

Using Effect Size—or Why the *P* Value Is Not Enough

GAIL M. SULLIVAN, MD, MPH
RICHARD FEINN, PhD

Statistical significance is the least interesting thing about the results. You should describe the results in terms of measures of magnitude—not just, does a treatment affect people, but how much does it affect them.

-Gene V. Glass¹

*The primary product of a research inquiry is one or more measures of effect size, not *P* values.*

-Jacob Cohen²

These statements about the importance of effect sizes were made by two of the most influential statistician-researchers of the past half-century. Yet many submissions to *Journal of Graduate Medical Education* omit mention of the effect size in quantitative studies while prominently displaying the *P* value. In this paper, we target readers with little or no statistical background in order to encourage you to improve your comprehension of the relevance of effect size for planning, analyzing, reporting, and understanding education research studies.

What Is Effect Size?

In medical education research studies that compare different educational interventions, effect size is the *magnitude of the difference between groups*. The absolute effect size is the difference between the average, or mean, outcomes in two different intervention groups. For example, if an educational intervention resulted in the improvement of subjects' examination scores by an average total of 15 of 50 questions as compared to that of another intervention, the absolute effect size is 15 questions or 3 grade levels (30%) better on the examination. Absolute effect size does not take into account the variability in scores, in that not every subject achieved the average outcome.

In another example, residents' self-assessed confidence in performing a procedure improved an average of 0.4 point on a Likert-type scale ranging from 1 to 5, after simulation training. While the absolute effect size in the first example

Gail M. Sullivan, MD, MPH, is Editor-in-Chief, *Journal of Graduate Medical Education*; Richard Feinn, PhD, is Assistant Professor, Department Psychiatry, University of Connecticut Health Center.

Corresponding author: Gail M. Sullivan, MD, MPH, University of Connecticut, 253 Farmington Avenue, Farmington, CT 06030-5215, gsullivan@ns01.uchc.edu

DOI: <http://dx.doi.org/10.4300/JGME-D-12-00156.1>

appears clear, the effect size in the second example is less apparent. Is a 0.4 change a lot or trivial? Accounting for variability in the measured improvement may aid in interpreting the magnitude of the change in the second example.

Thus, effect size can refer to the raw difference between group means, or absolute effect size, as well as standardized measures of effect, which are calculated to transform the effect to an easily understood scale. Absolute effect size is useful when the variables under study have intrinsic meaning (eg, number of hours of sleep). Calculated indices of effect size are useful when the measurements have no intrinsic meaning, such as numbers on a Likert scale; when studies have used different scales so no direct comparison is possible; or when effect size is examined in the context of variability in the population under study.

Calculated effect sizes can also quantitatively compare results from different studies and thus are commonly used in meta-analyses.

Why Report Effect Sizes?

The effect size is the main finding of a quantitative study. While a *P* value can inform the reader whether an effect exists, the *P* value will not reveal the size of the effect. In reporting and interpreting studies, both the substantive significance (effect size) and statistical significance (*P* value) are essential results to be reported.

For this reason, effect sizes should be reported in a paper's Abstract and Results sections. In fact, an estimate of the effect size is often needed before starting the research endeavor, in order to calculate the number of subjects likely to be required to avoid a Type II, or β , error, which is the probability of concluding there is no effect when one actually exists. In other words, you must determine what number of subjects in the study will be sufficient to ensure (to a particular degree of certainty) that the study has acceptable power to support the null hypothesis. That is, if no difference is found between the groups, then this is a true finding.

Why Isn't the *P* Value Enough?

Statistical significance is the probability that the observed difference between two groups is due to chance. If the *P* value is larger than the alpha level chosen (eg, .05), any observed difference is assumed to be explained by sampling variability. With a sufficiently large sample, a statistical test

TABLE 1

Index	Description ^b	Effect Size	Comments
Between groups			
Cohen's d^a	$d = M_1 - M_2 / s$ $M_1 - M_2$ is the difference between the group means (M); s is the standard deviation of either group	Small 0.2 Medium 0.5 Large 0.8 Very large 1.3	Can be used at planning stage to find the sample size required for sufficient power for your study
Odds ratio (OR)	<u>Group 1 odds of outcome</u> <u>Group 2 odds of outcome</u> If OR = 1, the odds of outcome are equally likely in both groups	Small 1.5 Medium 2 Large 3	For binary outcome variables Compares odds of outcome occurring from one intervention vs another
Relative risk or risk ratio (RR)	Ratio of probability of outcome in group 1 vs group 2; If RR = 1, the outcome is equally probable in both groups	Small 2 Medium 3 Large 4	Compares probabilities of outcome occurring from one intervention to another
Measures of association			
Pearson's r correlation	Range, -1 to 1	Small ± 0.2 Medium ± 0.5 Large ± 0.8	Measures the degree of linear relationship between two quantitative variables
r^2 coefficient of determination	Range, 0 to 1; Usually expressed as percent	Small 0.04 Medium 0.25 Large 0.64	Proportion of variance in one variable explained by the other

^a Adapted from Ferguson et al.⁹^b Based on Soper.⁷

will almost always demonstrate a significant difference, unless there is no effect whatsoever, that is, when the effect size is exactly zero; yet very small differences, even if significant, are often meaningless. Thus, reporting only the significant P value for an analysis is not adequate for readers to fully understand the results.

For example, if a sample size is 10,000, a significant P value is likely to be found even when the difference in outcomes between groups is negligible and may not justify an expensive or time-consuming intervention over another.

The level of significance by itself does not predict effect size. Unlike significance tests, effect size is independent of sample size. Statistical significance, on the other hand, depends upon both sample size and effect size. For this reason, P values are considered to be confounded because of their dependence on sample size. Sometimes a statistically significant result means only that a huge sample size was used.³

A commonly cited example of this problem is the Physicians Health Study of aspirin to prevent myocardial infarction (MI).⁴ In more than 22,000 subjects over an average of 5 years, aspirin was associated with a reduction in MI (although not in overall cardiovascular mortality) that was highly statistically significant: $P < .00001$. The study was terminated early due to the conclusive evidence,

means "heart attack"

and aspirin was recommended for general prevention. However, the effect size was very small: a risk difference of 0.77% with $r^2 = .001$ —an extremely small effect size. As a result of that study, many people were advised to take aspirin who would not experience benefit yet were also at risk for adverse effects. Further studies found even smaller effects, and the recommendation to use aspirin has since been modified.

Summary

Effect size helps readers understand the magnitude of differences found, whereas statistical significance examines whether the findings are likely to be due to chance. Both are essential for readers to understand the full impact of your work. Report both in the Abstract and Results sections.