1                                          Effect size

2                    Marie Delacre[1], Christophe Leys[1], Limin Liu[2], & Daniël Lakens[3]

3      [1] Université Libre de Bruxelles, Service of Analysis of the Data (SAD), Bruxelles, Belgium

4                                      [2] Université de Gant

5      [3] Eindhoven University of Technology, Human Technology Interaction Group, Eindhoven,

6                                       the Netherlands

7                                          Author Note

8          Correspondence concerning this article should be addressed to Marie Delacre, CP191,

9      avenue F.D. Roosevelt 50, 1050 Bruxelles. E-mail: marie.delacre@ulb.ac.be

# Abstract

*Keywords:* keywords

Word count: X

14                                                Effect size

15                                                  **Intro**

16      During decades, researchers in social science (Henson & Smith, 2000) and education

17   (Fan, 2001) have overestimated the ability of the null hypothesis (H0) testing to determine

18   the importance of their results. The standard for researchers in social science is to define H0

19   as the absence of effect (Meehl, 1990). For example, when comparing the mean of two

20   groups, researchers commonly test the H0 that there is no mean differences between groups

21   (Steyn, 2000). Any effect that is significantly different from zero will be seen as sole support

22   for a theory.

23      Such an approach has faced many criticisms among which the most relevant to our

24   concern is that the null hypothesis testing highly depends on sample size: for a given alpha

25   level and a given difference between groups, the larger the sample size, the higher the

26   probability of rejecting the null hypothesis (Fan, 2001; Kirk, 2009; Olejnik & Algina, 2000;

27   Sullivan & Feinn, 2012). It implies that even tiny differences could be detected as

28   statistically significant with very large sample sizes (McBride, Loftis, & Adkins, 1993)[1].

29      Facing this argument, it has become an adviced practice to report the $p$-value assorted

30   by a measure of the effect size, that is, a quantitative measure of the magnitude of the

31   experimenter effect (Cohen, 1965; Fan, 2001; Hays, 1963). This practice is also highly

32   endorsed by the *APA Publication Manual* (Association, 2010). However, we will argue that

33   very often, the contribution of the measures of effect size is misunderstood as a measure of

34   "the importance of an effect in real life" while it is not. Moreover, there are several situations

---

[1] This is especially problematic since these tiny differences might be due to other factors than the one of

interest: even under the assumption of random assignent (which is a necessary but not sufficient condition),

it is almost impossible to be sure that the only difference between two conditions is the one defined by the

factor of interest. Other tiny factors of no theoretical interest might slightly influence results, making the

probability of getting an actual zero effect very low. This is what Meehl (1990) calls 'systematic noise'

that call for effect size measures and in the current litterature, it's not always easy to know which measure using in which circumstances. Depending on the purpose of using a measure of effect size, different indicators of effect size have their own strenghts and weaknesses. Lastly, when associated with interential tests, the main measures of effect sizes are submitted to a range of assumptions that are unrealistic in many research designs. As consequences many estimations of effect size are inaccurate and alter the robustness of the statistical conclusions.

In sum the aim of this paper is fourfold: (1) Discuss the different levels of significance, following Bothe (2011)'s nomenclature, and distinguish their domain of validity; (2) Identify the different situations that justify the use of effect size; (3) Identify the strenghts and weaknesses of several measures of effect sizes in these situations; (4) Discuss the problem of the effect size assumptions and the impact of these assumptions violations on the robustness of the measures of effect size and the statistical conclusions.

**Dans cette structure, il va falloir expliquer les qualités requises d'une mesure de taille d'effet.**

#Levels of Significance

In their paper, Bothe (2011) distinguish between three levels of significance, namely Statistical significance, Practical significance and Clinical significance (with the adjunction of Personal significance). Statistical significance refers to the $p$-value. As stated before, this conclusion is highly dependent from the sample size. Practical significance refers to the magnitude of a change or a difference between groups. In other terms, it is any statistical indicator that assess mathematically the effect size. Laslty, the Clinical significance refers to the interpretation of treatment outcomes. This last level is not statistical nor mathematical, it is related to underlying theory that posits an empirical hypothesis.

It is important to understand the difference between these three concepts. Statistical

significance allows the researcher to determine whether the oberved departure from H0 can be attributed to something else than randomness (i.e. an actual effect). Practical significance is a mathematical indicator of effect size that is not necessarily related to the theoretical effect or at least that the relation is not straightforward. As stated by Kazdin « . . . clinical significance has been defined as whether an intervention "makes a real (e.g., genuine, noticeable) difference in everyday life to the clients or to others with whom the clients interact"" (Kazdin, 1999, p.332 ; cité par Bothe (2011)). Indeed, Pratical significance depends on the way a variable is converted into numerical indicator. For example, when assessing Self-Compassion, one can use a scale such as the Self-Compassion Scale (Kotsou & Leys, 2016; Neff, 2003). This scale inform the researcher about the level of self-compassion based on a ordinal scale that can yield different value depending on the influence of any independent variable. For example, some training program can improve subjects level of self-compassion (Jazaieri et al. (2013)). Yet, since the scale is ordinal, meaning that there is no standard unit to assess the construct, the relation between the mathematical effect size (i.e. Pratical significance) and the actual change in self-compassion (i.e. Clinical significance) will always remain unknown. Therefore, although, as we will see, pratcial significance is important to determine, it's relation with clinical significance has often to be addressed, and that is more a theoretical argument than a statistical one.

To further distinguish between important constructs, the authors suggest incorporating as definitive the existing notion that clinical significance may refer to measures selected or interpreted by professionals or with respect to groups of clients. The term personal significance is introduced to refer to goals, variables, measures, and changes that are of demonstrated value to individual clients.

AS a conclusion, *statistically significant* effect is not necessarily of *practical* interest. The *statistical* significance is the probability that findings have occured by chance (Stout & Ruble, 1995). The *practical* significance is the magnitude of findings and is assessed by

86 measures of **effect sizes**.

87      At the same time, a vast literature has developed that casts doubt on the credibility of

88 the assumptions of Student's *t*-test and classical *F*-test ANOVA (i.e. the assumptions that

89 two or more samples are independent, and that independent and identically distributed

90 residuals are normal and have equal variances between groups; Glass, Peckham, & Sanders,

91 1972) (CITER TOUTES MES REFERENCES). In a previous paper, We focused on the

92 assumptions of normality and equality of variances, and argued that these assumptions are

93 often unrealistic in the field of psychology. Bcp d'autres chercheurs avant nous étaient

94 arrivés à la même conclusion. Pourtant, beaucoup moins d'auteurs se sont penchés sur les

95 mesures de taille d'effet à utiliser en complément du test de welch. Il existe de la littérature

96 sur la question, mais pas vraiment d'accord (parce que grande confusion quant à la questino

97 suivante: à quoi sert la mesure de taille d'effet? ) Par ailleurs, s'il est de plus en plus

98 communément admis que les conditions d'application des tests de comparaison de moyennes

99 (dominant toujours la recherche) sont peu réalistes et rarement respectées, pourtant et que

100 de nombreux chercheurs recommandent d'utiliser le Welch au lieu du test de Student, peu de

101 littérature suggère quelle taille d'effet associer à ce test. Même Jamovi ne propose comme

102 mesure de taille d'effet que le d de Cohen, souffrant des mêmes limites que le test de Student.

103      Pour cette raison, nous proposons de structurer cet article comme suit: # 1) Bien

104 definir practical significance (donc donner une définition claire de la taille d'effet qui nous

105 convient) Expliquer un peu pourquoi c'est important d'avoir l'IC autour de l'effect size: 1)

106 Parce que l'estimation dépend du n (plus n est grand, plus précise est l'estimation) 2) parce

107 que la mesure de taille d'effet est un complément de la significativité statistique: comme le

108 dit

**2) Bien définir à quel objectif on tente de répondre via la mesure de taille d'effet (je les cite tous dans mon pwp)**

**3) Qualités MATHEMATISUES importantes d'une bonne mesure de taille d'effet et de l'IC**

**4) Revue sur les familles de tailles d'effet (r et d, et mesures les plus connues)**

**5) Simulations**

Association. (2010). *Publication manual of the american psychological association [apa] (6 ed.)* (American Psychological Association.). Washington, DC:

Bothe, R., A. K. (2011). Statistical, practical, clinical, and personal significance: Definitions and applications in speech-language pathology. *American Journal of Speech-Language Pathology, 20*, 233–242.

Cohen, J. (1965). Some statistical issues in psychological research. In *Handbook of clinical psychology* (B. B. Wolman., pp. 95–121). New York: McGraw-Hill.

Fan, X. (2001). Statistical significance and effect size in education research: Two sides of a coin. *Journal of Educational Research, 94*(5), 275–282. doi:10.1080/00220670109598763

Glass, G. V., Peckham, P. D., & Sanders, J. R. (1972). Consequences of failure to meet assumptions underlying the fixed effects analyses of variance and covariance. *Review of Educational Research, 42*(3), 237–288. doi:10.3102/00346543042003237

Hays, W. L. (1963). *Statistics for psychologists* (Holt, Rinehart & Winston.). New York.

Henson, R. I., & Smith, A. D. (2000). State of the art in statistical significance and effect size reporting: A review of the APA task force report and current trends. *Journal of Research and Development in Education, 33*(4), 285–296.

<sup>132</sup> Jazaieri, H., Jinpa, G. T., McGonigal, K., Rosenberg, E. L., Finkelstein, J.,

<sup>133</sup> Simon-Thomas, E., ... Goldin, P. R. (2013). Enhancing compassion: A randomized

<sup>134</sup> controlled trial of a compassion cultivation training program. *Journal of Happiness Studies*,

<sup>135</sup> *14*(4).

<sup>136</sup> Kirk, R. E. (2009). Practical significance: A concept whose time has come. *Educational*

<sup>137</sup> *and Psychological Measurement*, *56*(5), 746–759. doi:10.1177/0013164496056005002

<sup>138</sup> Kotsou, I., & Leys, C. (2016). Self-compassion scale (scs): Psychometric properties of

<sup>139</sup> the french translation and its relations with psychological well-being, affect and depression.

<sup>140</sup> *PloS one*, *11*(4). doi:10.1371/journal.pone.0152880

<sup>141</sup> McBride, G. B., Loftis, J. C., & Adkins, N. C. (1993). What do significance tests really

<sup>142</sup> tell us about the environment? *Environmental Management*, *17*(4), 423–432.

<sup>143</sup> Meehl, P. E. (1990). Appraising and amending theories: The strategy of Lakatosian

<sup>144</sup> defense and two principles that warrant it. *Psychological Inquiry*, *1*(2), 108–141.

<sup>145</sup> Neff, K. D. (2003). The development and validation of a scale to measure

<sup>146</sup> self-compassion. *Self and Identity*, *2*(3), 223–250. doi:10.1080/15298860309027

<sup>147</sup> Olejnik, S., & Algina, J. (2000). Measures of effect size for comparative studies:

<sup>148</sup> Applications, interpretations, and limitations. *Contemporary Educational Psychology*, *25*,

<sup>149</sup> 241–286. doi:10.1006/ceps.2000.1040

<sup>150</sup> Steyn, H. S. (2000). Practical significance of the difference in means. *Journal of*

<sup>151</sup> *Industrial Psychology*, *26*(3), 1–3.

<sup>152</sup> Stout, D. D., & Ruble, T. L. (1995). Assessing the practical signfance of empirical

<sup>153</sup> results in accounting education research: The use of effect size information. *Journal of*

<sup>154</sup> *Accounting Education*, *13*(3), 281–298.

155    Sullivan, G., & Feinn, R. (2012). Using effect size—or why the p value is not enough.

156    *Journal of Graduate Medical Education*, 279–282. doi:10.4300/JGME-D-12-00156.1