

Middlesex University Research Repository

An open access repository of
Middlesex University research

<http://eprints.mdx.ac.uk>

Turner, Anthony N. ORCID: <https://orcid.org/0000-0002-5121-432X>, Parmar, Nimai ORCID:
<https://orcid.org/0000-0001-5540-123X>, Jovonoski, Alex and Hearne, Gary ORCID:
<https://orcid.org/0000-0003-2146-4878> (2021) Assessing group-based changes in
high-performance sport. Part 1: null hypothesis significance testing and the utility of p values.
Strength & Conditioning Journal, 43 (3) . pp. 112-116. ISSN 1524-1602 [Article]
(doi:10.1519/SSC.0000000000000625)

Final accepted version (with author's formatting)

This version is available at: <https://eprints.mdx.ac.uk/31203/>

Copyright:

Middlesex University Research Repository makes the University's research available electronically.

Copyright and moral rights to this work are retained by the author and/or other copyright owners unless otherwise stated. The work is supplied on the understanding that any use for commercial gain is strictly forbidden. A copy may be downloaded for personal, non-commercial, research or study without prior permission and without charge.

Works, including theses and research projects, may not be reproduced in any format or medium, or extensive quotations taken from them, or their content changed in any way, without first obtaining permission in writing from the copyright holder(s). They may not be sold or exploited commercially in any format or medium without the prior written permission of the copyright holder(s).

Full bibliographic details must be given when referring to, or quoting from full items including the author's name, the title of the work, publication details where relevant (place, publisher, date), pagination, and for theses or dissertations the awarding institution, the degree type awarded, and the date of the award.

If you believe that any material held in the repository infringes copyright law, please contact the Repository Team at Middlesex University via the following email address:

eprints@mdx.ac.uk

The item will be removed from the repository while any claim is being investigated.

See also repository copyright: re-use policy: <http://eprints.mdx.ac.uk/policies.html#copy>

1 **Assessing group-based changes in high-performance sport.**

2 **Part 1: Null hypothesis significance testing and the utility of *p* values**

3
4 Anthony N Turner, Nimai Parmar, Alex Jovanoski, Gary Hearne

5 London Sport Institute, Middlesex University

6
7 **Abstract**

8 The role of a strength and conditioning coach (SCC) has evolved over the last 10 years to
9 accommodate the large influx of data now available. As such, today's SCC must extend their
10 skill set to include data analysis, understanding the validity and utility of *p* values, effect sizes,
11 confidence intervals, and terms such as the smallest worthwhile change, and minimal
12 difference. The aim of part one of this two-part review is to define and discuss the utility of
13 null hypothesis significance testing (NHST), *p* values, and error rates. In part two, we introduce
14 effect sizes, measures of variability, and confidence intervals, culminating in recommendations
15 as to which may be the most viable options within the context of performance-based sport, and
16 thus potential methods to report group-based changes.

17
18
19 **Introduction**

20 The role of a strength and conditioning coach (SCC) has evolved over the last 10 years, perhaps
21 to the point where the term strength and conditioning scientist may be just as apt. In principle,
22 this is because SCC's are likely to spend as much time behind a computer analysing the streams
23 of data they have just collected, as they are coaching athletes in the gym. While acknowledging
24 that evidence-based practice has always been at the root of this role, perhaps today's shift is
25 due to some of the following: (a) the need to demonstrate an objective approach to appease
26 various stakeholders, (b) an increase in academic scientific-based degrees in this discipline of
27 sport science, and (c) perhaps most causative, an influx of performance-based (affordable)
28 software and hardware that provide vast quantities of performance variables. These advances
29 mean that today's SCC must extend their skill set to include data analysis. While the evolution
30 of this role, of course, requires the development of many other skill sets (Stewart, Comfort, &
31 Turner, 2017), this paper focuses on analysing data, and in particular, assessing group changes
32 (differences) in performance in applied settings. This appears to be an area of interest at the
33 moment given the confusion in the best method to do this, from the implementation of *p*-values
34 (Greenland, et al., 2016; Wasserstein & Lazar, 2016), to reporting effect sizes (Cohen 1988,

35 1992), confidence intervals (Cumming G. , 2014; Cumming & Finch, 2001), and the utility of
36 analyses incorporating concepts such as the smallest worthwhile change (Hopkins, 2004) and
37 minimal difference (Weir, 2005).

38

39 The aim of this two-part review therefore, is to discuss each of these in turn and provide
40 recommendations to the reader as to which may be most viable within the context of
41 performance-based sport, and thus potential methods to report group-based changes over time.
42 Furthermore, such an analysis will also assist practitioners in critiquing relevant research
43 findings when considering adoption of new strategies to practice. Here, in part one, we first
44 define and discuss the utility of null hypothesis significance testing (NHST), p values, and error
45 rates, given that NHST is the most commonly taught approach to testing research questions
46 with statistical models, and thus the most well-known and used within the literature
47 (Wasserstein & Lazar, 2016); in essence, this will serve as the platform from which we can
48 develop our statistical approach to analyse our data. In part two then, through a series of worked
49 examples, we will introduce effect sizes, measures of variability, and confidence intervals,
50 culminating in recommendations as to which may be the most viable options for the applied
51 analysis of data relating to group-based changes in strength and conditioning.

52

53 **Null Hypothesis Significance Testing**

54 We should also start by recognising the work of Sir Ronald Fisher, who is considered a pioneer
55 of statistics and devised the p value that we will shortly explain; Fisher also coined the term
56 “*statistically significant*” (Fisher, 1925). Actually, the term “*significant*” is now recognised as
57 a poor choice of word and is consequently considered as one of the seven most misused words
58 in science (Ghose, 2013). What Fisher actually meant, was along the lines of *statistically*
59 *interesting and requires further scrutiny* (Wasserstein R. , 2019). However, its meaning was
60 taken literally and is one of many reasons around the confusion of what a p value is, why
61 statisticians have repeatedly asked us to refrain from using the phrase “*statistically significant*”
62 (Wasserstein, Schirm, & Lazar, 2019; Amrhein, Greenland, & McShane, 2019; Lakens, et al.,
63 2018), and why one journal even felt they had to ban it from use (Traimow & and Marks,
64 2015).

65

66 In actuality, the p value (and the NHST framework that commonly relies on it) can be a useful
67 tool, if we appreciate its true meaning and utility (Greenland, 2019), with it offering a first line
68 of defence against us being fooled by random error and any confirmation bias we may have

69 towards a theory (Lakens, 2019). In understanding the p value, we should first note that NHST
70 is part of frequentist statistics, which means that it is concerned with the interpretation of
71 probability, and specifically, long run probability; that is, what the likely results of a study
72 would be, if repeated over and over again (Greenland, et al., 2016; Wasserstein & Lazar, 2016).
73 So, when you appreciate the frequentist statistical framework, you can note that the results of
74 any single study, only tell you what would happen if it were infinitely repeated (Greenland, et
75 al., 2016) and do not actually relate to your single use study. As such, p values can never be
76 regarded as evidence of the error or effect in your study (Greenland, et al., 2016). An additional
77 misconception occurs when we don't appreciate the name and literal meaning of the test we
78 are conducting, i.e., NHST. An NHST does exactly that, it typically investigates a test statistic
79 obtained from a parameterized model against a null model, which centres on there being no
80 effect or difference noted in your data, i.e., the result will be zero (however, via chance and
81 random variation, we more than likely observe variability around zero). By way of example, if
82 we introduce a new exercise to an intervention group (to check it works) and compare it to a
83 control group, we are not actually testing the hypothesis that this new exercise will improve
84 performance (which is referred to as the alternate hypothesis), but rather, that there is no
85 difference between the groups (i.e., the null). This can be a confusing concept to grasp, because
86 really, what we want to know, is if our alternate hypothesis (does our new exercise intervention
87 work) is true or false. But as applied sport scientists, using NHST, we must appreciate that this
88 is not the question that we are answering. Instead, we are testing the probability of our data,
89 given our (null) hypothesis, which is written as $P(D|H_0)$. We are not testing our (alternate)
90 hypothesis, given our data, written as $P(H_1|D)$. For the latter, we would need to use Bayesian
91 statistics, but in any case, these are not the same thing and this explains why when reading
92 around this issue, you may see $P(D|H) \neq P(H|D)$. We should also point out that NHST need not
93 always be about null hypothesis (zero effect or difference) testing, and that hypotheses centring
94 on equivalence, non-inferiority, and superiority, can also be tested. For example, Lakens (2017)
95 provides examples where equivalence testing may be more advantageous, especially in
96 scenarios where researchers want to argue for the absence of an effect that is large enough to
97 be worthwhile to examine, and where researchers should also consider the effect size under the
98 alternative hypothesis (we discuss this concept more in part two). Equivalence testing is beyond
99 the scope of this text, so we recommend readers to the work of Lakens (2017); we suspect that
100 in a large number of sports performance based research, this may be more appropriate than null
101 hypothesis testing.

103 If we continue with the example of comparing the difference between two independent groups
104 against a null model, we would then choose the t -test as our statistical model, ensuring we have
105 met the model's underlying assumptions, such as independent groups, normal distribution of
106 the means, and characteristically similar samples for example. By ensuring we have met these
107 assumptions, we can *act* as if the only thing that differs between groups is our training
108 intervention (Lakens, 2019) – all the while acknowledging that our data will always contain
109 random error and thus noise, which we can never fully identify, control or exclude. Neyman
110 and Pearson (1933) suggest the term *act* as a way forward with NHST, given it does not imply
111 we truly believe the results directly relate to our single study, which in any case, did not
112 investigate $P(H_1|D)$. We then run our test and generate our test statistic, which in this case is
113 the t -statistic. The t -statistic we get, coupled with the sample size, is used to generate a p value,
114 which informs us of the probability of obtaining our result (or more extreme), assuming the
115 null hypothesis (and all the statistical model's assumptions) is true. Remember, that if the null
116 hypothesis is true, the difference between groups would be around zero (assuming that was the
117 threshold you decided on), and in this instance, all p values are equally likely (Caldwell &
118 Chevront, 2019).

119

120 Now, if you get a low p value, you can then say that your data (not your alternate hypothesis,
121 as you did not test this) is not compatible with the statistical model (and all its underlying
122 assumptions) and thus is *interesting and requires further scrutiny*. The next question that
123 logically arises then, is how low does the p value need to be, to spark your interest and for you
124 to be satisfied that you are not merely measuring noise? Well, in answering this, let's first
125 explain what the p value actually tells us. Say you obtained a p value of 0.03, this would mean
126 that if the null hypothesis were true (and all the assumptions made by the underlying model),
127 the probability of obtaining such a result (or more extreme) is 3%. This is now a good point to
128 also introduce the term alpha (α), which describes the error rate you settled on prior to
129 undertaking the research. By convention, we use an α level of 0.05, which implies that we
130 accept the probability that we will get a false-positive (which is called a Type I error) in 5% of
131 future studies, when using the same model and similar samples. This now brings us full circle
132 back to *statistical significance*. In our example, we got a p value of 0.03, which is less than our
133 pre-defined (conventional) α of 0.05. Given this, we historically conclude we have
134 "*statistically significant*" results. Furthermore, when our p value is less than our accepted 5%
135 error rate, if we choose to *act* as if there is in an effect when there really is not, in the long run,

136 you won't be wrong more than 5 % of the time. To note again, if the null is true, all p values
137 are equally likely.

138

139 The drive to steer researchers away from using “*statistically significant*” is because it creates a
140 dichotomy of evidence (McShane & Gal, 2017), whereby values on one side are important and
141 meaningful (i.e., statistically significant) and values on the other side are unimportant and
142 unhelpful (i.e., statistically non-significant). For example, if one strength training intervention
143 results in $p = 0.049$, while another $p = 0.051$, we tend to deem the former as being an effective
144 intervention and we would probably plan on implementing it with future athletes. The latter
145 however ($p = 0.051$), would be deemed as non-statistically significant and thus ineffective, and
146 we would therefore not plan to use it any longer. This thinking is of course incorrect and a
147 consequence of categorical thinking. Rather, it has been argued that the p value should be
148 treated and reported as a continuous quantity between 0 and 1, e.g., $p = 0.06$ (Greenland, et al.,
149 2016), with us acknowledging that it does not tell us which assumption is incorrect; it could be
150 the null or any of the model's underlying assumptions. Equally, any noted “effect” is subject
151 to the *statistical power* of a test, which is discussed later in this paper. Such dichotomous
152 thinking also drives publication bias (Franco, Malhotra, & Simonovits, 2014) and the “*file*
153 *drawer effect*” (Rosenthal, 1979), whereby in some cases, we only get to read of studies that
154 were statistically significant; this in turn can motivate p -hacking (the practice of flexibly
155 analysing data until the p value passes the “significance threshold”). Collectively, these
156 negative consequences have amongst other suggestions including total abandonment of NHST,
157 resulted in calls for the alpha level to be lowered to $p = 0.005$ (Benjamin, et al., 2018). It is
158 suggested that this change in critical threshold would help to reduce the number of published
159 false-positives and the generation of weak evidence.

160

161 **The appropriate use of p values**

162 Firstly, the p value is a continuous probability and should be reported as such. While we act as
163 though we have met all the model's underlying assumptions, this is actually very challenging
164 and explains why, when coupled with the random variation that occurs in all data, we have
165 profound study replication issues (Cumming G. , 2014; Caldwell & Chevront, 2019). For
166 example, even if the exact same study was repeated with similar samples, it would generate
167 different p values most of the time (Cumming G. , 2014; Amrhein, Greenland, & McShane,
168 2019). As such, it is difficult to profess that the difference between $p = 0.049$ and $p = 0.051$ is
169 anything other than random variation, as opposed to a discrepancy residing only with

170 incompatibility of the data with the null. Instead then, we should simply say that assuming the
171 null hypothesis were true (and all the assumptions made by the underlying model), the
172 probability of obtaining such a result (or more extreme) is 5.1%. Now given that in sport
173 science, unlike medicine perhaps, the consequence of making a Type 1 error (i.e., a false-
174 positive) will unlikely be fatal or lead to any health complications, let alone lead to injury, it is
175 probably okay to increase the *a priori* α level. Such thinking is in line with Lakens et al.,
176 (2018), who state that the α level should be adjusted based on the context at hand and on the
177 cost of false-negatives (Type II errors): a higher α would be used by those for whom false-
178 positives are relatively inconsequential, and lower α would be used by those for whom false-
179 positives could be disastrous. In sport, we argue that in some cases, an α level as high as 0.1
180 (i.e., a 10% error rate) could be justified. For example, in performance sport, success is based
181 on the smallest of margins (a statement that every Olympic Games final proves testament to)
182 and professional athletes are often butting up against their genetic ceiling. As such, it is more
183 important to reduce false-negatives and thus potential opportunities that may stimulate positive
184 adaptations.

185

186 Irrespective of the α level, no decision should be made solely on the strength of a *p* value
187 (Wasserstein & Lazar, 2016), given its indirect link to the alternate hypothesis, which is based
188 on many assumptions that cannot be individually accounted for. Therefore, we need to move
189 away from lazy dichotomous thinking (Gardner & Altman, 1986), and accept the uncertainty
190 of our data and embrace its variation (Wasserstein, Schirm, & Lazar, 2019). One way to do this
191 is to use confidence intervals, as well as acknowledging that one single study can never be
192 taken as conclusive evidence of a new theory, irrespective of how low a *p* value is. With respect
193 to the latter, this is why meta-analyses are so important to any field of study. That said, when
194 considering the aforementioned publication bias, it is also interesting to consider if results
195 derived from meta-analysis are in fact over inflated effect sizes, given that with far less
196 frequency do we read about an intervention that does not work. Finally, we need to identify the
197 magnitude of the effect, which NHST does not do. For example, rather than inferring an effect
198 occurred, it would be far more useful to actually quantify the magnitude of the effect, such that
199 coaches who are looking to adopt this new exercise can base decisions also on whether changes
200 were small, moderate or large. However, we should again be cautious about applying critical
201 thresholds to our data and perhaps consider the smallest effect size of interest (SESOI), which
202 we could calculate if we appreciate the variability in our data or a particular target we are

203 aiming toward. Effect sizes, including the SESOI, as well as confidence intervals, are discussed
204 in part 2 of this 2-part review.

205

206 **Statistical Power**

207 The final element to address as part of NHST, is statistical power, which is defined as the
208 ability of a test to detect an effect when one exists. Statistical power should thus be considered
209 by researchers and applied practitioners before they undertake any experiment. This is because
210 studies that are woefully underpowered can be a waste of resources as well as time for all those
211 involved (Caldwell & Cheuvront, 2019). For example, Caldwell and Cheuvront (2019)
212 illustrated the results of a data simulation test, involving 100,000 repetitions, to demonstrate
213 the distribution of p -values when the null hypothesis was false (i.e., there was an effect). When
214 they ran the simulation with 80% statistical power, 80,000 simulations (80%) correctly
215 identified a “statistically significant” effect, meaning that in 20,000 simulations (20%), a type
216 II error (false-negative) occurred. They then repeated the simulations, but this time with 50%
217 power, and unsurprisingly, 50,000 simulations (50%) correctly identified a “statistically
218 significant” effect and the remaining 50,000 simulations (50%) missed it, generating a type II
219 error. The natural conclusion to be reached here is, are studies worth doing (even from an
220 ethical perspective) when the chance of finding an effect, if one exists, is 50-50? Probably not
221 given you increase the risk of making erroneous conclusions if your decisions are based solely
222 on p -values. So, let’s now look at how we determine statistical power but first noting that again,
223 as with p -values and confidence intervals, this probability is defined over repetitions of the
224 same study design and so is a frequency probability (Greenland, et al., 2016).

225

226 Statistical power can be calculated using a host of statistical software, some are free such as
227 G* Power, whereby you need only enter your pre-determined α level, sample size, and the
228 SESOI. Given the requirement of these data, statistical power is considered a *conditional*
229 *probability* (Caldwell & Cheuvront, 2019). For example, using G*Power for the calculation of
230 statistical power, if we wanted to compare two independent groups (to see who could jump
231 highest for example), and we used the conventional α of 0.05 and the conventional power of
232 80%, as well as aiming to detect an effect size (or magnitude of difference between groups) of
233 half a standard deviation, we would need 64 participants per group. Increasing the α to 0.1,
234 reduces the number of participants to 51 per group. It is not hard to appreciate therefore, that
235 many studies undertaken in sport are likely underpowered and some true effects are missed.

236 This understanding should serve to justify the need for additional analysis such as effect sizes
237 and confidence intervals to make inferences of whether an effect or difference was observed
238 or not.

239

240 **Conclusion**

241 The roles and responsibilities of today's SCC means they must extend their skill set to include
242 data analysis. NHST, along with its derived p value, can be a useful tool for this, if we
243 appreciate its true meaning and utility, with it offering a first line of defence against us being
244 fooled by random error and any confirmation bias we may have towards a theory. Importantly,
245 we must note that NHST is part of the frequentist framework of statistics and thus refers to
246 long run probability, with results from any single study used to infer what would happen if the
247 study was repeated over and over again under identical conditions with different but identically
248 distributed samples. Also, we must appreciate that our ability to find an effect, when one exists,
249 is affected by statistical power – if this is too low, the utility of NHST is questionable when the
250 expected effect (or difference between groups) is hypothesised to be small. Finally, if we do
251 choose to use thresholds to limit the error within which we are happy to operate, we should
252 choose an alpha level that represents the context at hand and the risks associated with Type I
253 and II errors. Either way, we must recognise that the p value is a continuous variable, and thus
254 should be reported as such. Therefore, practitioners using p values should conclude with a
255 statement along the following lines (in this example let's say we got $p = 0.083$): *Assuming the*
256 *null hypothesis were true and all the assumptions made by the underlying model, the*
257 *probability of obtaining such a result or more extreme, is 8.3%. Furthermore, given our alpha*
258 *level of 0.1, if we choose to act as if there is in an effect when in fact there is not, in the long*
259 *run, we won't be wrong more than 10 % of the time.*

260

261 Winter et al., (2014) nicely summarise the essence of NHST via Karl Popper's principle of
262 falsifiability, that is, before something can be accepted the opposite has to be shown to be
263 untenable. So, in closing, it is prudent to again reinforce that we are analysing the probability
264 of our data, given our (null) hypothesis, which is written as $P(D|H_0)$. In performance-based
265 sport, however, we may determine that NHST isn't necessary or appropriate, going straight to
266 methods that determine the *practical significance* of our data using effect sizes, and embrace
267 the uncertainty of our data through confidence intervals. Through a series of worked examples,
268 these will be explored in Part 2.

269 **References**

- 270 1. Amrhein, V., Greenland, S., & McShane, B. (2019). Scientists rise up against statistical
271 significance. *Nature*, 567, 305.
- 272 2. Benjamin, D., Berger, J., Johannesson, M., Nosek, B., Wagenmakers, E., Berk, R., &
273 Cesarini, D. (2018). Redefine statistical significance. *Nature Human Behaviour*, 2, 6.
- 274 3. Caldwell, A., & Cheuvront, S. (2019). Basic statistical considerations for physiology: The
275 journal Temperature toolbox. *Temperature*, 6, 181-210.
- 276 4. Cumming, G. (2014). The new statistics: Why and how. *Psychological science*, 25(1), 7-
277 29.
- 278 5. Cumming, G., & Finch, S. (2001). A primer on the understand- ing, use, and calculation of
279 confidence intervals that are based on central and noncentral distributions. *Educational and*
280 *Psychological Measurement*, 61, 532-574.
- 281 6. Fisher, R. (1925). *Statistical Methods for Research Workers*. Edinburgh: Oliver & Boyd .
- 282 7. Franco, A., Malhotra, N., & Simonovits, G. (2014). Publication bias in the social sciences:
283 Unlocking the file drawer. *Science*, 345(6203), ., 345, 1502-1505.
- 284 8. Gardner, M., & Altman, D. (1986). Confidence intervals rather than P values: estimation
285 rather than hypothesis testing. *Br Med J*, 292, 746-750.
- 286 9. Ghose, T. (2013, April 2). "Just a Theory": 7 Misused Science Words. Retrieved from
287 Scientific American: [https://www.scientificamerican.com/article/just-a-theory-7-misused-](https://www.scientificamerican.com/article/just-a-theory-7-misused-science-words/)
288 [science-words/](https://www.scientificamerican.com/article/just-a-theory-7-misused-science-words/)
- 289 10. Greenland, S. (2019). Valid p-values behave exactly as they should: Some misleading
290 criticisms of p-values and their resolution with s-values. *The American Statistician*, 73,
291 106-114.
- 292 11. Greenland, S., Senn, S., Rothman, K., Carlin, J., Poole, C., Goodman, S., & Altman, D.
293 (2016). Statistical tests, P values, confidence intervals, and power: a guide to
294 misinterpretations. *European journal of epidemiology*, 31(4), 337-350.
- 295 12. Hopkins, W. (2004). How to interpret changes in an athletic performance test. *Sportscience*
296 , 8, 1-7.
- 297 13. Lakens, D. (2017). Equivalence tests: a practical primer for t tests, correlations, and meta-
298 analyses. *Social psychological and personality science*, 8(4), 355-362.
- 299 14. Lakens, D. (2019). *What is a p-value?* Retrieved from YouTube:
300 https://www.youtube.com/watch?v=RVxHlsIw_Do
- 301 15. Lakens, D., Adolphi, F., Albers, C., Anvari, F., Apps, M., Argamon, S., & Buchanan, E.
302 (2018). Justify your alpha. . *Nature Human Behaviour*, 2, 168-171.

- 303 16. McShane, B., & Gal, D. (2017). Statistical significance and the dichotomization of
304 evidence. *Journal of the American Statistical Association*, *112*, 885-895.
- 305 17. Neyman, J., & Pearson, E. (1933). On the problem of the most efficient tests of statistical
306 hypotheses. *Philosophical Transactions of the Royal Society of London*, *231*, 289-337.
- 307 18. Rosenthal, R. (1979). The file drawer problem and tolerance for null results. *Psychological*
308 *Bulletin*, *86*, 638-641.
- 309 19. Stewart, P., Comfort, P., & Turner, A. (2017). Strength and conditioning: Coach or
310 scientist? . In A. Turner, & P. Comfort, *Advanced Strength and Conditioning: An Evidence-*
311 *based Approach*. (pp. 1 - 11). London: Routledge.
- 312 20. Traimow, D., & Marks, M. (2015). Editorial. *Basic and Applied Social Psychology*,
313 *37*, 1-2.
- 314 21. Wasserstein, R. (2019). *Misunderstandings of 'statistical significance'*. Retrieved from
315 You Tube: https://www.youtube.com/watch?v=Wu4-OK_91EM&t=271s
- 316 22. Wasserstein, R., & Lazar, N. (2016). The ASA's statement on p-values: context, process,
317 and purpose. *The American Statistician*, *70*(2), 129-133.
- 318 23. Wasserstein, R., Schirm, A., & Lazar, N. (2019). Moving to a world beyond "p< 0.05". *The*
319 *American Statistician*, *73*, 1-19.
- 320 24. Weir, J. (2005). Quantifying test-retest reliability using the intraclass correlation coefficient
321 and the SEM. *J. Strength Cond. Res*, *19*(1), 231-240.
- 322 25. Winter, E., Abt, G., & Nevill, A. (2014). Metrics of meaningfulness as opposed to sleights
323 of significance. *Journal of Sport Sciences*, *31*(10), 901-902.
- 324