



HAL
open science

A Mean Difference is an Effect Size

Pierre Dragicevic

► **To cite this version:**

Pierre Dragicevic. A Mean Difference is an Effect Size. [Research Report] RR-9354, Inria Saclay Ile de France. 2020, pp.1-12. hal-02905210

HAL Id: hal-02905210

<https://inria.hal.science/hal-02905210v1>

Submitted on 23 Jul 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



A Mean Difference is an Effect Size

Pierre Dragicevic

**RESEARCH
REPORT**

N° 9354

July 2020

Project-Team Aviz

ISRN INRIA/RR--9354--FR+ENG

ISSN 0249-6399



A Mean Difference is an Effect Size

Pierre Dragicevic

Project-Team Aviz

Research Report n° 9354 — July 2020 — 12 pages

Abstract: Methodologists urge us to report effect sizes, but rarely explain what they mean by “effect size”. This can lead to counterproductive disputes about terminology. There is a narrow sense and a broad sense of the term “effect size”. The narrow sense refers to a family of unitless measures such as Cohen’s d , while the broad sense refers to any measure of interest, such as a mean difference in completion time expressed in seconds. Researchers in meta-analysis often use the narrow sense, while methodologists focusing on transparency in reporting generally prefer the broad sense. Researchers from the first group sometimes claim that those from the second group are misusing the term. They are not. The broad sense is older than the narrow one and even Jacob Cohen, who co-founded meta-analysis and popularized the term “effect size”, defined it broadly. It is OK to call a mean difference an effect size. When necessary, the term “effect size” can be easily made crisper with the widely-used qualifiers “standardized” and “unstandardized” (or “simple”).

Key-words: methodology, effect size, statistics, terminology, HCI.

**RESEARCH CENTRE
SACLAY – ÎLE-DE-FRANCE**

1 rue Honoré d'Estienne d'Orves
Bâtiment Alan Turing
Campus de l'École Polytechnique
91120 Palaiseau

Une différence de moyennes est une taille d'effet

Résumé : Les méthodologistes nous encouragent à reporter les tailles d'effet, mais expliquent rarement ce qu'ils entendent par "taille d'effet", ce qui peut conduire à des disputes contre-productives sur la terminologie. Il existe un sens étroit et un sens large de l'expression "taille d'effet". Le sens étroit fait référence à une famille de mesures sans unité, comme le d de Cohen, tandis que le sens large fait référence à toute mesure intéressante pour le chercheur, comme une différence moyenne de temps d'exécution exprimée en secondes. Les chercheurs en méta-analyse utilisent souvent le sens étroit, tandis que les méthodologistes promouvant la communication transparente préfèrent généralement le sens large. Les chercheurs du premier groupe affirment parfois que ceux du second groupe font un mauvais usage du terme. Ce n'est pas le cas. Le sens large est plus ancien que le sens étroit et même Jacob Cohen, qui a co-fondé la méta-analyse et popularisé le terme "taille d'effet", l'a défini de manière large. Il est correct d'appeler une différence de moyennes une taille d'effet. Si nécessaire, il est facile de rendre le terme "taille de l'effet" plus précis avec les qualificatifs très répandus "normalisée" et "non normalisée" (ou "simple").

Mots-clés : méthodologie, taille d'effet, statistiques, terminologie, IHM.

1 Introduction

This article, which first appeared as a [blog post](#) in July 2018, was prompted by a [meeting](#) we ran at the [CHI 2018 conference](#) to collect feedback on the current draft of the [transparent statistics guidelines](#). This draft has an FAQ and exemplar on effect sizes. During the meeting, a participant strongly objected to our use of the term “effect size” in the guidelines. This prompted me to investigate further to make sure we haven’t missed anything in deciding the terminology. Here is what I found.

2 We should report effect sizes. But what are effect sizes?

We are repeatedly told that effect sizes are important, and many methodologists urge us to report effect sizes. For example, [Ron Wasserstein](#), the executive director of the American Statistical Association, stated that:

In the post $p < 0.05$ era, scientific argumentation is not based on whether a p -value is small enough or not. Attention is paid to effect sizes and confidence intervals. Evidence is thought of as being continuous rather than some sort of dichotomy. [McCook, 2016]

Unfortunately, methodologists rarely explain what they exactly mean by “effect size”. The current draft of the transparent statistics guidelines ([section 2.1, Effect Size FAQ](#)) mentions that there is a **narrow sense** and a **broad sense** of the term “effect size”. Briefly, the narrow sense refers to a family of standardized (i.e., unitless) measures such as Cohen’s d , while the broad sense refers to any measure of interest, standardized or not. This includes simple and familiar metrics like **unstandardized mean differences**, e.g., a mean difference in completion time between two techniques, expressed in seconds.

The guidelines currently use “effect size” in a broad sense, and often mentions unstandardized mean differences as an example. At the CHI meeting, a participant who was manifestly knowledgeable about effect sizes strongly objected to this. If I recall correctly, the reasons were: 1) statisticians do not use the term “effect size” in that broad sense, 2) unstandardized mean differences are problematic, and calling them “effect sizes” may encourage HCI researchers to abuse them. There are two distinct questions here:

- a. Is it appropriate to use the term “effect size” to mean things like unstandardized mean differences?
- b. Which measure of effect size should be reported?

Both questions are important, but in this article, I will only focus on a).

3 Sources currently cited in the guidelines

The Effect Size FAQ provides a few references to justify its current use of the term “effect size”. The first reference is from [Geoff Cumming](#), a researcher in statistical cognition and prominent methodologist. Here are two excerpts from his 2013 book:

An effect is anything we might be interested in, and an effect size is simply the size of anything that may be of interest. [Cumming, 2013, p. 34]

[an effect size] can be as familiar as a mean, a difference between means, a percentage, a median, or a correlation. It may be a standardized value, such as Cohen's d (more on this later), or a regression coefficient, path coefficient, odds ratio, or percentage of variance explained. [Cumming, 2013, p. 38]

The second reference is from [Leland Wilkinson](#), who is a statistician, and also a prominent methodologist. Here is a quote from a paper he co-authored with a “Task Force on Statistical Inference” commissioned by the American Psychological Association (APA) in 1999:

Always present effect sizes for primary outcomes. If the units of measurement are meaningful on a practical level (e.g., number of cigarettes smoked per day), then we usually prefer an unstandardized measure (regression coefficient or mean difference) to a standardized measure (r or d). [Wilkinson, 1999, p. 599]

The Effect Size FAQ additionally cites Thomas Baguley, a psychology Professor and author of a solid introductory statistics book entitled “Serious Stats” [Baguley, 2012]. Baguley also employs the term “effect size” in a broad sense, and refers to unstandardized effect sizes as “simple effect sizes”:

A straightforward way to report the magnitude of an effect is to use the original units of measurement. [...] Such effect sizes are frequently labeled as raw or unstandardized effect sizes – terms that might suggest these metrics are inferior in some way. To avoid this suggestion I adopt the more neutral term simple effect size (see Frick, 1999). [Baguley, 2012, p. 239]

A standardized measure of effect is one which has been scaled in terms of the variability of the sample or population from which the measure was taken. In contrast, simple effect size (Frick, 1999) is unstandardized and expressed in the original units of analysis. [Baguley, 2009, p. 604].

I will get back to (Frick, 1999) later on. For now, I should point out that the current FAQ also cites an article by Peter Cummings, an epidemiology Professor and methodologist, who uses the term “effect size” to refer to standardized mean differences only:

For a continuous outcome, some researchers estimate the difference in the mean outcome values of 2 groups, such as a treated group and a control group, and divide that difference by the standard deviation (SD) of the outcome values; this converts the estimated effect to SD units. This has been called a standardized mean difference or effect size, and it has 3 variations. [Cummings, 2011, p. 592]

4 “Effect size” is often used in a broad sense

Following the CHI meeting, I dived deeper into the literature. I found that several general sources use “effect size” in a broad sense, including Wikipedia:

The term effect size can refer to a standardized measure of effect (such as r , Cohen's d , or the odds ratio), or to an unstandardized measure (e.g., the difference between group means or the unstandardized regression coefficients). [Wikipedia contributors, 2018]

and the latest edition of the APA Publication Manual:

Effect sizes may be expressed in the original units (e.g., the mean number of questions answered correctly; kg/month for a regression slope) and are often most easily understood when reported in original units. It can often be valuable to report an effect size not only in original units but also in some standardized or units-free unit (e.g., as a Cohen's d value) or a standardized regression weight. [American Psychological Association, 2010, p. 34]

The terms “unstandardized effect size” and “standardized effect size” are commonly used in the literature (see also [Richardson, 1996, Hentschke and Stüttgen, 2011, Kampenes et al., 2007]), and this automatically implies a broad definition of effect size.

5 “Effect size” is also often used in a narrow sense

I also found a number of sources that use “effect size” in a narrow sense. For the Cambridge Dictionary of Statistics, an effect size is a standardized mean difference:

Effect size: Most commonly the difference between the control group and experimental group population means of a response variable divided by the assumed common population standard deviation. Estimated by the difference of the sample means in the two groups divided by a pooled estimate of the assumed common standard deviation. Often used in meta-analysis. See also counternull-value. [Everitt and Skrondal, 2010, p. 148]

Other articles employ a similar definition:

In most ecological applications of meta-analysis to date, effect size has been defined as the difference between two treatments—experimental and control—standardized by the pooled within-treatment standard deviation. [Osenberg et al., 1997, p. 798]

The effect size is just the standardised mean difference between the two groups. [Coe, 2002, p. 3]

In other articles, typically articles with a focus on meta-analytic applications, effect sizes are not necessarily standardized mean differences, but they are necessarily standardized or unitless:

A concept which could seem puzzling is that the effect size needs to be dimensionless, as it should deliver the same information regardless of the system used to take the observations. [Ialongo, 2016, p. 151]

There are myriad effect sizes from which the researcher can choose. Useful reviews of the choices have been provided by Kirk (1996), Snyder and Lawson (1993) and Friedman (1968), among others. Effect sizes can be categorized into two broad classes: variance accounted-for measures (e.g. R^2 , η^2) and standardized differences (e.g. Cohen's d , Hedges' g). Kirk [1996] identifies a third, 'miscellaneous' class. [Thompson, 1999, p. 171]

6 Unstandardized effect sizes are often dismissed as uninteresting

Several survey articles do not explicitly define the term “effect size” but exclusively discuss standardized effect sizes [Levine et al., 2008, Ferguson, 2009, Huberty, 2002]. Other articles admit that effect sizes can be unstandardized, but consider that such measures are not useful for meta-analysis and are barely worth mentioning. This is the case for the Sage Dictionary of Statistics:

Effect size: a term used in meta-analysis and more generally to indicate the relationship between two variables. The normal implication of the term effect size is that it indicates the size of the difference between the means of the conditions or groups on the dependent variable. Such an approach does not readily allow direct comparisons between studies using different measuring instruments and so forth. Consequently effect size is normally reported as a more standardized index such as Cohen's d . [Cramer and Howitt, 2004, p. 55]

Similarly, one of the most cited books on meta-analysis states:

This chapter began by presenting alternative measures of the size of the treatment effect: the raw score mean difference, the standard score mean difference (d or δ), and the point biserial correlation (r or ρ). Because different authors use different measures of the dependent variable, the raw score difference is not usually reasonable for meta-analysis. [...] Thus, the usual statistics used to characterize the size of the treatment effect are d and r . [Schmidt and Hunter, 2014, p. 328]

And here is another example from a survey on effect size metrics:

The ideal effect size estimate should be a statistic that has at least three characteristics. First, it should measure the practical significance of a result [...] Second, it should be independent of sample size [...]. Third, it should be metric free [...] although not all effect size measures are metric free, metric-free effect size measures facilitate the comparison of results across different studies. [...] In an effort to fulfill these three goals, quite a few statistics have been proposed as effect size measures. [Ives, 2003, p. 493]

These articles helped me understand why some researchers would object to the way the term “effect size” is currently used in the transparent statistics guidelines. This term is used differently in the field of meta-analysis, and I thought maybe the term has originated from that field and

has been corrupted by others. As it turns out, the participant to our CHI meeting who objected to our use of the term “effect size” gives lectures on meta-analysis.

At this point, I started to wonder whether the guidelines should adopt another term like “effect magnitude” so that it does not participate in spreading the confusion. But then, I looked at the original publications from Jacob Cohen, prompted by the following quote from Robert W. Frick, the researcher from whom Baguley borrowed the term “simple effect size”:

To most researchers, ‘effect size’ is the difference between the means of two conditions. I will call this simple effect size. Contrary to the implications of the APA Publication Manual (1994), some methodologists (e.g. Cohen, 1988) include simple effect size as a measure of effect size. To most methodologists, however, ‘effect size’ is a measure in which the simple effect size is compared to a measure of variance. I will call these measures statistical effect size. [Frick, 1999, p. 184]

7 What Jacob Cohen meant by “effect size”

Jacob Cohen is a renown statistician and psychologist who contributed to laying the foundations of meta-analysis. Although I couldn’t obtain articles older than 1970 and I can’t definitely confirm this, it is plausible that he was the one who coined the term “effect size”, or at least popularized it. Here is how he defined it in his famous book “Statistical Power Analysis for the Behavioral Sciences”, originally published in 1969:

The “effect size” [is] the degree to which the phenomenon exists. [Cohen, 1988, p. 4]

His definition was broad and clearly included unstandardized mean differences:

Without intending any necessary implication of causality, it is convenient to use the phrase “effect size” to mean “the degree to which the phenomenon is present in the population,” or “the degree to which the null hypothesis is false.” Whatever the manner of representation of a phenomenon in a particular research in the present treatment, the null hypothesis always means that the effect size is zero. [...] Thus, in terms of the previous illustrations: [...] If the population of consumers preferring brand A has a median annual income \$700 higher than that of brand B, the ES is \$700. If the population median difference and hence the ES is \$1000, the effect of income on brand preference would be larger. [...] Thus, whether measured in one unit or another, whether expressed as a difference between two population parameters or the departure of a population parameter from a constant or in any other suitable way, the ES can itself be treated as a parameter which takes the value zero when the null hypothesis is true and some other specific nonzero value when the null hypothesis is false, and in this way the ES serves as an index of degree of departure from the null hypothesis. [Cohen, 1988, pp. 9–10]

These excerpts are also present in the revised version of the first edition [Cohen, 1977] and probably also in the original 1969 edition, although I wasn’t able to access it. In 1970, Cohen already wrote:

The effect size measures the extent to which the null hypothesis is false, i.e., the discrepancy between S_0 and S_1 . It can be measured in raw units as simply $S_1 - S_0$, or in standard error units $\delta = \frac{S_1 - S_0}{\sigma}$. [Cohen, 1970, p. 814]

In his book on meta-analysis, Cohen then explains that it is desirable to come up with “universal effect size indices” that facilitate comparisons across studies:

In the previous illustrations, ES was variously expressed as a departure in percent from 50, a departure in IQ units from 100, a product moment r , a difference between two medians in dollars, etc. It is clearly desirable to reduce this diversity of units as far as possible, consistent with present usage by behavioural scientists. From one point of view, a universal ES index, applicable to all the various research issues and statistical models used in their appraisal, would be the ideal. [Cohen, 1988, p.11]

The index of effect size, ES [is] the departure of the “true” (population) state of affairs from H_0 , as assumed or hypothesized by the investigator, measured in metric-free units appropriate to the statistical test. [Cohen, 1973, p.226]

It appears that for Cohen, an index of effect size is a type of effect size, that is unitless and therefore useful for the purposes of meta-analysis. But he never abandoned his broad definition of effect size, as can be read in his 1990 retrospective article:

Effect-size measures include mean differences (raw or standardized), correlations and squared correlation of all kinds, odds ratios, kappas—whatever conveys the magnitude of the phenomenon of interest appropriate to the research context. [Cohen, 1990, p. 1310]

8 Why the different meanings today?

I’m far from having done a complete literature survey and I can only speculate. But it seems possible that:

1. When he formalized meta-analysis, Jacob Cohen introduced the term “effect size” that he unambiguously meant in a broad sense. He however focused his research on standardized measures of effect sizes, which he referred to as “indices of effect size”.
2. Subsequent researchers writing about meta-analysis came to use the term “effect size” as a synonym for those “interesting” (for meta-analysis) types of effect size, and perhaps ended up forgetting Cohen’s original definition.
3. More recently, methodologists focusing on statistical communication (rather than meta-analysis) started to employ the term “effect size” with a similar meaning to Cohen’s original definition, albeit stripped from its reference to statistical significance.

Two methodologists, Ken Kelley and Kristopher Preacher, provide a good overview of the different meanings of “effect size” used in the literature since Jacob Cohen [Preacher and Kelley, 2011, Kelley and Preacher, 2012]. They conclude that a broad definition of effect size is the most useful:

In response to the need for a general, inclusive definition of effect size, we define effect size as any measure that reflects a quantity of interest, either in an absolute sense or as compared with some specified value. The quantity of interest might refer to variability, association, difference, odds, rate, duration, discrepancy, proportionality, superiority, or degree of fit or misfit. [...] Although standardized effect sizes can be valuable, they are not always to be preferred over an effect size that is wedded to the original measurement scale, which may already be expressed in meaningful units that appropriately address the question of interest. [Preacher and Kelley, 2011, p. 95]

Their definition clearly includes unstandardized mean differences:

For example, group mean differences in scores on a widely understood instrument for measuring depressive symptoms are already expressed on a metric that is understandable to depression researchers, and to standardize effects involving the scale would only confuse matters. [Preacher and Kelley, 2011, p. 95]

We have seen several methodologists who argue for adopting a broad definition of the term “effect size”. Meanwhile, I haven’t seen papers where a methodologist explicitly argues that the term “effect size” should be used in a narrow sense. In all the papers I’ve seen, the narrow definition seems to be used in a [stipulative](#) manner, without authors giving much thought to it.

9 So how should we use the term?

From all this, it seems clear that **“effect size” can be used both in a broad sense and in a narrow sense**, provided that the author is consistent and clarifies what they are talking about.

In my opinion, **the broad sense is more useful in general methodological discussions**. The terms “unstandardized (or simple) effect size” and “standardized effect size” are widely used and perfectly capable of resolving any ambiguity. Using “effect size” in a broad sense doesn’t commit anyone to any particular statistical reporting practice. Anyone is free to point out limitations and possible dangers of reporting unstandardized effect sizes or standardized effect sizes. When the truth of a statement depends on whether it refers to standardized or unstandardized effect sizes, the term “effect size” without a qualifier should simply be avoided.

Using “effect size” to mean only “standardized effect size” is problematic because it requires finding another term for unstandardized effect sizes. I don’t recall any user of the narrow definition proposing such a term, and even if they did, a different term could give the wrong impression that unstandardized and standardized effect sizes have nothing in common. Cohen and others were clear about the analogies between the two.

I referred to “the” broad sense“ and ”the“ narrow sense, but in reality there are many possible definitions of ”effect size“ in both cases. If one needs an explicit definition, the one from Preacher and Kelley [2011] quoted above is quite extensive and informed by a good knowledge of the past literature. Kelley and Preacher [2012] later provided an updated definition that is more precise but perhaps a bit more awkward:

Effect size is defined as a quantitative reflection of the magnitude of some phenomenon that is used for the purpose of addressing a question of interest. [...] The question of

interest might refer to central tendency, variability, association, difference, odds, rate, duration, discrepancy, proportionality, superiority, or degree of fit or misfit, among others. [Kelley and Preacher, 2012, p. 140]

Statistical terminology is notoriously messy [Grace-Martin, 2018]. The case of effect sizes is one among many where blanket statements containing the word "statisticians" should be taken with skepticism [Gelman, 2018]. Few HCI researchers appreciate the number of statistics topics on which there is no consensus. When there is no consensus on a topic, article authors should be free to choose the option they prefer, provided that their choice is made explicit and justified with references from the methodology literature. It is bad behavior for a reviewer to ignore these references and force their own opinion on the authors. As HCI Professor Shumin Zhai put it, "research results, paper writing, and reviewing are just not the right forum for statistical method discussion." [Robertson and Kaptein, 2016, p. v].

Now there is also the question of what types of effect sizes we should report. This article hasn't addressed that question. I trust you have already guessed from the quotes that this is another question for which there is no simple and universal answer.

10 Acknowledgements

Thanks to [Fanny Chevalier](#), [Steve Haroz](#), [Yvonne Jansen](#), and [Matthew Kay](#) for their feedback.

References

- American Psychological Association. *The Publication manual of the APA (6th ed.)*. Washington, DC, 2010.
- T. Baguley. Standardized or simple effect size: What should be reported? *British journal of psychology*, 100(3):603–617, 2009.
- T. Baguley. *Serious stats: A guide to advanced statistics for the behavioral sciences*. Macmillan International Higher Education, 2012.
- R. Coe. It's the effect size, stupid: What effect size is and why it is important. 2002.
- J. Cohen. Approximate power and sample size determination for common one-sample and two-sample hypothesis tests. *Educational and Psychological Measurement*, 30(4):811–831, 1970.
- J. Cohen. Brief notes: statistical power analysis and research results. *American Educational Research Journal*, 10(3):225–229, 1973.
- J. Cohen. *Statistical power analysis for the behavioral sciences* (rev. ed.), 1977.
- J. Cohen. *Statistical power analysis for the behavioral sciences* (2nd ed.), 1988.
- J. Cohen. Things I have learned (so far). *American psychologist*, 45(12):1304, 1990.
- D. Cramer and D. L. Howitt. *The Sage dictionary of statistics: a practical resource for students in the social sciences*. Sage, 2004.

- G. Cumming. *Understanding the new statistics: Effect sizes, confidence intervals, and meta-analysis*. Routledge, 2013.
- P. Cummings. Arguments for and against standardized mean differences (effect sizes). *Archives of pediatrics & adolescent medicine*, 165(7):592–596, 2011.
- B. S. Everitt and A. Skrondal. *The Cambridge Dictionary of Statistics; 4th ed.* Cambridge University Press, Leiden, 2010.
- C. J. Ferguson. An effect size primer: A guide for clinicians and researchers. *Professional Psychology: Research and Practice*, 40(5):532, 2009.
- R. W. Frick. Defending the statistical status quo. *Theory & Psychology*, 9(2):183–189, 1999.
- A. Gelman. Pizzagate: The problem’s not with the multiple analyses, it’s with the selective reporting of results. *Statistical Modeling, Causal Inference, and Social Science*, 2018. URL <http://andrewgelman.com/2018/02/27/no-researchers-not-typically-set-prove-specific-hypothesis-study-begins/>.
- K. Grace-Martin. Series on confusing statistical terms. *The Analysis Factor*, 2018. URL <https://www.theanalysisfactor.com/series-on-confusing-statistical-terms/>.
- H. Hentschke and M. C. Stüttgen. Computation of measures of effect size for neuroscience data sets. *European Journal of Neuroscience*, 34(12):1887–1894, 2011.
- C. J. Huberty. A history of effect size indices. *Educational and Psychological measurement*, 62(2):227–240, 2002.
- C. Ialongo. Understanding the effect size and its measures. *Biochemia medica: Biochemia medica*, 26(2):150–163, 2016.
- B. Ives. Effect size use in studies of learning disabilities. *Journal of Learning Disabilities*, 36(6):490–504, 2003.
- V. B. Kampenes, T. Dybå, J. E. Hannay, and D. I. Sjøberg. A systematic review of effect size in software engineering experiments. *Information and Software Technology*, 49(11-12):1073–1086, 2007.
- K. Kelley and K. J. Preacher. On effect size. *Psychological methods*, 17(2):137, 2012.
- T. R. Levine, R. Weber, H. S. Park, and C. R. Hullett. A communication researchers’ guide to null hypothesis significance testing and alternatives. *Human Communication Research*, 34(2):188–209, 2008.
- A. McCook. We’re using a common statistical test all wrong: statisticians want to fix that. *Retraction Watch*, 2016. URL <https://retractionwatch.com/2016/03/07/were-using-a-common-statistical-test-all-wrong-statisticians-want-to-fix-that/>.
- C. W. Osenberg, O. Sarnelle, and S. D. Cooper. Effect size in ecological experiments: the application of biological models in meta-analysis. *The American Naturalist*, 150(6):798–812, 1997.
- K. J. Preacher and K. Kelley. Effect size measures for mediation models: quantitative strategies for communicating indirect effects. *Psychological methods*, 16(2):93, 2011.

- J. T. Richardson. Measures of effect size. *Behavior Research Methods, Instruments, & Computers*, 28(1):12–22, 1996.
- J. Robertson and M. Kaptein. *Modern statistical methods for HCI*. Springer, 2016.
- F. L. Schmidt and J. E. Hunter. *Methods of meta-analysis: Correcting error and bias in research findings. 2nd Edition*. Sage publications, 2014.
- B. Thompson. If statistical significance tests are broken/misused, what practices should supplement or replace them? *Theory & Psychology*, 9(2):165–181, 1999.
- Wikipedia contributors. Effect size. *Wikipedia*, 2018. URL https://en.wikipedia.org/wiki/Effect_size#Standardized_and_unstandardized_effect_sizes.
- L. Wilkinson. Statistical methods in psychology journals: Guidelines and explanations. *American psychologist*, 54(8):594, 1999.

Contents

1	Introduction	3
2	We should report effect sizes. But what are effect sizes?	3
3	Sources currently cited in the guidelines	3
4	“Effect size” is often used in a broad sense	4
5	“Effect size” is also often used in a narrow sense	5
6	Unstandardized effect sizes are often dismissed as uninteresting	6
7	What Jacob Cohen meant by “effect size”	7
8	Why the different meanings today?	8
9	So how should we use the term?	9
10	Acknowledgements	10



**RESEARCH CENTRE
SACLAY – ÎLE-DE-FRANCE**

1 rue Honoré d'Estienne d'Orves
Bâtiment Alan Turing
Campus de l'École Polytechnique
91120 Palaiseau

Publisher
Inria
Domaine de Voluceau - Rocquencourt
BP 105 - 78153 Le Chesnay Cedex
inria.fr

ISSN 0249-6399