

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29

No Estimation without Inference:

A Response to the International Society of Physiotherapy Journal Editors

Keith Lohse, PhD¹

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

Acknowledgments: I would like to thank Dr. Emma Johnson, Dr. Kristin Sainani, and two anonymous reviewers for their detailed comments on earlier drafts of this commentary.

Date Submitted

2022-08-03

Keywords:

“physical therapy”; “statistical significance”; “inference”; “estimation”

Corresponding Author:

Keith Lohse, PhD, PStat; lohse@wustl.edu

Abstract

30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48

The International Society of Physiotherapy Journal Editors (ISPJE) recently published an editorial warning that many of their journals would soon prohibit the use of null hypothesis tests and instead require authors to interpret 95% confidence intervals relative to clinically important values. Although I encourage the reporting of confidence intervals and the discussing of uncertainty in the context of a research question, the ISPJE's proposed ban is illogical and there are several instances of flawed statistical reasoning in the editorial. In brief, the editorial: (1) fails to adequately grapple with the inherent connection between hypothesis testing and estimation, (2) presents several misleading arguments about the perceived flaws of hypothesis tests, and (3) presents an alternative to hypothesis testing that is, in itself, a form of hypothesis test – the minimal effects test – albeit done informally. If the editorials' arguments are taken at face value, then that will lower the statistical literacy in our field and readers will have a flawed understanding of p-values. Further, if the editorials' proposed ban is put into practice, I fear that could decrease the scientific integrity of our research as it removes quantitative benchmarks in favor of a more subjective interpretation of confidence intervals. Ultimately, I think that many of the ISPJE's concerns that led to the editorial are valid, but I think those problems are the result of questionable research practices stemming from poor methodological training for authors, reviewers, and editors. These problems will only be fixed through better and continuing education, not the banning of statistically valid methods.

49 Recently, Elkins et al.¹ (hereafter referred to as “the Editorial”) published an editorial on behalf of
50 the International Society of Physiotherapy Journal Editors (ISPJE), recommending that researchers stop
51 using null hypothesis significance tests and adopt “estimation methods”. Further, the editorial warns that
52 this is not merely an idea to consider, but a coming policy of journals: “the [ISPJE] will be expecting
53 manuscripts to use estimation methods *instead* of null hypothesis statistical tests” (emphasis added).
54 However, the Editorial is deeply flawed in its statistical reasoning. If the proposed policies were adopted,
55 they could damage the statistical literacy and scientific integrity of the field.

56 I detail each of my critiques below, but in short the Editorial: (1) fails to adequately grapple with
57 the inherent connection between hypothesis testing and estimation as methods of statistical inference, (2)
58 presents several misleading arguments about the flaws of statistical significance tests, and (3) presents an
59 alternative that is, in itself, a form of significance testing – the minimal effects test² (but the alternative
60 does this implicitly and muddles two-sided and one-sided hypothesis testing). Finally, I end with a short
61 list of more urgent problems that the ISPJE could work to address.

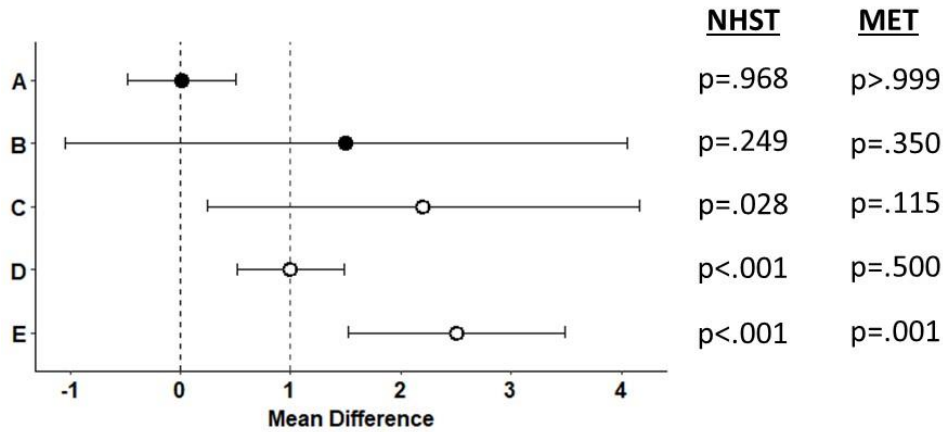
62 I commend the Editorial for encouraging researchers to think deeply about the statistical tools
63 available to them, to consider “practical significance” as well as “statistical significance”, and for
64 bringing important methodological discussions to the forefront of physical therapy research. However, the
65 central argument of the Editorial is illogical and I worry what coming policy changes might mean for how
66 authors interpret their data. I think the antidote to researchers making faulty decisions is not to ban p-
67 values, but to improve education. A rising tide lifts all boats, and if the baseline statistical literacy in our
68 field were higher, authors would make fewer mistakes, reviewers would be more apt to catch remaining
69 mistakes, and readers would be better equipped to make their own conclusions given the available data.
70 Editors then need to hold the line and ensure rigorous review, not ban valid statistical tools.

71 Hypothesis Testing and Estimation are Inescapably Intertwined

72 The Editorial presents hypothesis testing and estimation as two distinct methodological
73 approaches. However, these approaches are two sides of the same coin, as illustrated by a simple example
74 in Figure 1. When a 95% confidence interval excludes the null value, then one can reject the null
75 hypothesis at $p < .05$. This is because hypothesis tests and confidence intervals are based on the same
76 underlying mathematics: e.g., how big is the observed effect relative to the variability we would expect
77 due to sampling? Although typically we think of the null-hypothesis as an assumption of “no effect”, the
78 null hypothesis can assume zero or non-zero effects. So, as shown in the figure, we can ascertain the
79 probability of observing the data we did, assuming a null value of 0 or a null value of 1.

80 Hypothesis testing and estimation cannot be fully disentangled: estimation (frequentist or
81 Bayesian) asks about *plausible values* of the parameter in the population, hypothesis testing asks about
82 the plausibility of *a specific parameter value*. These are both inferences, because we are inferring
83 something about the population based on the data in our sample. In the frequentist paradigm, uncertainty
84 in the inference is accounted for with long-run error control; e.g., setting the Type 1 Error rate, $\alpha = 0.05$.
85 We can see this when running simulations as shown in Figure 1A-E: any confidence interval that does not
86 contain zero also has $p < 0.05$, for the null hypothesis significance test (NHST).

87 The 95% confidence interval shows values in the population that are *compatible* with what we
88 observed in the sample.³ That is, if you move outside of the confidence interval, any of those parameter
89 values (the “true” mean differences; Δ 's) would be statistically different from the mean difference
90 observed in the sample (\bar{x}_d) at the $p < 0.05$ level. Inside of the confidence interval, none of those parameter
91 values would be statistically different ($p > 0.05$) from the observed mean difference. Recall that the p-value
92 is the probability of observing data as extreme or more extreme, assuming that the null hypothesis is true,
93 formally written as $p(\geq \bar{x}_d | H_0)$.



94

95 **Figure 1.** 95% confidence intervals and corresponding p-values for testing $H_0: \Delta = 0$ (NHST, null
 96 hypothesis significance testing) and $H_0: \Delta \leq 1$ (MET, a one-sided minimal effects test). Open circles
 97 indicate mean differences with $p < 0.05$ for the NHST.

98

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

^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.

108 estimate in both the NHST and the MET^B, but the hypothesis test and the estimate are complementary and
109 connected.

110 **Misleading Arguments about flaws with Significance Tests**

111 The Editorial bases many arguments on a previous list of perceived problems from Herbert
112 (2019).⁶ The Herbert paper is in itself an editorial that presents informed arguments, but is not an
113 objective demonstration of any mathematical facts. So, reinforcing the Editorial's list through a citation to
114 Herbert does not provide an evidentiary foundation: it is layering opinion on top of opinion. Second, each
115 of the five “problems” outlined by the Editorial is either not really a problem inherent to p-values or the
116 problem is a true but misleading statement. I address each problem from the Editorial (in quotes) below:

117 **1. “A p-value is not the probability that a hypothesis is (or is not) true.”** – This is correct, but it does not
118 follow that this makes p-values, or even statistical significance tests, unhelpful or uninformative.

119 Knowing that the observed data are incompatible with some null value is a crucial step for many research
120 questions. For instance, hypothesis testing in early phase research can help us make decisions about
121 where to direct our resources, starting us down the road of replication and ultimately determining the
122 efficacy and effectiveness of an intervention.

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

126 Technical definitions of evidence can get a bit complicated and are debated.⁷⁻⁹ However, I would
127 invite readers to consider a simple example of absolute probability versus relative probability. If I find
128 that eating green jelly beans reduces post-surgical recovery time for the ACL by 10% relative to controls
129 with $p < 0.05$, then the most likely explanation is still that jelly beans have no effect on recovery and what I

^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

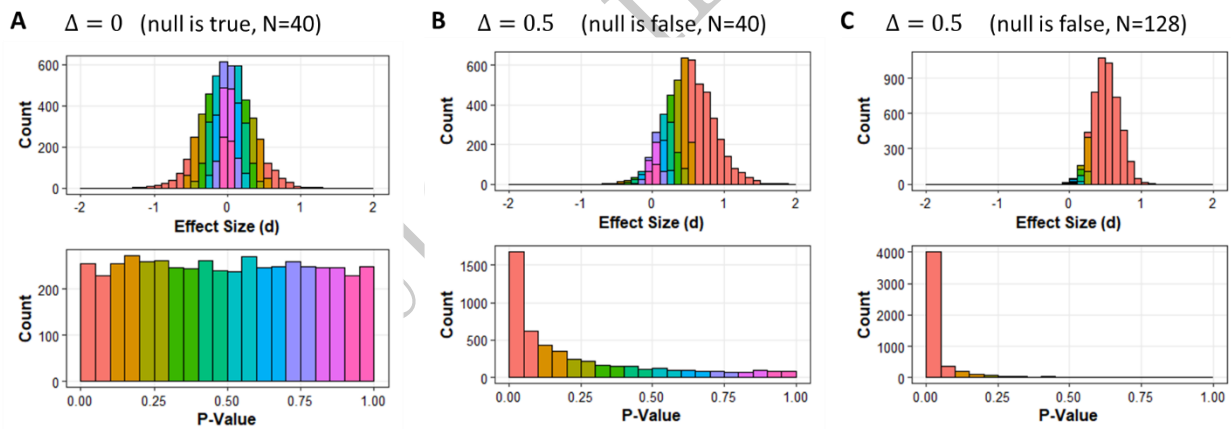
130 observed was chance fluctuation. That is, the null hypothesis is still the most likely explanation even
131 though p was <0.05 , because the baseline probability of “jelly bean efficacy” is very low and false
132 positives occur 5% of the time when $\alpha = 0.05$. Thus, the p-value is not in itself a measure of evidence,
133 because I would need additional *outside information* in order to change (or not change) my beliefs. As
134 Goodman and Royall⁹ write “The p-value is not adequate for inference because *the measurement of*
135 *evidence* requires are least three components: the observations, and two competing explanations for how
136 they were produced” (p. 1569; emphasis added).

137 Some researchers might think of the p-value as evidence against the null specifically, without the
138 need for comparison to a given alternative. But the p-value is calculated assuming that the null is true, so
139 again the Editorial is correct that we cannot simply flip the question around, assume the data, and get the
140 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
141 hypothesis being true, we would need Bayesian statistics in which we formalize some *prior* probability
142 about the null hypothesis.⁹ If we have a strong enough prior probability that the null is true, then the
143 current data in the sample may not lead us to change our beliefs based on the *posterior* distribution, no
144 matter how small the p-value. This was the case in my jelly bean example, where $p<0.05$ still did not
145 shake my belief in the null hypothesis. For any given prior distribution, however, there is a smaller
146 *likelihood* of observing highly discrepant effects (e.g., $|\bar{x}_d| \gg 0$), leading to a smaller relative probability
147 of 0 in the posterior distribution compared to the prior distribution.^c Updating the probability of 0 in the
148 posterior distribution reflects rational decision making in daily life. For instance, the first time I find jelly
149 beans reduce recovery time with $p<0.05$, I might rightly ignore that as a false positive. The fifth time I
150 find jelly beans reduce recovery time with $p<0.05$, I should take a long hard look at the ingredients and

^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 data) is proportional to the prior (our expectation before we saw the data) multiplied by the likelihood (which is the probability of the current data given the hypothesis): $posterior \propto likelihood \cdot prior$.

151 maybe my study procedures; as $p < 0.05$ is not always a sign that the null is wrong, but that some other
152 assumption has been violated.

153 Still, the p-value does not need to be a measure of evidence for it to be useful. Critically, small p-
154 values are *relatively* less likely to be observed when the null hypothesis is true compared to when an
155 alternative hypothesis is true. Thus, in a practical sense, p-values can help us make decisions about what
156 effects to study, assuming that we are testing at least some real effects. As shown in Figure 2A, p-values
157 have a uniform distribution under the null hypothesis, with 5% of p-values necessarily below 0.05.
158 However, if the null is not true, then we will see a shift in the distribution of p-values, with small p-values
159 becoming more common. An example of this is shown in Figure 2B, where the null is false and 34% of p-
160 values are below 0.05. However, correctly rejecting the null hypothesis only 34% of the time is not ideal,
161 so consider Figure 2C, where I have now tripled the sample size and 80% of p-values are below 0.05.
162 That is, with 64 people per group, we now have 80% statistical power to detect a $\Delta = 0.5$.



163

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

171

172 This is where the concept of a decision is important to distinguish from the term “evidence”.⁹
173 Without knowing the actual *evidence* against the null-hypothesis, if I *decide* to reject the null when
174 $p < 0.05$, then I will only be wrong 5% of the time (i.e., the Type 1 error rate). Similarly, if I have 80%
175 statistical power and a reasonable estimate for the smallest effect size of interest, then I only have a 20%
176 chance of missing an effect of that size (i.e., the Type 2 error rate). Mathematically, these probabilities are
177 robust if we accept the null-hypothesis as true and make minimal other assumptions, which is very helpful
178 when limited outside information is available. See Goodman quoting Neyman and Pearson about
179 hypothesis testing, “Without hoping to know whether each separate hypothesis is true or false, we may
180 search for rules to govern our behaviour with regard to them, in following which we insure that, in the
181 long run of experience, we shall not often be wrong.”¹⁰

182 So, p-values are not a measure of evidence, but they are useful tools for helping us make the
183 correct decision. If we want a proper measure of evidence for one hypothesis versus another, then we can
184 do more work, but we also need to make more assumptions and/or bring in outside information. This can
185 be both a feature and bug of *using* hypothesis tests. We can control long run error rates with minimal
186 information, but if we do that so habitually that we forget other information is available, then that is on us
187 not the p-value.

188 **3. “Statistically significant findings are not very replicable.”** – This is misleading. First, it is difficult to
189 precisely define replication,^{11,12} but if we think about “being replicable” as the probability that a
190 statistically significant result represents a real, non-zero effect then we would expect more statistically
191 significant findings to “replicate” provided that hypothesis tests have adequate statistical power,
192 researchers have not engaged in p-hacking, there is not selective reporting of results, etc. Thus, not all
193 statistically significant findings will replicate,¹³ but statistically significant findings in well-designed
194 studies are more likely to replicate.^{14–16} Second and by any definition, threats to replicability are also
195 going to affect confidence intervals (the Editorial’s proposed solution) as much as they affect p-values,
196 because, again, the p-value is intrinsically linked to the confidence interval. Thus, the Editorial is correct

197 in a practical sense: many statistically significant findings in the current literature do not replicate.
198 However, a lack of replication is the fault of poor study design and questionable research practices, not
199 the use of hypothesis tests as a method of inference.

200 **4. “In most clinical trials, the null hypothesis must be false.”** – This is arguably true but very
201 misleading. It is true that real treatment effects are unlikely to be precisely 0 (e.g., they might be +0.001),
202 but it begs the question: do we really care if the true effect is 0 or 0.001? And will we ever have the
203 statistical precision to discern that difference? All measurement has some error, so I would argue that
204 many effects are functionally 0 even if the (unknowable) true value is not actually zero. But, in a strict
205 mathematical sense I will concede the Editorial is correct, if we accept a hyper-precise definition, the
206 null-hypothesis of $H_0: \Delta = 0.\overline{00}$ will usually be false. However, if we accept that definition, then all
207 point-estimates are usually false and almost no value will be precisely the minimum clinically important
208 difference either, which is the Editorial’s proposed point-estimate in their alternative.

209 In response¹⁷ to an independent critique by Lakens¹⁸, this hyper-precise definition does seem to
210 be the argument that the editorial is making.^D They claim, “The assertion that the null hypothesis is false
211 in most clinical trials does not require empirical evidence, because it is self-evidently true” and “The null
212 hypothesis may often be approximately true, but it is rarely if ever exactly true”. The Editorial seems to
213 miss the point that the null is a useful *model*: testing against 0 is still useful for things that are
214 approximately 0. As an analogy, I have successfully gotten many places using maps, but none of those
215 maps was a photo-realistic version of reality.

216 Scientists are often working on the frontiers of human knowledge; this is costly work where we
217 need to explore a lot of different ideas and many them do not pan out. That is, many tested “effects” are

^D I was very excited to see the Lakens commentary¹⁸ and others¹⁹, 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*²⁰, 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.

218 functionally zero.¹⁴ So, simply because a point estimate of precisely 0 is unlikely to be true does not mean
219 that it is unhelpful to ask. It should be a very low bar to show that your clinical treatment has a non-zero
220 effect! Further, the Editorial is specifically critiquing this “nil” hypothesis (i.e., $H_0 = 0$), when we could
221 hypothesize any value, or avoid the point-null entirely with a one-sided test (i.e., $H_0 \leq 0$).^{2,5} So, if
222 assuming $H_0 = 0$ is not desirable, we can set that null value to be anything we want (i.e., $H_0: \Delta \leq 0.4$
223 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
224 example in Figure 1).

225 **5. “Researchers need information about the size of effects.”** – This is a true statement, but it is not a
226 problem with p-values nor null hypothesis significance tests. To my knowledge, no statistician has ever
227 recommended that applied researchers ignore measures of effect size (either raw or standardized).
228 Estimates of effect size are integral to any results section. I would even take this one step further and
229 encourage authors to share their data whenever possible²¹, enabling other researchers to calculate their
230 own effect sizes as there can be limitations with and confusion about standardized effects sizes, and there
231 is no one-size-fits-all solution to effect sizes^{22–24}.

232 **The Editorial’s “Alternative” is a Hypothesis Test – The Minimal Effects Test**

233 After detailing the potential problems with the NHST, the Editorial proposes an alternative
234 solution in which they encourage authors to compare their 95% confidence interval to some minimum
235 clinically meaningful value (which I will write as δ).^E Estimation is a good practice and I would
236 encourage researchers to report 95% confidence intervals and interpret their upper and lower limits in
237 context, when appropriate. However, what the Editorial is suggesting is effectively an MET where
238 $H_0: \Delta \leq \delta$. That is, if the test is to see if the 95% confidence interval does not contain δ , then that is
239 mathematically equivalent to an MET assuming $H_0: \Delta \leq \delta$ and finding $p < 0.025$. Note $p < 0.025$, not

^E I caution that it is difficult to find a single measure of δ ; it changes as a function of the study population, the study context, and has its own uncertainty due to sampling error.^{19,25}

240 $p < 0.05$, because most METs are one-sided hypothesis tests whereas confidence intervals are two sided
241 (see Figure 1 and Footnote A). After heavily critiquing hypothesis testing as a method of inference, the
242 Editorial ends up effectively proposing a hypothesis test. This is clearly an illogical proposition.

243 I want to emphasize that it is valid for the Editorial to recommend that authors consider their 95%
244 confidence interval relative to some clinically meaningful value. However, this is not an “alternative” to
245 conducting a null hypothesis significance test, it is in fact mathematically identical to conducting a null
246 hypothesis test with a carefully chosen null hypothesis. Both are valid.

247 I would add, however, that there can also be advantages to explicitly framing this as a hypothesis
248 test rather than the informal interpretation of a confidence interval. First, it encourages researchers to
249 explicitly commit to a specific δ while the study is being designed, rather than simply obtaining an
250 estimate of the effect and then comparing it to candidate δ 's post hoc. Second, it requires researchers to
251 think carefully about the direction of the test and the desired α -level, whereas simply invoking a 95%
252 confidence interval implicitly uses a two-tailed test and $\alpha = 0.05$, which may not be best suited to the
253 research question.

254 Finally, it is also important to stress that history provides us with several examples of how
255 authors will view their data through rose-tinted glasses when quantitative statistical safeguards are
256 removed. For instance, when *Basic and Applied Social Psychology* banned p -values, authors were found
257 to overstate their conclusions well beyond what would have been considered if “statistical significance”
258 had been a benchmark.²⁶ In sport and exercise science, “magnitude-based inference” was leveraged as a
259 niche method that allowed authors to interpret differences as meaningful when they had very little
260 statistical support (e.g., p 's > 0.25).^{27–29} Statistical significance in an NHST does not necessarily need to be
261 the benchmark nor 0.05 the default value^{30–33}, but it is always important to have a statistically sound
262 framework for dealing with uncertainty.

263 **Virtues of Hypothesis Testing**

264 One of the great virtues of null hypothesis significance testing is Type I error control while
265 making minimal assumptions about the nature of the data or the world at large. If we set $\alpha = 0.05$, then
266 we can be confident we will only get data greater than or equal to what we observed 5% of the time when
267 the null is true. Importantly, this works for a wide range of statistics and types of tests, including F - and
268 χ^2 -statistics that have multiple degrees of freedom from models asking questions about multiple effects
269 simultaneously. For instance, in a randomized controlled trial with three arms, I could conduct an
270 omnibus F -test and obtain a p -value to see if there is any evidence of a difference between groups overall,
271 before conducting additional post-hoc tests to compare specific groups. This situation is not covered by
272 the Editorial and the Editorial's confidence interval alternative is not easily applied here, although one
273 could plausibly adjust the width of the confidence intervals to control for multiple comparisons.

274 **Bigger Threats to Statistical Integrity**

275 Misinterpretation and misuse of p -values are threats to statistical integrity. However, questionable
276 research practices such as p -hacking, sub-group analyses, flexible stopping rules, selective exclusion of
277 outliers, selective reporting, and hypothesizing after results are known are much larger threats.³⁴⁻³⁸
278 Furthermore, these questionable research practices have consistently negative consequences regardless of
279 the method of inference. For instance, although the term " p -hacking" connotes the NHST, these
280 questionable research practices pose an equal threat to confidence intervals because again confidence
281 intervals and p -values are based on the same underlying mathematics. Similarly, switching to a fully
282 Bayesian method of analysis is not an antidote for poor study design, small samples, and questionable
283 research practices. As others have argued,^{39,40} p -values get a disproportionate amount of attention in
284 popular discussions of research methodology. I encourage the ISPJE to instead focus their attention on
285 methods for improving data/code sharing, transparency, and replicability through tools like
286 preregistration, results-blind review, registered reports, or even "data papers" whose primary function is
287 to report a study and archive the data, without drawing inferences from limited samples.

288 It is fair to say that p-values are often mis-used and mis-interpreted, and “statistical significance”
289 may not ultimately be the best term for applied researchers to use.⁴¹ However, it is incorrect to present
290 these human errors as inherent flaws in hypothesis testing. For instance, if someone mis-interprets $p>0.05$
291 as evidence of “no difference”, then I would argue the correct action is to teach them about equivalence
292 tests and non-inferiority designs, not ban p -values. Similarly, there are times when Bayesian inference is
293 what authors are really interested in (e.g., what is the probability that the null is true, given the
294 evidence?), and in those cases Bayesian inference can and should be used. However, Bayesian analysis is
295 not a panacea and needs to be used thoughtfully like any statistical tool. So, although a simple heuristic of
296 $p<0.05$ may well be overused as “the” test in physical therapy research, frequentist hypothesis tests are
297 still valid and useful tools for physical therapy researchers. Moreover, the scientific integrity of the field
298 has much larger concerns, and both p-values and confidence intervals will be corrupted by p -hacking,
299 under-powered subgroup analyses, surrogate outcomes, and other questionable research practices.

300 In conclusion, I agree with the Editorial on the importance of reporting effect sizes and
301 interpreting them in context. However, the Editorial makes numerous statistical faux pas that could harm
302 the statistical literacy in our field, if readers take them at face value, and harm the scientific integrity of
303 our field, if put into editorial practice.

304

305 **Acknowledgments:** I would like to thank Dr. Emma Johnson, Dr. Kristin Sainani, and two anonymous
306 reviewers for their detailed comments on earlier drafts of this commentary.

307

308 **Funding Sources:** None

309

310 **Data Sharing and Supplementary Material Accessibility Statement:** R code for all analyses and
311 simulations presented in this commentary are included as a digital supplement on SportRxiv

312 (<https://sportrxiv.org/index.php/server/preprint/view/178/version/211>).

313

314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337

References

1. Elkins, M. R. *et al.* Statistical inference through estimation: recommendations from the International Society of Physiotherapy Journal Editors. *Phys. Ther.* **102**, pzac066 (2022).
2. Murphy, K. R. & Myers, B. Testing the hypothesis that treatments have negligible effects: Minimum-effect tests in the general linear model. *J. Appl. Psychol.* **84**, 234–248 (1999).
3. Rafi, Z. & Greenland, S. Semantic and cognitive tools to aid statistical science: replace confidence and significance by compatibility and surprise. *BMC Med. Res. Methodol.* **20**, 244 (2020).
4. Cohen, J. The earth is round ($p < .05$). *Am. Psychol.* **49**, 997–1003 (1994).
5. Lakens, D. The Practical Alternative to the p Value Is the Correctly Used p Value. *Perspect. Psychol. Sci.* **16**, 639–648 (2021).
6. Herbert, R. Research Note: Significance testing and hypothesis testing: meaningless, misleading and mostly unnecessary. *J. Physiother.* **65**, 178–181 (2019).
7. Lakens, D. Why P values are not measures of evidence. *Trends Ecol. Evol.* **37**, 289–290 (2022).
8. Muff, S., Nilsen, E. B., O’Hara, R. B. & Nater, C. R. Response to ‘Why P values are not measures of evidence’ by D. Lakens. *Trends Ecol. Evol.* **37**, 291–292 (2022).
9. Goodman, S. N. & Royall, R. Evidence and scientific research. *Am. J. Public Health* **78**, 1568–1574 (1988).
10. Goodman, S. N. Toward Evidence-Based Medical Statistics. 1: The P Value Fallacy. *Ann. Intern. Med.* **130**, 995–1004 (1999).
11. Collaboration, O.-S. Estimating the reproducibility of psychological science. *Science* **349**, aac4716 (2015).
12. Patil, P., Peng, R. D. & Leek, J. T. What Should Researchers Expect When They Replicate Studies? A Statistical View of Replicability in Psychological Science. *Perspect. Psychol. Sci.* **11**, 539–544 (2016).

- 338 13. Scheel, A. M., Schijen, M. R. M. J. & Lakens, D. An Excess of Positive Results: Comparing the
339 Standard Psychology Literature With Registered Reports. *Adv. Methods Pract. Psychol. Sci.* **4**,
340 25152459211007468 (2021).
- 341 14. Ioannidis, J. P. Why most published research findings are false. *PLoS Med.* **2**, e124 (2005).
- 342 15. Anderson, S. F. & Maxwell, S. E. Addressing the “Replication Crisis”: Using Original Studies to
343 Design Replication Studies with Appropriate Statistical Power. *Multivar. Behav. Res.* **52**, 305–324
344 (2017).
- 345 16. Nosek, B. A. *et al.* Replicability, robustness, and reproducibility in psychological science. *Annu. Rev.*
346 *Psychol.* **73**, 719–748 (2022).
- 347 17. Elkins, M. R. *et al.* Correspondence: Response to Lakens. *J. Physiother.* **68**, 214 (2022).
- 348 18. Lakens, D. Correspondence: Reward, but do not yet require, interval hypothesis tests. *J. Physiother.*
349 **68**, 213–214 (2022).
- 350 19. Tenan, M. & Caldwell, A. A Critical Review of Phyiotherapy Editor’s Comments on Statistical
351 Practice.
- 352 20. Elkins, M. R. *et al.* Statistical inference through estimation: recommendations from the International
353 Society of Physiotherapy Journal Editors. *J. Physiother.* **68**, 1–4 (2022).
- 354 21. Borg, D. N. *et al.* Sharing data and code: a comment on the call for the adoption of more transparent
355 research practices in sport and exercise science. (2020).
- 356 22. Caldwell, A. & Vigotsky, A. D. A case against default effect sizes in sport and exercise science.
357 *PeerJ* **8**, e10314 (2020).
- 358 23. McGrath, R. E. & Meyer, G. J. When effect sizes disagree: the case of r and d. *Psychol. Methods* **11**,
359 386 (2006).
- 360 24. Levine, T. R. & Hullett, C. R. Eta Squared, Partial Eta Squared, and Misreporting of Effect Size in
361 Communication Research. *Hum. Commun. Res.* **28**, 612–625 (2002).
- 362 25. Dabija, D. I. & Jain, N. B. Minimal Clinically Important Difference of Shoulder Outcome Measures
363 and Diagnoses: A Systematic Review. *Am. J. Phys. Med. Rehabil.* **98**, 671–676 (2019).

- 364 26. Fricker Jr, R. D., Burke, K., Han, X. & Woodall, W. H. Assessing the statistical analyses used in
365 basic and applied social psychology after their p-value ban. *Am. Stat.* **73**, 374–384 (2019).
- 366 27. Sainani, K. L. The Problem with " Magnitude-based Inference". *Med. Sci. Sports Exerc.* **50**, 2166–
367 2176 (2018).
- 368 28. Sainani, K. L., Lohse, K. R., Jones, P. R. & Vickers, A. Magnitude-based inference is not Bayesian
369 and is not a valid method of inference. *Scand. J. Med. Sci. Sports* **29**, 1428 (2019).
- 370 29. Lohse, K. R. *et al.* Systematic review of the use of “magnitude-based inference” in sports science and
371 medicine. *PloS One* **15**, e0235318 (2020).
- 372 30. Benjamin, D. J. *et al.* Redefine statistical significance. *Nat. Hum. Behav.* **2**, 6–10 (2018).
- 373 31. Lakens, D. *et al.* Justify your alpha. *Nat. Hum. Behav.* **2**, 168–171 (2018).
- 374 32. Amrhein, V. & Greenland, S. Remove, rather than redefine, statistical significance. *Nat. Hum. Behav.*
375 **2**, 4–4 (2018).
- 376 33. McShane, B. B., Gal, D., Gelman, A., Robert, C. & Tackett, J. L. Abandon statistical significance.
377 *Am. Stat.* **73**, 235–245 (2019).
- 378 34. Simmons, J. P., Nelson, L. D. & Simonsohn, U. Life after p-hacking. in *Meeting of the society for*
379 *personality and social psychology, New Orleans, LA* 17–19 (2013).
- 380 35. Simmons, J. P., Nelson, L. D. & Simonsohn, U. False-positive psychology: undisclosed flexibility in
381 data collection and analysis allows presenting anything as significant. (2016).
- 382 36. Sun, X. *et al.* Credibility of claims of subgroup effects in randomised controlled trials: systematic
383 review. *Bmj* **344**, (2012).
- 384 37. Kerr, N. L. HARKing: Hypothesizing after the results are known. *Personal. Soc. Psychol. Rev.* **2**,
385 196–217 (1998).
- 386 38. Rosenthal, R. The file drawer problem and tolerance for null results. *Psychol. Bull.* **86**, (1979).
- 387 39. Borg, D. N., Lohse, K. R. & Sainani, K. L. Ten common statistical errors from all phases of research,
388 and their fixes. *PM&R* **12**, 610–614 (2020).

- 389 40. Leek, J. T. & Peng, R. D. Statistics: P values are just the tip of the iceberg. *Nature* **520**, 612–612
390 (2015).
- 391 41. Wasserstein, R. L., Schirm, A. L. & Lazar, N. A. Moving to a world beyond “ $p < 0.05$ ”. *The*
392 *American Statistician* vol. 73 1–19 (2019).
- 393

Article in Press