

**An opinion regarding equivalence testing for
evaluating measurement agreement**

ABSTRACT

The novel statistical approach ‘equivalence testing’ has been proposed in order to statistically examine agreement between different physical activity measures. By using this method, researchers argued that it is possible to determine whether a method is significantly equivalent to another method. Recently, equivalence testing was supported with the use of 90% confidence interval, obtained from a mixed ANOVA, which I believe is a more robust approach. This paper further discusses the use of this method in comparison to a more well-established statistical analysis (i.e. mixed design ANOVA), as well as various limitations and arbitrary assumptions in order to perform this analysis. The paper concludes with some remarks and considerations for future use in similar approaches.

Keywords: Mixed design ANOVA; p-value; confidence interval; methods’ comparison

1. INTRODUCTION

I recently came across the paper of Dixon and colleagues (2018) regarding the novel statistical approach ‘equivalence testing’ in order to statistically examine agreement between different physical activity (PA) measures and to evaluate the validity of a new method. The researchers suggest that the use of standard statistical tests of mean differences (e.g. ANOVA, t-test) is employed in similar research approaches, which generally focus in

21 significant differences, rather than equivalences. I have to mention that this test is proposed
22 only for group-level measurement agreement, because for individual-level agreement other
23 tests are widely accepted as more adequate and valid (i.e. Bland-Altman plots, Mean
24 Absolute Percentage Error and Root Mean Square Error of Approximation).

25

26 **2. RATIONALE**

27

28 Initially I came across this statistical technique in Lee, Kim and Welk's [1] study, with the
29 exception that there were not mentioned any p-values. They stated that 'in traditional
30 hypothesis testing, the focus is on testing for a significant difference', however by 'using an
31 equivalence test, it is possible to determine whether a method is significantly equivalent to
32 another method' [1, p. 1843]. Since then, a number of studies have used this method in
33 order to evaluate the agreement between methods in sport science [2-4].

34 In that initial approach [1], as well as some of following studies [4], equivalence was
35 supported with the use of 90% confidence interval (CI), obtained from a mixed ANOVA. I
36 believe that was a more robust approach, taking into consideration the misuse of p-values in
37 order to support statistical hypotheses [5]. In fact, the American Statistical Association
38 released specific guidelines on the use of p-values stating, among else, that p-values do not
39 measure the probability that the studied hypothesis is true, do not measure the size of an
40 effect or the importance of a result and not provide a good measure of evidence regarding a
41 model or hypothesis. For this reason, the use of methods that emphasize estimation over
42 testing was suggested, such as confidence, credibility, or prediction intervals and Bayesian
43 methods [6].

44 The basic assumption made in order to justify equivalence testing was that standard
45 statistical tests of mean differences are designed to detect differences, not equivalence and
46 failure to reject the null hypothesis of no difference does not necessarily provide evidence of
47 equivalence [7]. I am not convinced that this statement is correct. It is widely accepted that
48 hypothesis testing is an important activity of empirical research. The initial null hypothesis
49 (H_0) assumes that population means are equivalent and only if there is strong evidence to
50 the contrary (alternative hypothesis; H_A), it can be assumed that there are differences among
51 group means [8].

52 Furthermore, multivariate inferential procedures (i.e. repeated measures ANOVA) include
53 hypothesis tests that allow several variables to be studied by preserving the significance
54 level without inflating type I error rate [9]. The sample size is an issue, however with the
55 correct use of appropriate tests, such as Pillai's trace, small or unequal sample sizes are not
56 considered problematic, because the greatest protection against type I errors is offered [10].

57 Additionally, mixed-model designs are recommended in most cases because they can
58 control for the repeated nature of the data (i.e. collection of data from PA monitors for
59 multiple activities) [11]. This is not possible in equivalence testing, even though this
60 approach might have limited value, because a single regression model is fitted to the
61 average of the estimates throughout the range of all activities and not each activity
62 separately [7].

63 Lastly, I would like to add some concluding remarks. First, the confidence intervals for
64 equivalence suggested by the writers (i.e. 10% and 2%) are somewhat arbitrary, an issue
65 also highlighted by Dixon and colleagues [7]. This is acceptable, since equivalence bounds
66 in sport science are not set by regulations, as it happens for drug development (e.g.,
67 differences up to 20% are not considered to be clinically relevant). Such general regulations
68 about what constitutes a meaningful effect seem unlikely to emerge, even though these
69 could be extremely helpful and of increased value, especially in sport and exercise medicine.
70 However, these intervals remain arbitrary and no statistically-based justification has been
71 proposed in order to justify them.

72

73 **3. CONCLUDING REMARKS**

74

75 In future similar researches, I believe it would be more appropriate to compare the results
76 derived from different statistical methods (i.e. equivalence testing vs mixed design ANOVA)
77 and not only present the results from a single method. This approach could provide evidence
78 of similarities and differences between the methods, so that the readers can understand
79 more adequately what extra the new method has to offer. However, in order to correctly
80 address these results, CI and effect sizes (a set of statistics that indicates the relative
81 magnitude of the differences between means [8] should also be calculated for all methods
82 and not simply rely on p-values, as it happens nowadays with equivalence testing. Lastly, in
83 order this attempt to introduce equivalence testing in sport and exercise science to be
84 successful, the following considerations should be taken into account: a) Develop easy-to-
85 use and accessible software; b) Express equivalence bounds in standardized effect sizes
86 rather than raw scores; c) Related articles should discuss both power analyses and
87 statistical tests for dependent t-tests, repeated measures or mixed design ANOVA and meta-
88 analyses; d) Guidance should be provided on how to set equivalence boundaries, given that
89 there are often no specific theoretical limitations on how small effects are predicted to be nor
90 cost-benefit boundaries of when effects are too small to be practically meaningful [12].

91 The interesting article of Dixon and colleagues [7] adds further to our understanding
92 regarding the adequate use of equivalence testing for evaluating measurement agreement in

93 sport science. While it is exciting to see increased attention to the development and
94 dissemination of new statistical approaches, scientists should be cautious about making and
95 adopting statistical recommendations, because these could be considered as another 'trend'.
96

97 **COMPETING INTERESTS**

98

99 Authors have declared that no competing interests exist.

100

101 **REFERENCES**

102

- 103 1. Lee J-M, Kim Y, Welk GJ. Validity of consumer-based physical activity monitors. *Med Sci*
104 *Sports Exerc.* 2014;46(9):1840-48. doi:10.1249/MSS.0000000000000287
- 105 2. Florez-Pregonero A, Meckes N, Buman M, Ainsworth. Wearable monitors criterion validity
106 for energy expenditure in sedentary and light activities. *J Sport Health Sci.* 2017;6(1):103-
107 10. doi:10.1016/j.jshs.2016.10.005
- 108 3. Kim Y, Welk GJ. Criterion validity of competing accelerometry-based activity monitoring
109 devices. *Med Sci Sports Exerc.* 2015;47(11):2456-63.
110 doi:10.1249/MSS.0000000000000691
- 111 4. Morris CE, Wessel PA, Tinius RA, Schafer MA, Maples JM. Validity of activity trackers in
112 estimating energy expenditure during high-intensity functional training. *Res Q Exercise*
113 *Sport.* 2019 [Epub ahead of print]. doi:10.1080/02701367.2019.1603989
- 114 5. Baker M. Statisticians issue warning over misuse of P values. *Nature.*
115 2016;531(7593):151. doi:10.1038/nature.2016.19503
- 116 6. Wasserstein RL, Lazar NA. The ASA's statement on p-values: Context, process, and
117 purpose. *Am Stat.* 2016;70(2):129-33. doi:10.1080/00031305.2016.1154108
- 118 7. Dixon PM, Saint-Maurice PF, Kim Y, Hibbing P, Bai Y, Welk GJ. A primer on the use of
119 equivalence testing for evaluating measurement agreement," *Med Sci Sports Exerc.*
120 2018;50(4):837-45. doi:10.1249/MSS.0000000000001481
- 121 8. Tabachnick BG, Fidell LS. *Using multivariate statistics.* 5th ed. Boston, MA: Allyn & Bacon;
122 2007.
- 123 9. Rencher AC. *Methods of Multivariate Analysis.* 2nd ed. New York, NY: John Wiley & Sons,
124 Inc; 2002.
- 125 10. Park E, Cho M, K, CS. Correct use of repeated measures analysis of variance. *Korean J*
126 *Lab Med.* 2009;29:1-9. doi:10.3343/kjlm.2009.29.1.1
- 127 11. Welk GJ, McClain J, Ainsworth BE. Protocols for evaluating equivalency of
128 accelerometry-based activity monitors. *Med Sci Sports Exerc.* 2012;44(Suppl1):S39-49.
129 doi:10.1249/MSS.0b013e3182399d8f
- 130 12. Lakens D. Equivalence tests: A practical primer for t-tests, correlations, and meta-
131 analyses. *Soc Psychol Pers Sci.* 2017;8(4):355-62. doi:10.1177/1948550617697177