to see we all largely agree. Interestingly, however, I only became aware of these commentaries after writing my own because I did not see the editorial until it was re-published in *Physical Therapy*[1] in June, 2022, whereas my more astute colleagues responded to the original publication in the *Journal of Physiotherapy*[21], in January 2022. The editorial has been re-published in four different journals to date. While I can appreciate trying to spread one's message, this creates confusion.

208 m/s for improvement in gait speed, $H_0: \Delta \leq 30\%$ change on a pain scale, or $H_0: \Delta \leq 1$ in the hypothetical

209 example in Figure 1).

210 **5. "*Researchers need information about the size of effects*."** – This is a true statement, but it is not a

211 problem with p-values nor null hypothesis significance tests. To my knowledge, no statistician has ever

212 recommended that applied researchers ignore measures of effect size (either raw or standardized).

213 Estimates of effect size are integral to any results section. I would even take this one step further and

214 encourage authors to share their data whenever possible[22], enabling other researchers to calculate their

215 own effect sizes as there can be limitations with and confusion about standardized effects sizes, and there

216 is no one-size-fits-all solution to effect sizes[23–25].

217 **The Editorial's "Alternative" is a Hypothesis Test – The Minimal Effects Test**

218 After detailing the potential problems with the NHST, the Editorial proposes an alternative

219 solution in which they encourage authors to compare their 95% confidence interval to some minimum

220 clinically meaningful value (which I will write as $\delta$).[E] Estimation is a good practice and I would

221 encourage researchers to report 95% confidence intervals and interpret their upper and lower limits in

222 context, when appropriate. However, what the Editorial is suggesting is effectively an MET where

223 $H_0: \Delta \leq \delta$. That is, if the test is to see if the 95% confidence interval does not contain $\delta$, then that is

224 mathematically equivalent to an MET assuming $H_0: \Delta \leq \delta$ and finding $p < 0.025$. Note p<0.025, not

225 p<0.05, because most METs are one-sided hypothesis tests whereas confidence intervals are two sided

226 (see Figure 1 and Footnote A). After heavily critiquing hypothesis testing as a method of inference, the

227 Editorial ends up effectively proposing a hypothesis test. This is clearly an illogical proposition.

228 I want to emphasize that it is valid for the Editorial to recommend that authors consider their 95%

229 confidence interval relative to some clinically meaningful value. However, this is not an "alternative" to

---

[E] I caution that it is difficult to find a single measure of $\delta$; it changes as a function of the study population, the study context, and has its own uncertainty due to sampling error.[20,26]

230      conducting a null hypothesis significance test, it is in fact mathematically identical to conducting a null

231      hypothesis test with a carefully chosen null hypothesis. Both are valid.

232      I would add, however, that there are also advantages to explicitly framing this as a hypothesis test

233      rather than the informal interpretation of a confidence interval. First, it encourages researchers to

234      explicitly commit to a specific $\delta$ while the study is being designed, rather than simply obtaining an

235      estimate of the effect and then comparing it to candidate $\delta$'s post hoc. Second, it requires researchers to

236      think carefully about the direction of the test and the desired $\alpha$-level, whereas simply invoking a 95%

237      confidence interval implicitly uses a two-tailed test and $\alpha = 0.05$, which may not be best suited to the

238      research question.

239      Finally, it is also important to stress that history provides us with several examples of how

240      authors will view their data through rose-tinted glasses when quantitative statistical safeguards are

241      removed. For instance, when *Basic and Applied Social Psychology* banned p-values, authors were found

242      to overstate their conclusions well beyond what would have been considered if "statistical significance"

243      had been a benchmark.[27] In sport and exercise science, "magnitude-based inference" was leveraged as a

244      niche method that allowed authors to interpret differences as meaningful when they had very little

245      statistical support (e.g., *p*'s >0.25).[28–30] Statistical significance in an NHST does not necessarily need to be

246      the benchmark nor 0.05 the default value[31–34], but it is always important to have a statistically sound

247      framework for dealing with uncertainty.

248      **Virtues of Hypothesis Testing**

249      One of the great virtues of null hypothesis significance testing is Type I error control while

250      making minimal assumptions about the nature of the data or the world at large. If we set $\alpha = 0.05$, then

251      we can be confident we will only get data greater than or equal to what we observed 5% of the time when

252      the null is true. Importantly, this works for a wide range of statistics and types of tests, including *F*- and

253      $\chi^2$-statistics that have multiple degrees of freedom from models asking questions about multiple effects

254  simultaneously. For instance, in a randomized controlled trial with three arms, I could conduct an

255  omnibus $F$-test and obtain a $p$-value to see if there is any evidence of a difference between groups overall,

256  before conducting additional post-hoc tests to compare specific groups. This situation is not covered by

257  the Editorial and the Editorial's confidence interval alternative is not easily applied here, although one

258  could plausibly adjust the width of the confidence intervals to control for multiple comparisons.

259  **Bigger Threats to Statistical Integrity**

260      Misinterpretation and misuse of p-values are threats to statistical integrity. However, questionable

261  research practices such as p-hacking, sub-group analyses, flexible stopping rules, selective exclusion of

262  outliers, selective reporting, and hypothesizing after results are known are much larger threats.[35–39]

263  Furthermore, these questionable research practices have consistently negative consequences regardless of

264  the method of inference. For instance, although the term "$p$-hacking" connotes the NHST, these

265  questionable research practices pose an equal threat to confidence intervals because again confidence

266  intervals and p-values are based on the same underlying mathematics. Similarly, switching to a fully

267  Bayesian method of analysis is not an antidote for poor study design, small samples, and questionable

268  research practices. As others have argued,[40,41] p-values get a disproportionate amount of attention in

269  popular discussions of research methodology. I encourage the ISPJE to instead focus their attention on

270  methods for improving data/code sharing, transparency, and replicability through tools like

271  preregistration, results-blind review, registered reports, or even "data papers" whose primary function is

272  to report a study and archive the data, without drawing inferences from limited samples.

273      It is entirely valid to say that p-values are often mis-used and mis-interpreted, and "statistical

274  significance" may not ultimately be the best term for applied researchers to use.[42] However, it is incorrect

275  to present these human errors as inherent flaws in hypothesis testing. For instance, if someone mis-

276  interprets $p>0.05$ as evidence of "no difference", then I would argue the correct action is to teach them

277  about equivalence tests and non-inferiority designs, not ban p-values. Similarly, there are times when

278  Bayesian inference is what authors are really interested in (e.g., what is the probability that the null is

279    true, given the evidence?), and in those cases Bayesian inference can and should be used. However,

280    Bayesian analysis is not a panacea and needs to be used thoughtfully like any statistical tool. So, although

281    a simple heuristic of $p<0.05$ may well be overused as "the" test in physical therapy research, frequentist

282    hypothesis tests are still valid and useful tools for physical therapy researchers. Moreover, the scientific

283    integrity of the field has much larger concerns, and both p-values and confidence intervals will be

284    corrupted by p-hacking, under-powered subgroup analyses, surrogate outcomes, and other questionable

285    research practices.

286        In conclusion, I agree with the Editorial on the importance of reporting effect sizes and

287    interpreting them in context. However, the Editorial makes numerous statistical faux pas that could harm

288    the statistical literacy in our field, if readers take them at face value, and harm the scientific integrity of

289    our field, if put into editorial practice.

290

300 **References**

301   1.   Elkins, M. R. *et al.* Statistical inference through estimation: recommendations from the International

302        Society of Physiotherapy Journal Editors. *Phys. Ther.* **102**, pzac066 (2022).

303   2.   Murphy, K. R. & Myors, B. Testing the hypothesis that treatments have negligible effects: Minimum-

304        effect tests in the general linear model. *J. Appl. Psychol.* **84**, 234–248 (1999).

305   3.   Rafi, Z. & Greenland, S. Semantic and cognitive tools to aid statistical science: replace confidence

306        and significance by compatibility and surprise. *BMC Med. Res. Methodol.* **20**, 244 (2020).

307   4.   Cohen, J. The earth is round ($p < .05$). *Am. Psychol.* **49**, 997–1003 (1994).

308   5.   Lakens, D. The Practical Alternative to the p Value Is the Correctly Used p Value. *Perspect. Psychol.*

309        *Sci.* **16**, 639–648 (2021).

310   6.   Herbert, R. Research Note: Significance testing and hypothesis testing: meaningless, misleading and

311        mostly unnecessary. *J. Physiother.* **65**, 178–181 (2019).

312   7.   Lakens, D. Why P values are not measures of evidence. *Trends Ecol. Evol.* **37**, 289–290 (2022).

313   8.   Muff, S., Nilsen, E. B., O'Hara, R. B. & Nater, C. R. Response to 'Why P values are not measures of

314        evidence' by D. Lakens. *Trends Ecol. Evol.* **37**, 291–292 (2022).

315   9.   Goodman, S. N. & Royall, R. Evidence and scientific research. *Am. J. Public Health* **78**, 1568–1574

316        (1988).

317 10.   Goodman, S. N. Toward Evidence-Based Medical Statistics. 1: The P Value Fallacy. *Ann. Intern.*

318        *Med.* **130**, 995–1004 (1999).

319 11.   Collaboration, O. S. Estimating the reproducibility of psychological science. *Science* **349**, aac4716

320        (2015).

321 12.   Patil, P., Peng, R. D. & Leek, J. T. What Should Researchers Expect When They Replicate Studies?

322        A Statistical View of Replicability in Psychological Science. *Perspect. Psychol. Sci.* **11**, 539–544

323        (2016).

324   13. Scheel, A. M., Schijen, M. R. M. J. & Lakens, D. An Excess of Positive Results: Comparing the

325       Standard Psychology Literature With Registered Reports. *Adv. Methods Pract. Psychol. Sci.* **4**,

326       25152459211007468 (2021).

327   14. Ioannidis, J. P. Why most published research findings are false. *PLoS Med.* **2**, e124 (2005).

328   15. Anderson, S. F. & Maxwell, S. E. Addressing the "Replication Crisis": Using Original Studies to

329       Design Replication Studies with Appropriate Statistical Power. *Multivar. Behav. Res.* **52**, 305–324

330       (2017).

331   16. Nosek, B. A. *et al.* Replicability, robustness, and reproducibility in psychological science. *Annu. Rev.*

332       *Psychol.* **73**, 719–748 (2022).

333   17. Elkins, M. R. *et al.* Correspondence: Response to Lakens. *J. Physiother.* **68**, 214 (2022).

334   18. Correspondence: Reward, but do not yet require, interval hypothesis tests. *J. Physiother.* **68**, 213–214

335       (2022).

336   19. Lakens, D. Correspondence: Reward, but do not yet require, interval hypothesis tests. *J. Physiother.*

337       **68**, 213–214 (2022).

338   20. Tenan, M. & Caldwell, A. A Critical Review of Phyiotherapy Editor's Comments on Statistical

339       Practice.

340   21. Elkins, M. R. *et al.* Statistical inference through estimation: recommendations from the International

341       Society of Physiotherapy Journal Editors. *J. Physiother.* **68**, 1–4 (2022).

342   22. Borg, D. N. *et al.* Sharing data and code: a comment on the call for the adoption of more transparent

343       research practices in sport and exercise science. (2020).

344   23. Caldwell, A. & Vigotsky, A. D. A case against default effect sizes in sport and exercise science.

345       *PeerJ* **8**, e10314 (2020).

346   24. McGrath, R. E. & Meyer, G. J. When effect sizes disagree: the case of r and d. *Psychol. Methods* **11**,

347       386 (2006).

348   25. Levine, T. R. & Hullett, C. R. Eta Squared, Partial Eta Squared, and Misreporting of Effect Size in

349       Communication Research. *Hum. Commun. Res.* **28**, 612–625 (2002).

350    26.  Dabija, D. I. & Jain, N. B. Minimal Clinically Important Difference of Shoulder Outcome Measures

351          and Diagnoses: A Systematic Review. *Am. J. Phys. Med. Rehabil.* **98**, 671–676 (2019).

352    27.  Fricker Jr, R. D., Burke, K., Han, X. & Woodall, W. H. Assessing the statistical analyses used in

353          basic and applied social psychology after their p-value ban. *Am. Stat.* **73**, 374–384 (2019).

354    28.  Sainani, K. L. The Problem with" Magnitude-based Inference". *Med. Sci. Sports Exerc.* **50**, 2166–

355          2176 (2018).

356    29.  Sainani, K. L., Lohse, K. R., Jones, P. R. & Vickers, A. Magnitude-based inference is not Bayesian

357          and is not a valid method of inference. *Scand. J. Med. Sci. Sports* **29**, 1428 (2019).

358    30.  Lohse, K. R. *et al.* Systematic review of the use of "magnitude-based inference" in sports science and

359          medicine. *PloS One* **15**, e0235318 (2020).

360    31.  Benjamin, D. J. *et al.* Redefine statistical significance. *Nat. Hum. Behav.* **2**, 6–10 (2018).

361    32.  Lakens, D. *et al.* Justify your alpha. *Nat. Hum. Behav.* **2**, 168–171 (2018).

362    33.  Amrhein, V. & Greenland, S. Remove, rather than redefine, statistical significance. *Nat. Hum. Behav.*

363          **2**, 4–4 (2018).

364    34.  McShane, B. B., Gal, D., Gelman, A., Robert, C. & Tackett, J. L. Abandon statistical significance.

365          *Am. Stat.* **73**, 235–245 (2019).

366    35.  Simmons, J. P., Nelson, L. D. & Simonsohn, U. Life after p-hacking. in *Meeting of the society for*

367          *personality and social psychology, New Orleans, LA* 17–19 (2013).

368    36.  Simmons, J. P., Nelson, L. D. & Simonsohn, U. False-positive psychology: undisclosed flexibility in

369          data collection and analysis allows presenting anything as significant. (2016).

370    37.  Sun, X. *et al.* Credibility of claims of subgroup effects in randomised controlled trials: systematic

371          review. *Bmj* **344**, (2012).

372    38.  Kerr, N. L. HARKing: Hypothesizing after the results are known. *Personal. Soc. Psychol. Rev.* **2**,

373          196–217 (1998).

374    39.  Rosenthal, R. The file drawer problem and tolerance for null results. *Psychol. Bull.* **86**, (1979).

375    40. Borg, D. N., Lohse, K. R. & Sainani, K. L. Ten common statistical errors from all phases of research,

376        and their fixes. *PM&R* **12**, 610–614 (2020).

377    41. Leek, J. T. & Peng, R. D. Statistics: P values are just the tip of the iceberg. *Nature* **520**, 612–612

378        (2015).

379    42. Wasserstein, R. L., Schirm, A. L. & Lazar, N. A. Moving to a world beyond "p< 0.05". *The*

380        *American Statistician* vol. 73 1–19 (2019).

381