

1 What measure of effect size when comparing two groups based on their means?

2 Marie Delacre¹, Christophe Leys¹, Limin Liu², & Daniël Lakens³

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

10

Abstract

11

12 *Keywords:* keywords

13 Word count: X

What measure of effect size when comparing two groups based on their means?

Intro

During decades, researchers in social science (Henson & Smith, 2000) and education (Fan, 2001) have overestimated the ability of the null hypothesis (H0) testing to determine the importance of their results. The standard for researchers in social science is to define H0 as the absence of effect (Meehl, 1990). For example, when comparing the mean of two groups, researchers commonly test the H0 that there is no mean differences between groups (Steyn, 2000). Any effect that is significantly different from zero will be seen as sole support for a theory.

Such an approach has faced many criticisms among which the most relevant to our concern is that the null hypothesis testing highly depends on sample size: for a given alpha level and a given difference between groups, the larger the sample size, the higher the probability of rejecting the null hypothesis (Fan, 2001; Kirk, 2009; Olejnik & Algina, 2000; Sullivan & Feinn, 2012). It implies that even tiny differences could be detected as statistically significant with very large sample sizes (McBride, Loftis, & Adkins, 1993)¹.

Facing this argument, it has become an advised practice to report the p -value assorted by a measure of the effect size, that is, a quantitative measure of the magnitude of the experimenter effect (Cohen, 1965; Fan, 2001; Hays, 1963). This practice is also highly endorsed by the *APA Publication Manual* (Association, 2010). However, limited studies properly report effect size in the last several decades.

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 assignment (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'

First, there is a high confusion between the effect size and other related concept such as the clinical significance of a result (i.e. the relevance of an effect in real life). Moreover, there are several situations that call for effect size measures and in the current literature, it's not always easy to know which measure using in which circumstances.

Second, when associated with inferential 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 the context of comparing two groups based on their means, Cohen's d_s is the dominant effect size measure used by researchers (Peng, Chen, Chiang, & Chiang, 2013). We will argue that, like Student's t -test, this measure rely on the often untenable assumptions of normality and homogeneity of variances.

In sum the aim of this paper is threefold: 1. Clearly define what is (and what is not) a measure of effect size; 2. Listing the different situations that call for effect sizes measure and reviewing which measure is appropriate in which circumstance; 3. Define different properties of a good effect size estimator and discuss the impact of assumptions violations on the robustness of the measures of effect size.

Levels of Significance

Measures of effect sizes aim at communicating the **practical** significance of an effect. It refers to the *magnitude* of the difference between distributions, groups, means... (Bothe, 2011). Voir le word, essayer d'écrire un truc clair.

very often, the contribution of the measures of effect size is misunderstood as a measure of "the importance of an effect in real life" while it is not.

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. Lastly, 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 observed departure from H_0 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, Practical 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 informs the researcher about the level of self-compassion based on an ordinal scale that can yield different values 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. Practical significance) and the actual change in self-compassion (i.e. Clinical significance) will always remain unknown. Therefore, although, as we will see, practical significance is important to determine, its 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 occurred by chance (Stout & Ruble, 1995). The practical significance is the magnitude of findings and is assessed by measures of effect sizes.

At the same time, a vast literature has developed that casts doubt on the credibility of the assumptions of Student's t -test and classical F -test ANOVA (i.e. the assumptions that two or more samples are independent, and that independent and identically distributed residuals are normal and have equal variances between groups; Glass, Peckham, & Sanders, 1972) (CITER TOUTES MES REFERENCES). In a previous paper, We focused on the assumptions of normality and equality of variances, and argued that these assumptions are often unrealistic in the field of psychology. Bcp d'autres chercheurs avant nous étaient arrivés à la même conclusion. Pourtant, beaucoup moins d'auteurs se sont penchés sur les mesures de taille d'effet à utiliser en complément du test de welch. Il existe de la littérature sur la question, mais pas vraiment d'accord (parce que grande confusion quant à la question suivante: à quoi sert la mesure de taille d'effet?) Par ailleurs, s'il est de plus en plus communément admis que les conditions d'application des tests de comparaison de moyennes (dominant toujours la recherche) sont peu réalistes et rarement respectées, pourtant et que de nombreux chercheurs recommandent d'utiliser le Welch au lieu du test de Student, peu de littérature suggère quelle taille d'effet associer à ce test. Même Jamovi ne propose comme mesure de taille d'effet que le d de Cohen, souffrant des mêmes limites que le test de Student.

Pour cette raison, nous proposons de structurer cet article comme suit: # 1) Bien définir practical significance (donc donner une définition claire de la taille d'effet qui nous convient) Expliquer un peu pourquoi c'est important d'avoir l'IC autour de l'effect size: 1) Parce que l'estimation dépend du n (plus n est grand, plus précise est l'estimation) 2) parce que la mesure de taille d'effet est un complément de la significativité statistique: comme le 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.

Kirk, R. E. (2009). Practical significance: A concept whose time has come. *Educational and Psychological Measurement*, 56(5), 746–759. doi:10.1177/0013164496056005002

McBride, G. B., Loftis, J. C., & Adkins, N. C. (1993). What do significance tests really tell us about the environment? *Environmental Management*, 17(4), 423–432.

Meehl, P. E. (1990). Appraising and amending theories: The strategy of Lakatosian defense and two principles that warrant it. *Psychological Inquiry*, 1(2), 108–141.

Olejnik, S., & Algina, J. (2000). Measures of effect size for comparative studies: Applications, interpretations, and limitations. *Contemporary Educational Psychology*, 25, 241–286. doi:10.1006/ceps.2000.1040

Peng, C.-Y., Chen, L.-T., Chiang, H.-M., & Chiang, Y.-C. (2013). The Impact of APA and AERA Guidelines on Effect size Reporting. *Contemporary Educational Psychology*, 82(1), 22–50. doi:10.1080/00220973.2012.745471

Steyn, H. S. (2000). Practical significance of the difference in means. *Journal of Industrial Psychology*, 26(3), 1–3.

Sullivan, G., & Feinn, R. (2012). Using effect size—or why the p value is not enough. *Journal of Graduate Medical Education*, 279–282. doi:10.4300/JGME-D-12-00156.1