**_Title:_ The numbers will love you back in return – I promise**

**_Submission type:_ Invited Commentary**

**_Authors:_** _M. Buchheit._[1]

[1] Performance Department, Paris Saint-Germain Football Club, Saint-Germain-en-Laye, France

**_Running Head:_** Magnitude-based inferences

**_Contact details:_**
Martin Buchheit
Performance Department, Paris Saint-Germain Football Club,
4a avenue du président Kennedy
78100 Saint-Germain-en-Laye, France
Tel.: +33 1 61 07 10 77
E-mail: mbuchheit@psg.fr

**_Abstract word count:_ 248**

**_Text-only word count:_ 1727**

**_Number of Tables:_ 0**

**_Number of Figures:_ 3**

**_Disclosures:_** nothing to disclose

## 1. Abstract

The first sport science-oriented and comprehensive paper on magnitude-based inferences (MBI) was published 10 years ago in the first issue of this journal. While debate continues, MBI is today well-established in sports science and in other fields, particularly clinical medicine where practical/clinical significance often takes priority over statistical significance. In this commentary, some reasons why both academics and sport scientists should abandon null hypothesis significance testing (NHST) and embrace MBI are reviewed. Apparent limitations and future areas of research are also discussed. The following arguments are presented: P values and in turn, study conclusions, are sample-size dependent, irrespective of the size of the effect; significance doesn't inform on magnitude of effects, yet magnitude is what matters the most; MBI allows authors to be honest with their sample size and better acknowledge trivial effects; the examination of magnitudes *per se* helps provide better research questions; MBI can be applied to assess changes in individuals; MBI improves data visualisation; and lastly, MBI is supported by spreadsheets freely available on the internet. Finally, recommendations to define the smallest important effect and improve the presentation of standardized effects are presented.

**Keywords:** magnitude-based inferences; null hypothesis significance testing; sample size; trivial effect; smallest important effect.

## 2. Introduction

I discovered magnitude-based inferences (MBI) in 2008 while reading Impellizzeri et al.'s repeated-sprint testing paper in professional soccer.[1] I found the first figure of their paper to be simply fascinating. First, changes in repeated-sprint performance were compared in reference to a typical threshold representative of a smallest important or meaningful change (later to be termed the smallest worthwhile change, SWC[2]). Second, instead of a classical 'yes or no' type response, the authors reported both quantitatively and qualitatively the probabilities for these changes to be 'real' (Figure 1). Never had I previously read of anything more meaningful to that day. The message displayed within that figure spoke to both sport scientists and practitioners alike. In France, as in most other countries at that time, statistical lectures exclusively sang the praises of null hypothesis significance testing (NHST). The understanding of these statistics took long hours of self-driven and motivated learning. However, this new statistical approach, driven largely by Will G. Hopkins's efforts, has changed my life, both as an academic and practitioner in elite sport.[3]

Perhaps analogous to spirituality and religion, where individuals follow their own God, editors and reviewers of most journals (even those with high impact factors[4]) can find it difficult to think outside their bubble, believing only what they were taught in graduate school. Since authors driven by their H-index know that providing everything other than a P value increase dramatically their chances of seeing their paper rejected,[3] they simply stick to NHST to facilitate reviews and expedite publication. Fortunately, things have progressively moved on in some sports science journals.[5, 6] While some may view such occurrences as a coincidence, the first sport science-oriented and comprehensive paper on MBI was published 10 years ago in the first issue of our journal,[7] and remains one of the most cited papers on the topic, together with the 2009 update in another journal.[8] Somewhat unexpectedly last year, some authors claimed that MBI had questionable theoretical foundations and suffered from apparently high rates of type I errors (i.e., false positives), which lead them to advise researchers against using MBI.[9] In March this year, Hopkins and Batterhamm[10] provided evidence to dismiss the critiques and to reassure researchers and practitioners that MBI is in reality superior to NHST. While the debate will likely continue, MBI is today a well-established analytical approach in sports science and in other fields, particularly clinical medicine where practical/clinical significance often takes priority over statistical significance.

While the present work is only an invited commentary, and should not be considered as journal policy, I personally wish that MBI is influential with other scientists, as it has been to me. I take this opportunity to put forth the following recommendations, limitations and future areas of research, to assist researchers and practitioners to make better decisions with our numbers.

**Reasons why academics should abandon NHST and embrace MBI (using the probable effect of a new nutritional supplement on performance as an example).**

1. **P values and in turn, study conclusions, are sample-size dependent (the greater the n, the lower the P), irrespective of the size of the effect.** While it can be concluded that the nutritional supplement is ineffective with a sample of 12 athletes (P>0.05), the same comparison may turn useful with $n = 14$ (P<0.05). In other words, the drop-out of a few athletes, or the lucky involvement of 2 more subjects can induce a 180 degree change in a study conclusion.[11] This sample size issue explains also a large portion of the publication bias in research,[12] where only significant results tend to be submitted (among the studies with small sample size only those showing large effects –more likely significant (P<0.05)– are submitted and published).[10]

2. **Significance doesn't inform on magnitude of effects, yet magnitude is what matters the most.[13]** With a large enough sample size, even very small, trivial or non-practical effects can turn significant (P<0.05). In practice, with 200 athletes showing a 0.01% improvement in performance, NHST would suggest that the nutritional supplement works, while the effects may in fact be negligible. In my experience, coaches and athletes (and probably most of our readers

too) are first interested in knowing what kind of performance benefits may be expected from the supplement (i.e., how much, the actual magnitude), and how likely this magnitude is of practical importance (i.e., likelihood of the effect to be greater than the SWC).

3. **MBI allows authors to be honest with their sample size and better acknowledge trivial effects.** While a P>0.05 is often interpreted as a lack of an effect/difference, it is actually impossible to be confident that this is the right interpretation of the data analysis (sample size issue, Type II error resulting from low statistical power). The beauty of MBI is that it allows deciphering between clear (confidence limits within the SWC) and unclear (CL overlapping the SWC) trivial effects (Figure 1). This can't be touched by NHST. An unclear effect/difference is not to be interpreted as a lack of an effect, but suggests the need to increase sample size to improve precision.

4. **The examination of magnitudes *per se* helps provide better research questions**. Considering that the size of an effect matters more than a simple yes or no answer (NHST), typical hypotheses that do not have clear foundations (e.g., "We hypothesized that the new supplement would be beneficial for performance") can be replaced by a simpler and more relevant statement: "Our aim was to quantify the performance benefit of that supplement, if any".

5. **MBI is supported by spreadsheets freely available on the internet (e.g., [14])**

**Reasons why magnitude-based inferences are the essential statistical tool for practitioners in the field**

1. **MBI can be applied to assess changes in individuals.** In essence, conventional statistics allow analysis of population-based responses, which are impractical for monitoring performance changes in individuals (Figure 2). While individual score changes can be assessed in various ways (e.g., Z-scores,[15] standard difference score[16]), MBI additionally allow us to assess the likelihood of these changes to be true for any given athlete, once the typical error of the test of interest and the SWC are known.[17, 18]

2. **MBI improves data visualisation.** MBI principles should be applied to graphical reports produced by sport scientists, where shaded trivial areas and confidence limits (or typical errors for individual data) are presented systematically to acknowledge the fact that not all changes are worthwhile and that some uncertainty always remains (Figures 1, 2 and 3).

**An apparent limitation of MBI** is that, in contrast to NHST, researchers have to define *a priori* both the magnitude of the smallest important effect and the thresholds used to qualify likelihoods (e.g., *very likely*).[19] My view is that instead of being a limitation, this forces researchers to adopt a conscious process when analysing their data. "NHST is easy, but misleading. MBI is hard but honest" (W.G. Hopkins, personal communication). The importance of an appropriate SWC definition is often overlooked[20, 21] and may directly impact decisions: while a larger SWC may lead to more conservative decisions, a smaller SWC increases the chance of effects/differences being substantial. In fact, the most appropriate SWC is variable-dependent and based on either theoretical or practical considerations. While for individual athlete performance, a third of the performance coefficient of variation (CV) is generally suggested, and a fifth of the between-athlete SD is often used for performance variables in team sports.[17] A limitation however of using the SD for standardization is that the SWC may be affected by group homogeneity; for that reason, performance clues may be sometimes used instead, e.g., based on empirical observations of direct performance benefits, such as a distance of 20–50 cm that one soccer player needs to be ahead of the opponent to win a ball, corresponding to a 1% improvement in 20-m sprint time.[22] For physiological data with no direct link to performance (e.g., heart rate variability), using multiples of the within-athlete SD is a relevant option. In contrast, when an association with performance can be established for a physiological variable (i.e., submaximal HR), the actual change in this variable

145    that relates to the smallest important change in performance is often preferred.[23] There are some
146    variables however for which the most appropriate SWC remains to be determined. For match running
147    performance data in team sports for example, which are neither related to actual physical capacities nor
148    match outcomes,[24] using the between-athlete SD is questionable, but using within-athlete variation is
149    not easy either. In fact, the magnitude of within-athlete variations may depend on both tracking variables
150    and intensity zones.[25]

151    **How standardized changes/differences are presented is crucial for a better understanding of**
152    **magnitudes.** While percentages are commonly used to report changes/differences both in research and
153    field practice, there are no clear thresholds to interpret their magnitudes, and they often bias the
154    comparison of variables that differ in units[26] (e.g., in terms of athlete trainability, while a 3% increase
155    in sprinting speed may be considered remarkable,[22] the same improvement in maximal oxygen update
156    may be relatively negligible). For these reasons, using Cohen's effect size principle (d) is generally the
157    first step toward standardization (Figure 3).[27] However, if we consider that the actual method of SWC
158    determination may be variable-dependent (Cohen's d vs. within-athletes CV vs. performance clues), the
159    same approach could be applied to standardize the changes in different variables. The thresholds for
160    small, moderate, large and very large standardized changes (Cohen's d) being 0.2, 0.6, 1.2 and 2,
161    respectively, means that any change of 1x, 3x, 6x and 10x SWC can be considered as small, moderate,
162    large and very large, respectively (Figure 3). Reporting effects/changes as multiples of the SWC[28] is
163    relevant for at least two reasons: i) in manuscripts, the changes in all variables can be easily aggregated
164    into a single figure with a single shaded trivial area (Figure 3) and ii) for coaches and athletes, the
165    message cannot be simpler than: "the effect is *x* times greater than what generally matters to you guys".

166

167    **3.  Conclusion.**

168    The introduction of MBI into sports science nearly 15 ago represents one of the most important
169    analytical progressions in our field. While there are still areas that need to be developed, there is no
170    doubt that we should all be leaning toward a more mature and conscious process of analysing and
171    presenting our data.[19] "The numbers are where the discussion should start, not end."[29]

172

179 **References**

180

181　1.　Impellizzeri, F.M., et al., Validity of a Repeated-Sprint Test for Football. *Int J Sports Med*,
182　　　2008;29:899-905.
183　2.　Hopkins, W.G., Statistical vs clinical or practical significance. *Sportscience*
184　　　2002;6:http://www.sportsci.org/jour/0201/Statistical_vs_clinical.ppt.
185　3.　Buchheit, M. *Any Comments?* 2013; Available from:
186　　　https://herearemycomments.wordpress.com/category/best-hopeless-comments.
187　4.　Tressoldi, P.E., et al., High impact = high statistical standards? Not necessarily so. *PLoS One*,
188　　　2013;8(2):e56180.
189　5.　Atkinson, G., A.M. Batterham, and W.G. Hopkins, Sports performance research under the
190　　　spotlight. *Int J Sports Med*, 2012;33(12):949.
191　6.　Winter, E.M., G.A. Abt, and A.M. Nevill, Metrics of meaningfulness as opposed to sleights of
192　　　significance. *J Sports Sci*, 2014;32(10):901-2.
193　7.　Batterham, A.M. and W.G. Hopkins, Making meaningful inferences about magnitudes. *Int J*
194　　　*Sports Physiol Perform*, 2006;1(1):50-7.
195　8.　Hopkins, W.G., et al., Progressive statistics for studies in sports medicine and exercise
196　　　science. *Med Sci Sports Exerc*, 2009;41(1):3-13.
197　9.　Welsh, A.H. and E.J. Knight, "Magnitude-based inference": a statistical review. *Med Sci Sports*
198　　　*Exerc*, 2015;47(4):874-84.
199　10.　Hopkins, W.G. and A.M. Batterham, Error Rates, Decisive Outcomes and Publication Bias with
200　　　Several Inferential Methods. *Sports Med*, 2016.
201　11.　McCormack, J., B. Vandermeer, and G.M. Allan, How confidence intervals become confusion
202　　　intervals. *BMC Med Res Methodol*, 2013;13:134.
203　12.　Kuhberger, A., A. Fritz, and T. Scherndl, Publication bias in psychology: a diagnosis based on
204　　　the correlation between effect size and sample size. *PLoS One*, 2014;9(9):e105825.
205　13.　Cohen, J., Things I have learned (so far). *American Psychologist*, 1994;45:1304-1312.
206　14.　Hopkins, W.G. *A spreadsheet for deriving a confidence interval, mechanistic inference and*
207　　　*clinical inference from a P value.* Sportscience, 2007. 11, 16-20. DOI:
208　　　http://newstats.org/xcl.xls.
209　15.　McGuigan, M.R., S.J. Cormack, and N.D. Gill, Strength and Power Profiling of Athletes:
210　　　Selecting Tests and How to Use the Information for Program Design. *Strength and*
211　　　*Conditioning Journal*, 2013;37(6):7-14.
212　16.　Pettitt, R.W., The standard difference score: a new statistic for evaluating strength and
213　　　conditioning programs. *J Strength Cond Res*, 2010;24(1):287-91.
214　17.　Hopkins, W.G., How to Interpret Changes in an Athletic Performance Test. *Sportscience*,
215　　　2004; 8:1-7.
216　18.　Al Haddad, H., B.M. Simpson, and M. Buchheit, Monitoring changes in jump and sprint
217　　　performance: best or average values? *Int J Sports Physiol Perform*, 2015;10(7):931-4.
218　19.　Schaik, P.V. and M. Weston, Magnitude-based inference and its application in user research.
219　　　*International Journal of Human-Computer Studies*, 2016(88):38-50.
220　20.　Atkinson, G., Does size matter for sports performance researchers? *J Sports Sci*,
221　　　2003;21(2):73-4.
222　21.　Buchheit, M., A. Rabbani, and H.T. Beigi, Predicting changes in high-intensity intermittent
223　　　running performance with acute responses to short jump rope workouts in children. *J Sports*
224　　　*Sci Med*, 2014;13(3):476-82.
225　22.　Haugen, T. and M. Buchheit, Sprint Running Performance Monitoring: Methodological and
226　　　Practical Considerations. *Sports Med*, 2015.
227　23.　Buchheit, M., Monitoring training status with HR measures: do all roads lead to Rome? *Front*
228　　　*Physiol*, 2014;27(5):73.

229  24.  Mendez-Villanueva, A. and M. Buchheit, Physical capacity-match physical performance
230       relationships in soccer: simply, more complex. *Eur J Appl Physiol*, 2011;111(9):2387-9.
231  25.  Buchheit, M., et al., Integrating different tracking systems in football: multiple camera semi-
232       automatic system, local position measurement and GPS technologies. *J Sports Sci*,
233       2014;32(20)(20):1844-1857.
234  26.  Buchheit, M. and A. Rabbani, 30-15 Intermittent Fitness Test vs. Yo-Yo Intermittent Recovery
235       Test Level 1: Relationship and Sensitivity to Training. *Int J Sports Physiol Perform*,
236       2014;9(3):522-524.
237  27.  Cohen, J., *Statistical Power Analysis for the Behavioral Sciences*. 1988, Hillsdale: Lawrence
238       Erlbaum. 599.
239  28.  Buchheit, M., et al., Physiological, psychometric, and performance effects of the Christmas
240       break in Australian football. *Int J Sports Physiol Perform*, 2015;10(1):120-3.
241  29.  Nuzzo, R., Scientific method: statistical errors. *Nature*, 2014;506(7487):150-2.
242  30.  Buchheit, M. and A. Mendez-Villanueva, Effects of age, maturity and body dimensions on
243       match running performance in highly trained under-15 soccer players. *J Sports Sci*,
244       2014;32(13):1271-8.

245

246

247

248    **Figures Legends**

249

250    **Figure 1.** Example of possible decisions when interpreting changes using magnitude-based inferences.
251    Note the clear vs. unclear cases (based on confidence limits, in relation to the shaded trivial area),
252    which i) is one of the extreme beauty of magnitude-based inferences and ii) provide no insight through
253    null hypothesis significance testing. Note also how, for clear effects, the likelihood of changes
254    increases as the confidence limits shrink.

255    **Figure 2.** Individual changes in submaximal heart rate in a professional soccer player when running at
256    12 km/h throughout 1.5 competitive seasons (% of maximal heart rate). The shaded area represents
257    trivial changes (1%).[23] The error bars represent the typical error of measurement (3%).[23] The numbers
258    of * indicate the likelihood for the changes to be substantial, with 1 symbols referring to possible
259    changes, 2 to likely, 3 to very likely and 4 to almost certain changes.

260    **Figure 3.** Differences in various anthropometric, physiological and performance measures between
261    two groups of young soccer players differing by their maturity status (0.9 ± 0.3 vs. -0.2 ± 0.4 years
262    from predicted peak height velocity)[30] when expressed in percentages (A), using Cohen's effect size
263    principle (B) and as a factor of variable-specific smallest worthwhile differences (SWD) (C):[28] 0.2 x
264    between-athletes SD for height, MAS and matches tracking data; performance-related changes for
265    HRR and MSS (7[23] and 2[22]%, respectively). The numbers of * indicate the likelihood for the between-
266    group differences to be substantial, with 1 symbols referring to possible difference, 2 to likely, 3 to
267    very likely and 4 to almost certain differences. Note that that magnitude of the between-group
268    differences and their likelihood varies between the panels. My suggestion is to use the method used in
269    panel C (with a variable-specific SWD). MSS: maximal sprinting speed, MAS: maximal aerobic
270    speed, HRR: heart rate recovery after submaximal exercise, D>16 km/h: distance ran above 16 km/h
271    during matches, #HIR: number of high intensity runs during matches.

272