Running head: EFFECT SIZE

1

Why Hedges'  $g_s^*$  based on the non-pooled standard deviation should be reported with Welch's t-test

Marie Delacre<sup>1</sup>, Daniel Lakens<sup>2</sup>, Christophe Ley<sup>3</sup>, Limin Liu<sup>3</sup>, & Christophe Leys<sup>1</sup>

- $^{1}$  Université Libre de Bruxelles, Service of Analysis of the Data (SAD), Bruxelles, Belgium
- <sup>2</sup> Eindhoven University of Technology, Human Technology Interaction Group, Eindhoven, the Netherlands
- $^3$  Universite it Gent, Department of Applied Mathematics, Computer Science and Statistics, Gent, Belgium

## Author Note

A Shiny app to compute point estimators and confidence intervals based on descriptive statistics is available from https://effectsize.shinyapps.io/deffsize/. We would like to thank Aaron Caldwell for his help creating the figures in the Shiny App. Daniël Lakens was funded by VIDI Grant 452-17-013 from the Netherlands Organisation for Scientific Research.

Correspondence concerning this article should be addressed to Marie Delacre, CP191, avenue F.D. Roosevelt 50, 1050 Bruxelles. E-mail: marie.delacre@ulb.ac.be

Why Hedges'  $g_s^*$  based on the non-pooled standard deviation should be reported with Welch's t-test

Effect sizes are an important outcome of empirical research. Moving beyond decisions about statistical significance, there is a strong call for researchers to report and interpret effect sizes and associated confidence intervals. This practice is highly endorsed by the American Psychological Association (APA) and the American Educational Research Association (American Educational Research Association, 2006; **APA\_2010?**).

In "between-subject" designs where individuals are randomly assigned into one of two independent groups and group scores are compared based on their means, the dominant estimator of effect size is Cohen's  $d_s$ , where the sample mean difference is divided by the pooled sample standard deviation (Peng, Chen, Chiang, & Chiang, 2013; Shieh, 2013). This estimator is available in many statistical software packages, such as SPSS and Stata. However, computing the pooled sample standard deviation assumes that both sample variances are estimates of a common population variance, which is known as the homogeneity of variance assumption. It has been widely argued that there are many fields in psychology where this assumption is ecologically unlikely (Delacre, Lakens, & Leys, 2017; Erceg-Hurn & Mirosevich, 2008; R. J. Grissom, 2000). The question how to deal with the assumption of equal variances has been widely explored in the context of hypothesis testing, and it is becoming increasingly common to by default report a t-test that does not assume equal variances, such as Welch's t-test.

However, the question which effect size to report when equal variances are not assumed has received less attention. One possible reason is that researchers have not found consensus on which of the available options should be used (Shieh, 2013). Even within the very specific context of an estimate for the standardized sample mean difference there is little agreement about which estimator is the best choice. In this article, we will review the main candidates that have been proposed in the literature in the d family of effect sizes,

without (Cohen's  $d_s$ , Glass's  $d_s$ , Shieh's  $d_s$  and Cohen's  $d_s^*$ ) and with correction for bias (Hedges'  $g_s$ , Glass's  $g_s$ , Shieh's  $g_s$  and Hedges'  $g_s^*$ ). We provide an R package and Shiny app to compute relevant effect size measures and their confidence intervals.

Before reviewing the most important effect size measures in the d-family, we will first list the different purposes effect size measures serve, and discuss the relationship between effect sizes, statistical, and practical significance. Based on a detailed description of the good properties an effect size measure should possess, we will evaluate these properties in the Monte Carlo simulations we performed to compare the different effect size estimators with correction for bias.

# Three purposes of effect size estimators

The effect size is a measure of the magnitude of an effect. In the context of the comparison of two groups based on their means, when the null hypothesis is the absence of effect, d-family effect size estimators estimate the magnitude of the differences between parameters of two populations groups are extracted from [e.g. the mean; Peng and Chen (2014)]. Such a measure can be used for three different purposes.

First, effect size measures can be used for *interpretative* purposes. They allow researchers to assess the practical significance of a result (i.e. statements about the relevance of an effect in real life). In order to assess the meaningfulness of an effect, we should be able to relate this effect size estimate with behaviors/meaningful consequences in the real world (Andersen\_et\_al\_2007?). This typically involves an analysis of the costs (determined by a specific context) and the benefits (in part determined by the size of the effect). It is important to remember an effect size is just a mathematical indicator of the magnitude of a difference, which depends on the way a variable is converted into numerical indicator. An effect size in itself is not a measure of the importance or the relevance of an effect for real life [even if benchmarks for small, medium, or large effect

sizes might have contributed to such a misinterpretation; Stout and Ruble (1995)].

Second, effect size measures can be used for *comparative* purposes. They allow researchers to assess the stability of results across designs, analyses, and sample sizes. This includes statistically comparing and combining the results from two or more studies in a meta-analysis.

Third, effect size measures can be used for inferential purposes. Hypothesis tests and confidence intervals based on the same statistical quantity are directly related: if the area of the null hypothesis is out of the  $(1-\alpha)$ -confidence interval, then the hypothesis test would also result in a p-value below the nominal alpha level. At the same time, the interval provides extra information about the precision of the sample estimate for inferential purposes (Altman, 2005; Ellis, 2015), and which effect sizes are excluded. The narrower the interval, the higher the precision, and the wider the confidence interval, the more the data lack precision. Effect size measures are also indirectly related to the hypothesis tests as effect sizes from previous studies can be used in an a-priori power analysis when planning a new study (Lakens, 2013; Prentice & Miller, 1990; Stout & Ruble, 1995; Sullivan & Feinn, 2012; Wilkinson & the Task Force on Statistical Inference, 1999).

## Properties of a good effect size estimator

The empirical value of an estimator (called the *estimate*) depends on the sample value. Different samples extracted from the same population will lead to different sample estimates of the population value. The *sampling distribution* of the estimator is the distribution of all estimates, based on all possible samples of size n extracted from one population. Studying the sampling distribution is useful, as it allows us to assess the qualities of an estimator. More specifically, three desirable properties a good estimator should possess for inferential purposes are: *unbiasedness*, *consistency* and *efficiency* (Wackerly, Mendenhall, & Scheaffer, 2008).

An estimator is unbiased if the distribution of estimates is centered around the true population parameter. On the other hand, an estimator is positively (or negatively) biased if the distribution is centered around a value that is higher (or lower) than the true population parameter (see Figure 1). In other words, examining the bias of an estimator tells us if estimates are on average accurate. The bias of a point estimator  $\hat{\delta}$  can be computed as

$$\delta_{bias} = E(\hat{\delta}) - \delta \tag{1}$$

where  $E(\hat{\delta})$  is the expectation of the sampling distribution of the estimator and  $\delta$  is the true (population) parameter.

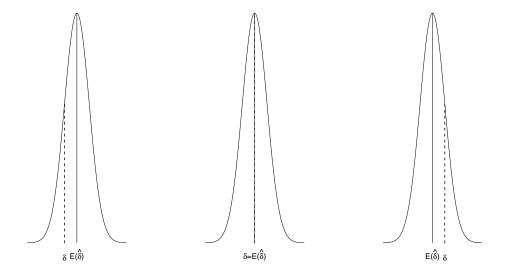


Figure 1. Sampling distribution for a positively biased (left), an unbiased (center) and a negatively biased estimator (right)

As we can see in Tables 1 and 2 the bias is directly related to the population effect size. The larger the population effect size, the larger the bias. It is therefore also interesting to examine the *relative bias*, defined as the ratio between the bias and the population effect size:

$$\delta_{relative\ bias} = \frac{E(\hat{\delta}) - \delta}{\delta} \tag{2}$$

While the bias informs us about the quality of estimates on average, in particular their capacity of lying close to the true value, it says nothing about individual estimates. Imagine a situation where the distribution of estimates is centered around the real parameter but with such a large variance that some point estimates are very far from the center. This would be problematic, since any single estimate might be very far from the true population value. Therefore it is not only essential for an estimator to be unbiased, but it is also desirable that the variability of its sampling distribution is small. Ideally, sample estimates are close to the true population parameter. Among two unbiased estimators  $\hat{\delta}_1$  and  $\hat{\delta}_2$ , we therefore say that  $\hat{\delta}_1$  is more efficient than  $\hat{\delta}_2$  if

$$Var(\hat{\delta}_1) \le Var(\hat{\delta}_2)$$
 (3)

where  $Var(\hat{\delta})$  is the variance of the sampling distribution of the estimator  $\hat{\delta}$ . Among all unbiased estimators, the more efficient estimator will be the one with the smallest variance <sup>1</sup>. The variance of an estimator  $\hat{\delta}$  is a function of its size (the larger the estimator, the larger the variance) and, therefore, we might be interested in evaluating the *relative* variance as the ratio between the variance and the square of the population estimator:

relative 
$$var(\hat{\delta}_1) = \frac{Var(\hat{\delta})}{\delta^2}$$
 (4)

Note that both unbiasedness and efficiency are very important when choosing an estimator. In some situations, it might be better to have a slightly biased estimator with low variance, (so that each estimate remains relatively close to the true parameter and one might be able to apply bias correction techniques) rather than an unbiased estimator with a large variance (Raviv, 2014).

Finally, the last property of a good point estimator is *consistency*. Consistency means that the bigger the sample size, the closer the estimate is to the population parameter. In other words, the estimates *converge* to the true population parameter.

<sup>&</sup>lt;sup>1</sup> The Cramér-Rao inequality provides a theoretical lower bound for the variance of unbiased estimators. An estimator reaching this bound is therefore optimally efficient.

### Different measures of effect sizes

The d-family effect sizes are commonly used for mean differences between groups or conditions. The population effect size is defined as

$$\delta = \frac{\mu_1 - \mu_2}{\sigma} \tag{5}$$

where both populations follow a normal distribution with mean  $\mu_j$  in the  $j^{th}$  population (j=1,2) and standard deviation  $\sigma$ . There exist different estimators of this effect size measure. For all, the mean difference is estimated by the difference  $\bar{X}_1 - \bar{X}_2$  of both sample means. When the equality of variances assumption is assumed,  $\sigma$  is estimated by pooling both sample standard deviations  $(S_1 \text{ and } S_2)$ . When the equality of variances assumption cannot be assumed, alternatives to the pooled standard deviation are available. In the next section we will present effect sizes that assume equal variances between groups (Cohen's  $d_s$  and Hedges'  $g_s$ ), and effect sizes that do not assume equal variances (Glass's  $d_s$ , Shieh's  $d_s$ , Hedges'  $d_s^*$ , Glass's  $g_s$ , Shieh's  $g_s$ , and Hedges'  $g_s^*$ ). For each effect size we will provide information about their theoretical bias, variance, and consistency.

## When variances are equal between groups

When we have good reasons to assume equality of variances between groups then the most common estimator of  $\delta$  is Cohen's  $d_s$ , where the sample mean difference is divided by a pooled error term (Cohen, 1965):

Cohen's 
$$d_s = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{(n_1 - 1) \times S_1^2 + (n_2 - 1) \times S_2^2}{n_1 + n_2 - 2}}}$$

where  $S_j$  is the standard deviation, and  $n_j$  the sample size of the  $j^{th}$  sample (j = 1, 2). The reasoning behind this measure is to make use of the fact that both samples share the same population variance (Keselman, Algina, Lix, Deering, & Wilcox, 2008), which means a more accurate estimation of the population variance can be achieved by pooling both

estimates of this parameter (i.e.  $S_1$  and  $S_2$ ). Since the larger the sample size, the more accurate the estimate, we give more weight to the estimate based on the larger sample size. Cohen's  $d_s$  is directly related to Student's t-statistic:

$$t_{Student} = \frac{Cohen's \ d_s}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}} \leftrightarrow Cohen's \ d_s = t_{Student} \times \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}$$
 (6)

Under the assumption of normality and equal variances between groups, Student's *t*-statistic follows a *t*-distribution with known degrees of freedom

$$df_{Student} = n_1 + n_2 - 2 \tag{7}$$

and noncentrality parameter <sup>2</sup>

$$ncp_{Student} = \frac{\delta_{Cohen}}{\sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

where  $\delta_{Cohen} = \frac{\mu_1 - \mu_2}{\sigma_{pooled}}$  and  $\sigma_{pooled} = \sqrt{\frac{(n_1 - 1) \times \sigma_1^2 + (n_2 - 1) \times \sigma_2^2}{n_1 + n_2 - 2}}$ . The relationship described in equation 6 and the theoretical distribution of Student's t-statistic allow us to determine the sampling distribution of Cohen's  $d_s$ , and therefore, its expectation and variance when the assumptions of normality and equal variances are met. All these equations are provided in Table 1. For interested readers, Supplemental Material 1 provides a detailed examination of the theoretical bias and variance of Cohen's  $d_s$  based on Table 1, as well as the bias and variance of all other estimators described later, based on Tables 2 and 3), with the goal to determine which parameters influence the bias and variance of different estimators. The main results will be discussed in the section "Monte Carlo Simulations: assessing the bias, efficiency and consistency of 5 estimators" below.

While Cohen's  $d_s$  is a consistent estimator, its bias and variance are substantial with small sample sizes, even under the assumptions of normality and equal variances (Lakens,

<sup>&</sup>lt;sup>2</sup> Under the null hypothesis of no differences between sample means, Student's t-statistic will follow a central t-distribution with  $n_1 + n_2 - 2$  degrees of freedom. However, when the null hypothesis is false, the distribution of this quantity will not be centered, and a noncentral t-distribution will arise.

2013). In order to compensate for Cohen's  $d_s$  bias with small sample sizes, Hedges and Olkin (1985) defined a bias-corrected version:

Hedges' 
$$g_s = Cohen's \ d_s \times \frac{\Gamma(\frac{df_{Student}}{2})}{\sqrt{\frac{df_{Student}}{2}} \times \Gamma(\frac{df_{Student}-1}{2})}$$

where  $df_{Student}$  has been defined in equation 7 and  $\Gamma()$  is the gamma function [for integers,  $\Gamma(x)$  is the factorial of x minus 1:  $\Gamma(x) = (x-1)!$ ; Goulet-Pelletier and Cousineau (2018)]. This equation can be approximated as follows:

$$Hedges' g_s = Cohen's d_s \times \left(1 - \frac{3}{4N - 9}\right)$$

where N is the total sample size. Hedges'  $g_s$  is theoretically unbiased when the assumptions of normality and equal variances are met (see Table 1). Moreover, while the variance of both Cohen's  $d_s$  and Hedges'  $g_s$  depend on the same parameters (i.e. the total sample size (N) and the sample sizes ratio  $\binom{n_1}{n_2}$ ), Hedges'  $g_s$  is less variable, especially with small sample sizes  $^3$ .

While the pooled error term is the best choice when variances are equal between groups (R. J. Grissom & Kim, 2001), it may not be well advised for use with data that violate this assumption (Cumming, 2013; R. J. Grissom & Kim, 2001; Grissom R. J. & Kim, 2005; Kelley, 2005; Shieh, 2013). When variances are unequal between groups, the expression in equation 5 is no longer valid because both groups do not share a common population variance. If we pool the estimates of two unequal population variances, the estimator of effect size will be smaller as it should be in case of positive pairing (i.e. the group with the larger sample size is extracted from the population with the larger variance) and larger as it should be in case of negative pairing (i.e. the group with the larger sample

 $<sup>^3</sup>$  In Table 1, one can see that the variance of Hedges'  $g_s$  equals the variance of Cohen's  $d_s$ , multiplied by  $\left[\frac{\Gamma(\frac{df}{2})}{\sqrt{\frac{df}{2}}\times\Gamma(\frac{df-1}{2})}\right]^2.$  This term is always less than 1 and tends to 1 when the sample sizes tends to infinity  $(52 \leq \left[\frac{\Gamma(\frac{df}{2})}{\sqrt{\frac{df}{2}}\times\Gamma(\frac{df-1}{2})}\right]^2 < 1 \text{ for } 3 \leq df < \infty).$  As a consequence, the larger the total sample size, the smaller the difference between the variance of Cohen's  $d_s$  and Hedges'  $g_s$ .

size is extracted from the population with the smaller variance). Because the assumption of equal variances across populations is rarely realistic in practice (Cain, Zhang, & Yuan, 2017; Delacre, Lakens, & Leys, 2017; Delacre, Leys, Mora, & Lakens, 2019; Erceg-Hurn & Mirosevich, 2008; Glass, Peckham, & Sanders, 1972; R. J. Grissom, 2000; Micceri, 1989; Yuan, Bentler, & Chan, 2004), both Cohen's  $d_s$  and Hedges'  $g_s$  should be abandoned in favor of an alternative robust to unequal population variances.

Table 1

Expentency, bias and variance of Cohen's  $d_s$  and Hedges'  $g_s$  under the assumptions that independent residuals are normally distributed with equal variances across groups.

Variance	$\frac{N \times df}{n_1 n_2 \times (df - 2)} + \delta_{Cohen}^2 \left[ \frac{df}{df - 2} - c_f^2 \right]$	$\approx \frac{N \times df}{n_1 n_2 \times (df - 2)} + \delta_C^2 ohen \left[ \frac{df}{df - 2} - \left( \frac{1}{1 - \frac{3}{4N - 9}} \right)^2 \right]$	$Var(Cohen's d_s) \times \left[ \frac{\Gamma(\frac{df}{2})}{\sqrt{\frac{df}{2} \times \Gamma(\frac{df-1}{2})}} \right]^2$	$\approx Var(Cohen's\ d_s) \times \left[1 - \frac{3}{4N-9}\right]^2$
Expectation	$\delta_{Cohen}  imes c_f$	$pprox rac{\delta_{Cohen}}{\left(1-rac{3}{4N-9} ight)}$	$\delta_{Cohen}$	
дþ	N-2		N-2	
	Cohen's $d_s$		${\rm Hedges'}\ g_s$	

Note.  $\delta_{Cohen} = \frac{\mu_1 - \mu_2}{\sigma_{pooled}}$  and  $c_f = \frac{\sqrt{\frac{df}{2}} \times \Gamma(\frac{df-1}{2})}{\Gamma(\frac{df}{2})}$ ; Cohen's  $d_s$  is a biased estimator, because its expectation differs from the population effect size. Moreover, the larger the population estimator  $(\delta_{Cohen})$ , the larger the bias. Indeed, the bias is the difference between the expectation and  $\delta_{Cohen}$ :  $\delta_{bias} = \delta_{Cohen} \times (c_f - 1)$ . On the other hand, Hedges'  $g_s$  is an unbiased estimator, because its expectation equals  $\delta_{Cohen}$ ; equations in Table 1 require  $df \geq 3$  (i.e.  $N \geq 5$ ).

# When variances are unequal between populations

In his review, Shieh (2013) mentions three options available in the literature to deal with the case of unequal variances: (A) the Glass's  $d_s$ , (B) the Shieh's  $d_s$  and (C) the Cohen's  $d_s^*$ .

Glass's  $d_s$ . When comparing one control group with one experimental group, Glass, McGav, and Smith (2005) recommend using the standard deviation of the control group as standardizer. This yields

$$Glass's d_s = \frac{\bar{X}_e - \bar{X}_c}{S_c}$$

where  $\bar{X}_e$  and  $\bar{X}_c$  are the sample means of the experimental and control groups, and  $S_c$  is the sample SD of the control group. One argument in favour of using  $S_c$  as standardizer is the fact that it is not affected by the experimental treatment. When it is easy to identify which group is the "control" one, it is therefore convenient to compare the effect size estimation of different designs studying the same effect (Cumming, 2013). However, defining this group is not always obvious (Coe, 2002). This could induce large ambiguity because depending on the chosen SD as standardizer, measures could be substantially different (Shieh, 2013). The distribution of Glass's  $d_s$  is defined as in Algina, Keselman, and Penfield (2006):

Glass's 
$$d_s \sim \sqrt{\frac{1}{n_c} + \frac{\sigma_e^2}{n_e \times \sigma_c^2}} \times t_{df,ncp}$$
 (8)

where  $n_c$  and  $n_e$  are the sample sizes of the control and experimental groups, and df and ncp are defined as follows:

$$df = n_c - 1 \tag{9}$$

$$ncp = \frac{\delta_{Glass}}{\sqrt{\frac{1}{n_c} + \frac{\sigma_e^2}{n_e \times \sigma_c^2}}}$$

where  $\delta_{Glass} = \frac{\mu_c - \mu_e}{\sigma_c}$  and  $\mu_c$  and  $\mu_e$  are respectively the mean of the populations control and experimental groups are extracted from. Thanks to equation 8, we can compute its

theoretical expectation and variance when the assumption of normality is met (see Table 2), and therefore determine which factors influence bias and variance, and how they do so (see Supplemental Material 1).

Table 2

Expentency, bias and variance of Glass's d<sub>s</sub> and Cohen's d<sub>s</sub>\* and Shieh's d<sub>s</sub> under the assumption that independent residuals are normally distributed.

Variance	$\frac{df}{df-2} \times \left(\frac{1}{n_c} + \frac{\sigma_c^2}{n_e \sigma_c^2}\right) + \delta_{Glass}^2 \left(\frac{df}{df-2} - c_f^2\right)$	$\frac{df}{df - 2} \times \frac{2\left(\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2^2}\right)}{\sigma_1^2 + \sigma_2^2} + \left(\delta_{Cohen}^*\right)^2 \left(\frac{df}{df - 2} - c_f^2\right)$	$\approx \frac{df}{df - 2} \times \frac{2\left(\frac{\sigma_1^2}{n_1 + \frac{\sigma_2^2}{n_2}}\right)}{\sigma_1^2 + \sigma_2^2} + \left(\delta_{Cohen}^*\right)^2 \left[\frac{df}{df - 2} - \left(\frac{4 \ df - 1}{4(df - 1)}\right)^2\right]$	$\frac{df}{(df-2)N} + \delta_{Shieh}^2 \left( \frac{df}{df-2} - c_f^2 \right)$
Expectation	$\delta_{Glass}  imes c_f$	$\delta^*_{Cohen}  imes c_f$	$\approx \delta_{Cohen}^*  imes rac{4df-1}{4(df-1)}$	$\delta_{Shieh}  imes c_f$
df	$n_c - 1$	$\frac{(n_1-1)(n_2-1)(\sigma_1^2+\sigma_2^2)^2}{(n_2-1)\sigma_1^4+(n_1-1)\sigma_2^4}$		$\frac{\left(\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)^2}{\frac{(\sigma_1^2/n_1)^2}{n_1 - 1} + \frac{(\sigma_2^2/n_2)^2}{n_2 - 1}}$
	Glass's $d_s$	Cohen's $d_s^*$		Shieh's $d_s$

Note.  $c_f = \frac{\sqrt{\frac{2}{2}} \times 1(\frac{-2}{4})}{\Gamma(\frac{d}{2})}$ ;  $\delta_{Glass} = \frac{\mu_c - \mu_e}{\sigma_c}$ ,  $\delta_{Shieh} = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2}{\sigma_1^2 N_1} + \frac{\sigma_2^2}{\sigma_2^2 N_1}}}$  and  $\delta_{Cohen}^* = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2 + \sigma_2^2}{\sigma_2^2}}}$ ; all estimators are biased estimators, because their expectations differ from the population effect size  $\delta$ . Moreover, the larger the population estimator  $(\delta)$ , the Note.  $c_f = \frac{\sqrt{\frac{df}{2} \times \Gamma(\frac{df-1}{2})}}{\Gamma(\frac{df}{2})}; \ \delta_{Glass} = \frac{\mu_c - \mu_e}{\sigma_c}, \ \delta_{Shieh} = \frac{1}{\sqrt{\sigma_c}}$ 

larger the bias. Indeed, the bias is the difference between the expectation and  $\delta$ :  $\delta_{bias} = \delta \times (c_f - 1)$ .

equations require  $df \geq 3$  and at least 2 subjects per group.

Shieh's  $d_s$ . Kulinskaya and Staudte (2007) were the first to recommend the use of a standardizer that takes the sample sizes allocation ratios into account, in addition to the variance of both samples. Shieh (2013), following Kulinskaya and Staudte (2007), proposed a modification of the exact SD of the sample mean difference:

Shieh's 
$$d_s = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{S_1^2/q_1 + S_2^2/q_2}}; \quad q_j = \frac{n_j}{N} (j = 1, 2)$$

where  $N = n_1 + n_2$ . Shieh's  $d_s$  is directly related with Welch's t-statistic:

$$Shieh's \ d_s = \frac{t_{Welch}}{\sqrt{N}} \leftrightarrow t_{welch} = Shieh's \ d_s \times \sqrt{N}$$
 (10)

The exact distribution of Welch's t-statistic is more complicated than the exact distribution of Student's t-statistic, but it can be approximated, under the assumption of normality, by a t-distribution with degrees of freedom and noncentrality parameters (Shieh, 2013; Welch, 1938):

$$df_{Welch} = \frac{\left(\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)^2}{\frac{(\sigma_1^2/n_1)^2}{n_1 - 1} + \frac{(\sigma_2^2/n_2)^2}{n_2 - 1}}$$

$$ncp_{Welch} = \delta_{Shieh} \times \sqrt{N} = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$
(11)

where  $\delta_{Shieh} = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2}{n_1/N} + \frac{\sigma_2^2}{n_2/N}}}$ . The relationship described in equation 10 and the theoretical distribution of Welch's t-statistic allow us to approximate the sampling distribution of Shieh's  $d_s$ . Based on the sampling distribution of Shieh's  $d_s$ , we can estimate its theoretical expectation and variance under the assumption of normality (see Table 2), and thereby determine which factors influence bias and variance, and how they do so (see Supplemental Material 1).

As demonstrated in the Appendix, when variances and sample sizes are equal across groups, the biases and variances of Cohen's  $d_s$  and Shieh's  $d_s$  are identical except for multiplication by a constant. The same is true for the estimators  $\delta_{Cohen}$  and  $\delta_{Shieh}$ :

$$\delta_{Cohen} = 2 \times \delta_{Shieh}$$
 (considering  $\sigma_1 = \sigma_2$  and  $n_1 = n_2$ ) (12)

$$Bias_{Cohen's d_s} = 2 \times Bias_{Shieh's d_s}$$
 (considering  $\sigma_1 = \sigma_2$  and  $n_1 = n_2$ ) (13)

$$Var_{Cohen's d_s} = 4 \times Var_{Shieh's d_s}$$
 (considering  $\sigma_1 = \sigma_2$  and  $n_1 = n_2$ ) (14)

We can deduce from equations 12, 13 and 14 that relative to their respective population effect size, Cohen's  $d_s$  and Shieh's  $d_s$  are equally accurate. In other words, their relative bias and variance are identical. This is a good illustration of our motivation to favor relative bias and variance (previously defined in equations 2 and 4) over the most commonly used raw bias and variance (previously defined in equations 1 and 3).

When sample sizes are not equal, according to the statistical properties of Welch's statistic under heteroscedasticity, Shieh's  $d_s$  accounts for the allocation ratio of sample sizes to each condition. The lack of generality caused by taking this specificity of the design into account has led Cumming (2013) to question its usefulness in terms of interpretability: when the mean difference  $(\bar{X}_1 - \bar{X}_2)$ ,  $S_1$ , and  $S_2$  remain constant, Shieh's  $d_s$  will vary as a function of the sample sizes allocation ratio (unlike Cohen's  $d_s^*$  that we will define below). At the population level,  $\delta_{Shieh}$  also depends on the sample sizes allocation ratio, as illustrated in the following shiny application: https://effectsize.shinyapps.io/ShiehvsCohen/.

Cohen's  $d_s^*$ . An effect size estimator based on the sample mean difference divided by the square root of the non pooled average of both variance estimates was suggested by Welch (1938). Here, we indicate the difference between Cohen's  $d_s$  (based on the pooled standard deviations) and Cohen's  $d_s^*$  with an asterisk. This yields:

Cohen's 
$$d_s^* = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{\left(S_1^2 + S_2^2\right)}{2}}}$$

where  $\bar{X}_j$  is the mean and  $S_j$  is the standard deviation of the  $j^{th}$  sample (j = 1,2). We know the distribution of Cohen's  $d_s^*$  (Huynh, 1989):

Cohen's 
$$d_s^* \sim \sqrt{\frac{2(n_2 \times \sigma_1^2 + n_1 \times \sigma_2^2)}{n_1 n_2(\sigma_1^2 + \sigma_2^2)}} \times t_{df^*,ncp^*}$$
 (15)

Where  $df^*$  and  $ncp^*$  are defined as follows:

$$df^* = \frac{(n_1 - 1)(n_2 - 1)(\sigma_1^2 + \sigma_2^2)^2}{(n_2 - 1)\sigma_1^4 + (n_1 - 1)\sigma_2^4}$$

$$ncp^* = \delta_{Cohen}^* \times \sqrt{\frac{n_1 n_2(\sigma_1^2 + \sigma_2^2)}{2(n_2 \sigma_1^2 + n_1 \sigma_2^2)}} = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$
(16)

where  $\delta_{Cohen}^* = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2 + \sigma_2^2}{2}}}$ . Using equation 15 we can compute its theoretical expectation and variance when the assumption of normality is met (see Table 2), and therefore determine which factors influence bias and variance, and how they do so (see Supplemental Material 1). This estimator has been widely criticized, because it results in a variance term of an artificial population (i.e. since the variance term does not estimate the variance of one or the other group, the composite variance is an estimation of the variance of an artificial population) and is therefore very difficult to interpret (R. J. Grissom & Kim, 2001), and unless both sample sizes are equal, the variance term does not correspond to the variance of the mean difference (Shieh, 2013).

However, we will show throughout the simulation section that this estimator exhibits very good inferential properties. Moreover, it has a constant value across sample sizes ratios, as shown in the Shiny App at https://effectsize.shinyapps.io/ShiehvsCohen/.

Glass's  $\mathbf{g_s}$ , Shieh's  $\mathbf{g_s}$  and Hedges'  $\mathbf{g_s^*}$ . As for Cohen's  $d_s$ , a Hedges' correction can be applied in order to compensate for the bias of Glass's  $d_s$ , Shieh's  $d_s$  and Cohen's  $d_s^*$  with small sample sizes (see Table 2). This correction has the following general form:

$$g_s = d_s \times \frac{\Gamma(\frac{\nu}{2})}{\sqrt{\frac{\nu}{2}} \times \Gamma(\frac{\nu-1}{2})}$$

where the distinct values of  $\nu$  are provided in equation 9 for Glass's  $g_s$ , in equation 16 for Hedges'  $g_s^*$  and in equation 11 for Shieh's  $g_s$ . The three corrected estimators are theoretically unbiased when the assumption of normality is met. Their variance is a function of the same parameters as their biased equivalent. However, due to the correction they have a smaller variance, especially with small sample sizes, as shown in Table 3. In summary:

• The variances of Hedges'  $g_s^*$  and Shieh's  $g_s$  depend on the total sample size (N), their respective population effect size  $(\delta)$ , and the interaction between the sample sizes ratio and the SD-ratio  $\left(\frac{n_1}{n_2} \times \frac{\sigma_1}{\sigma_2}\right)$ .

• The variance of Glass's  $g_s$  also depends on N,  $\delta$  and  $\frac{n_c}{n_e} \times \frac{\sigma_c}{\sigma_e}$ . In addition, there is also a main effect of the SD-ratio  $\left(\frac{\sigma_c}{\sigma_e}\right)$  on its variance.

How these parameters influence the variance of the estimators will be summarized and illustrated in Monte Carlo simulations below.

Table 3

Expectation, bias and variance of Glass's d<sub>s</sub> and Cohen's d<sub>s</sub>\* and Shieh's d<sub>s</sub> under the assumption that independent residuals are normally distributed.

	2	$\left(\frac{1}{2}\right)^2$	2
Variance	$Var(Glass's\ d_s) \times \left(\frac{\Gamma(\frac{df}{2})}{\sqrt{\frac{df}{2}} \times \Gamma(\frac{df-1}{2})}\right)$	$Var(Cohen's d_s^*) \times \left(\frac{\Gamma(\frac{df}{2})}{\sqrt{\frac{df}{2}} \times \Gamma(\frac{df-1}{2})}\right)$	$Var(Shieh's\ d_s) \times \left( \frac{\Gamma(\frac{df}{2})}{\sqrt{\frac{df}{2}} \times \Gamma(\frac{df-1}{2})} \right)$
Expectation	$\delta_{glass}$	$\delta^*_{Cohen}$	$\delta_{Shieh}$
Jp	$n_c - 1$	$\frac{(n_1-1)(n_2-1)(\sigma_1^2+\sigma_2^2)^2}{(n_2-1)\sigma_1^4+(n_1-1)\sigma_2^4}$	$\approx \frac{\left(\frac{\sigma_1^2 + \frac{\sigma_2^2}{n_1}}{n_1 + \frac{\sigma_2^2}{n_2 - 1}}\right)^2}{\frac{(\sigma_1^2/n_1)^2}{n_1 - 1} + \frac{(\sigma_2^2/n_2)^2}{n_2 - 1}}$
	Glass's $g_s$	${\rm Cohen's}\ g_s^*$	${\rm Shieh's}\ g_s$

Note.  $c_f = \frac{\sqrt{\frac{df}{2}} \times \Gamma(\frac{df-1}{2})}{\Gamma(\frac{df}{2})}; \delta_{Glass} = \frac{\mu_c - \mu_e}{\sigma_c}, \delta_{Shieh} = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_f^2}{n_1/N} + \frac{\sigma_s^2}{n_2/N}}} \text{ and } \delta_{Cohen}^* = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_f^2 + \sigma_s^2}{n_2}}}; \text{ all estimators are unbiased}$  estimators, because their expectations equal the population effect size  $\delta$ ; equations require  $df \geq 3$  and at least 2 subjects per

group.

### Monte Carlo Simulations

Assessing the bias, efficiency and consistency of 5 estimators.

#### Method.

We performed Monte Carlo simulations using R (version 3.5.0) to assess the bias, efficiency and consistency of Cohen's  $g_s$ , Glass's  $g_s$  (using respectively the sample SD of the first or second group as a standardizer), Hedges'  $g_s^*$  and Shieh's  $g_s$ .

A set of 100,000 datasets was generated for 1,008 scenarios as a function of different criteria. In 252 scenarios, samples were extracted from a normally distributed population (in order to ensure the reliability of our calculation method) and in 756 scenarios, samples were extracted from non normal population distributions. In order to assess the quality of estimators under realistic deviations from the normality assumption, we referred to the review of Cain, Zhang, and Yuan (2017). Cain et al. (2017) investigated 1,567 univariate distributions from 194 studies published by authors in Psychological Science (from January 2013 to June 2014) and the American Education Research Journal (from January 2010 to June 2014). For each distribution, they computed Fisher's skewness

$$G_1 = \frac{\sqrt{n(n-1)}}{n-2} \frac{m_3}{\sqrt{(m_2)^3}}$$

and kurtosis

$$G_2 = \frac{n-1}{(n-2)(n-3)} \times \left[ (n+1) \left( \frac{m_4}{(m_2)^2} - 3 \right) + 6 \right]$$

where n is the sample size and  $m_2$ ,  $m_3$  and  $m_4$  are respectively the second, third and fourth centered moments. They found values of kurtosis from G2 = -2.20 to 1,093.48. According to their suggestions, throughout our simulations, we kept constant the population kurtosis value at the 99th percentile of their distribution of kurtosis, i.e. G2=95.75. Regarding skewness, we simulated population parameter values which correspond to the 1st and 99th percentile of their distribution of skewness, i.e. respectively G1 = -2.08 and G1 = 6.32. We also simulated samples extracted from population where G1 = 0, in order to assess the

main effect of high kurtosis on the quality of estimators. All possible combinations of population skewness and kurtosis and the number of scenarios for each combination are summarized in Table 4.

Table 4

Number of combinations of skewness and kurtosis in our simulations.

		Kurtosis		
		0	95.75	TOTAL
	0	252	252	504
Skewness	-2.08	/	252	252
	6.32	/	252	252
	TOTAL	252	756	1008

Note. Fisher's skewness (G1) and kurtosis (G2) are presented in Table 4. The 252 combinations where both G1 and G2 equal 0 correspond to the normal case.

For the 4 resulting combinations of skewness and kurtosis (see Table 4), all other parameter values were chosen in order to illustrate the consequences of factors identified as playing a key role on the variance of unbiased estimators. We manipulated the population mean difference  $(\mu_1 - \mu_2)$ , the sample sizes (n), the sample size ratio  $(n\text{-ratio} = \frac{n_1}{n_2})$ , the population SD-ratio (i.e.  $\frac{\sigma_1}{\sigma_2}$ ), and the sample size and population variance pairing  $(\frac{n_1}{n_2} \times \frac{\sigma_1}{\sigma_2})$ . In our scenarios,  $\mu_2$  was always 0 and  $\mu_1$  varied from 1 to 4, in steps of 1 (so does  $\mu_1 - \mu_2$ )<sup>4</sup>. Moreover,  $\sigma_1$  always equals 1, and  $\sigma_2$  equals .1, .25, .5, 1, 2, 4 or 10, and

 $<sup>^4</sup>$  In the original plan, we had added 252 simulations in which  $\mu_1$  and  $\mu_2$  were both null. We decided not to

therefore, so were the ratios of  $\frac{\sigma_1}{\sigma_2}$ . The simulations for which both  $\sigma_1$  and  $\sigma_2$  equal 1 are the particular case of homoscedasticity, or equal population variances across groups. The sample sizes of both groups  $(n_1 \text{ and } n_2)$  were 20, 50 or 100. When sample sizes of both groups are equal, the n-ratio equals 1 (this is known as a balanced design). All possible combinations of n-ratio and population SD-ratio were simulated in order to distinguish scenarios where both sample sizes and population variances are unequal across groups (with positive pairing when the group with the largest sample size is extracted from the population with the largest SD, and negative pairing when the group with the smallest sample size is extracted from the population with the smallest SD) and scenarios with no pairing between sample sizes and variances (sample sizes and/or population SD are equal across all groups). In sum, the simulations grouped over different sample sizes yield 4 conditions (a, b, c and d) based on the *n*-ratio, population SD-ratio, and sample size and population variance pairing, as summarized in Table 5. We chose to divide scenarios into these 4 conditions because analyses in Supplemental Material 1 revealed main and interaction effects of sample sizes ratio and SD-ratio on the bias and variance of some estimators.

present the results of these simulations in the main article, because the relative bias and the relative variance appeared to us to be very useful to fully understand the comparison of the estimators, and computing them is impossible when the real mean difference is zero. Indeed, for these specific configurations, both relative bias and relative variance would have infinite values due to the presence of the population effect size term in their denominator. However, these extra simulations were included in the simulation checks, in Supplemental Material 2.

Table 5
4 conditions based on the n-ratio and the SD-ratio.

			n-ratio	
		1	>1	<1
	1	a	b	b
SD-ratio	>1	С	d	d
	<1	С	d	d

### Results.

Before presenting the comparison of the estimators for each condition, it is useful to make some general comments.

1) We previously discussed the fact that raw bias and variances are sometimes misleading. They can give the illusion of huge differences between two estimators, even if these differences only reflect a change of unit (i.e. different population effect sizes). To better understand this, imagine a sample of 15 people for whom we know the height (in meters) and we compute a sample variance of 0.06838. If we convert sizes to centimeters and compute the sample variance again, we find a measure of 683.8 (i.e. 10<sup>4</sup> larger). Both measures represent the same amount of variability, but they are expressed in different units. The same issue due to a change in scales occurs when comparing the estimates of different population measures. To avoid this possible confusion, we will only present the relative bias and relative variance in all Figures (and anytime we will mention the biases and variances in the results section, we will be referring to relative bias and variance). For interested readers, illustrations

of the raw bias and variance are available on Github: https://anonymous.4open.science/r/deffectsize-B1D2/.

- 2) For the sake of readability, the vertical axis differs across plots.
- 3) Throughout this section, we will *compare* the relative bias and variance of different estimators, but we do not present bias and variance in absolute terms. We chose very extreme (although realistic) conditions, and we know that none of the parametric measures of effect size will be robust against such extreme conditions. Our goal is therefore to study the robustness of the estimators against normality violations only in comparison with the robustness of other indicators, but not in absolute terms.

After these general remarks, we will analyze each condition separately. In all Figures presented below, for different sub-conditions, the averaged relative bias and relative variance of five estimators are presented. When describing the Glass's  $g_s$  estimators, we will systematically refer to the "control group" as the condition the standardizer is based on (i.e. the first group when using  $S_1$  as standardizer, the second group when using  $S_2$  as standardizer). The other condition will be referred to as the "experimental group."

When variances are equal across groups.

Figures 2 and 3 represent configurations where the equality of variances assumption is met. According to our expectations, one observes that the bias of all estimators is approximately zero as long as the normality assumption is met (first column in both Figures)<sup>5</sup>. However, the more the data generation process deviations from the normality assumption (i.e. when moving from left to right in the Figures), the larger the bias in the estimators.

<sup>&</sup>lt;sup>5</sup> When looking at relative bias for all estimators, the maximum departure from zero is 0.0064 when sample sizes are equal across groups, and 0.0065 with unequal sample sizes.

We will observe that Glass's  $g_s$  should always be avoided when the equality of variance assumption is met. Hedges'  $g_s$ , Hedges'  $g_s^*$  and Shieh's  $g_s$  perform equally well as long as the sample size ratio is close to 1 (condition a; see Figure 2). However, when designs are highly unbalanced (condition b; see Figure 3), Shieh's  $g_s$  is not consistent anymore, while Hedges'  $g_s^*$  remain consistent, Hedges's  $g_s$  is a better estimator. For interested readers, these findings are detailed in the three paragraphs below.

Figure 2 illustrates scenarios where both population variances and sample sizes are equal across groups (condition a). One can first notice that all estimators are consistent, as their bias and variance decrease when the total sample size increases. For any departure from the normality assumption, both bias and variance of Hedges'  $g_s$ , Shieh's  $g_s$  and Hedges'  $g_s^*$  are similar<sup>6</sup> and smaller than the bias and variance of Glass's  $g_s$  estimates using either  $S_1$  or  $S_2$  as a standardizer. Moreover, when samples are extracted from skewed distributions, Glass's  $g_s$  will show different bias and variance as a function of the chosen standardizer  $(S_1 \text{ or } S_2)$ , even if both  $S_1$  and  $S_2$  are estimates of the same population variance, based on the same sample size. This is due to non-null correlations of opposite sign between the mean difference  $(\bar{X}_1 - \bar{X}_2)$  and respectively  $S_1$  and  $S_2$ . In Supplemental Material 3, we detailed in which situation a non-null correlation occurs between the sample mean difference  $(\bar{X}_1 - \bar{X}_2)$  and the standardizer of compared estimators as well as the way this correlation impacts the bias and variance of estimators.

Figure 3 illustrates scenarios where population variances are equal across groups, but

<sup>&</sup>lt;sup>6</sup> While the bias and variance of Cohen's  $d_s$ , Cohen's  $d_s^*$  and Shieh's  $d_s$  are identical, the bias and variance of Hedges'  $g_s^*$  and Shieh's  $g_s$  (these last two having identical bias and variance). Indeed, because of the sampling error, differences remain between sample variances, even when population variances are equal between groups. Since the Hedges' correction applied to Cohen's  $d_s$  does not contain the sample variances (unlike the correction applied on both other estimators), the bias and variance of Hedges'  $g_s$  are slighly different from the bias and variance of Hedges'  $g_s^*$  and Shieh's  $g_s$ .

sample sizes are unequal (condition b). For any departures from the normality assumptions, Hedges'  $g_s$  shows the smallest bias and variance. Hedges'  $g_s$  and Hedges'  $g_s^*$  are consistent estimators (i.e. the larger the sample sizes, the lower the bias and the variance), unlike Shieh's  $g_s$  and Glass's  $g_s$ . The bias of Glass's  $g_s$  does not depend either on the size of the experimental group or on the total sample size. The only way to decrease the bias of Glass's  $g_s$  is therefore to add subjects in the control group. On the other hand, the variance of Glass's  $g_s$  depends on both sample sizes, but not in an equivalent way: in order to reduce the variance, it is much more efficient to add subjects in the control group and when the size of the experimental group decreases so does the variance, even when the total sample size is increased. Regarding Shieh's  $g_s$ , for a given sample size ratio, the bias and variance will decrease when sample sizes increase. However, there is a large effect of the sample sizes ratio such that when the sample sizes ratio moves away from 1 by adding subjects, bias and variance might increase.<sup>7</sup> On the other hand, when the sample sizes ratio moves closer to 1 by adding subjects, the bias will decrease.

When samples are extracted from skewed distributions and have unequal sizes (the two last columns in Figure 3), for a constant total sample size, Glass's  $g_s$ , Shieh's  $g_s$  and Hedges'  $g_s^*$  will show different bias and variance depending on which group is the largest one (e.g. when distributions are right-skewed, the bias and variance of all these estimators when  $n_1$  and  $n_2$  are respectively 50 and 20 are not the same as their bias and variance when  $n_1$  and  $n_2$  are respectively 20 and 50). This is due to a non-null correlation of opposite sign between the mean difference ( $\bar{X}_1 - \bar{X}_2$ ) and their respective standardizers depending on

<sup>&</sup>lt;sup>7</sup> Regarding variance, in Supplemental Material 1, we mentioned that when the population effect size is zero, the larger the total sample size, the lower the variance, whether the sample sizes ratio is constant or not. We also mentioned that this is no longer true when the population effect size is not zero. In our simulations the effect size is never zero. The effect size effect is partially visible in Figure 3 because we do not entirely remove the effect size effect when we divide the variance by  $\delta^2$ . This is due to the fact that one term, in the equation of the variance computation, does not depend on the effect size.

which group is the largest one, as detailed in Supplemental Material 3. One observes that under these configurations, the bias and variance of Glass's  $g_s$  are sometimes a bit smaller and sometimes much larger than the bias and variance of Shieh's  $g_s$  and Cohen's  $d_s^*$ .

When variances are unequal across groups.

Figures 4 to 9 represent configurations where the equality of variances assumption is not met. According to our expectations, one observes that the bias of all estimators is approximately zero as long as the normality assumption is met (first column in all Figures), and the further from the normality assumption (i.e. when moving from left to right in Figures), the larger the bias<sup>9</sup>. It might be considered surprising that the bias of Hedges'  $g_s$  remains very small throughout these conditions. As discussed in the section "Different measures of effect size," Hedges'  $g_s$  should be avoided when population variances and sample sizes are unequal across groups, because of the pooled error term. When pooling the estimates of two unequal population variances, the resulting estimator will be smaller (in case of positive pairing) or larger (in case of negative pairing) than it should be. At the same time, when pooling two unequal population variances, the population effect size will also be smaller (in case of positive pairing) or larger (in case of negative pairing) as it should be. As a consequence, the distortion cannot be seen through the difference between the expected estimator and the population effect size measure. For this reason, the bias

<sup>&</sup>lt;sup>8</sup> Supplemental Material 3 shows that when the  $\mu_1 - \mu_2 > 0$  (like in our simulations), all other parameters being equal, an estimator is always less biased and variable when choosing a standardizer that is positively correlated with  $\bar{X}_1 - \bar{X}_2$ . Supplemental Material 3 also shows that the smaller  $n_c$ , the larger the magnitude of correlation between  $S_c$  and  $\bar{X}_1 - \bar{X}_2$ . When  $cor(S_c, \bar{X}_1 - \bar{X}_2)$  is positive, the positive effect of increasing the magnitude of the correlation is counterbalanced by the negative effect of reducing  $n_c$ . On the other hand, when  $cor(S_c, \bar{X}_1 - \bar{X}_2)$  is negative, the negative effect of increasing the magnitude of the correlation is amplified by the negative effect of decreasing  $n_c$ . This explains why the difference between Glass's  $g_s$  and other estimators is larger when Glass's  $g_s$  is the least efficient estimator.

<sup>&</sup>lt;sup>9</sup> When looking at the relative bias for all estimators, the maximum departure from zero is 0.0173 when sample sizes are equal across groups, and 0.0274 when both sample sizes and variances differ across groups.

and variance of Hedges'  $g_s$  will not be taken into account in the following comparisons.

We will observe that when variances are unequal across populations, Glass's  $g_s$  sometimes performs better, but also sometimes performs much worst than Shieh's  $g_s$  and Hedges'  $g_s^*$ , both in terms of bias and variance. The performance of Glass's  $g_s$  highly depends on parameters that we cannot control (i.e. a triple interaction between the n-ratio, the SD-ratio and the correlation between the standardizer and the mean difference) and for this reason, we do not recommend using it. When the sample sizes ratio is close to 1, Shieh's  $g_s$  and Hedges'  $g_s^*$  are both appropriate but the further the sample sizes ratio is from 1, the larger the bias of Shieh's  $g_s$  in order that in the end, the measure that we believe performs best across scenarios is Hedges'  $g_s^*$ .

Figures 4 and 5 are dedicated to scenarios where population variances are unequal between groups and sample sizes are equal (condition c). In Figure 4, scenarios are subdivided as a function of the sample sizes and one can notice that all estimators are consistent, as their bias and variance decrease when the total sample size increases. In Figure 5, scenarios are subdivided as a function of the SD-ratio. Because the comparison pattern remains very similar for all sample sizes, we present only scenarios when sample sizes equal 20. One should first notice that for all estimators in Figure 5, the relative variance seems to be much larger when  $S_2 > S_1$ . <sup>10</sup> This information should not be taken into account because it is only an artefact of our simulation conditions combined with the way we computed the relative variance. <sup>11</sup>

When samples are extracted from skewed distributions, the bias and variance of

<sup>&</sup>lt;sup>10</sup> The difference between the variance of estimators when the second group is 10 times larger than the first group was so large that we decided to not present it, for the sake of readability of the Figures.

<sup>&</sup>lt;sup>11</sup> We previously mentioned that when dividing the variance by  $\delta^2$ , we do not entirely remove the effect size effect. Actually, we introduce  $\delta^2$  in the denominator of the first term, in the equation of the variance computation. Because we performed our simulations in order that  $\sigma_1$  always equals 1, the smaller  $S_2$ , the larger the population effect size and therefore, the smaller the relative variance.

Glass's  $g_s$  are sometimes smaller and sometimes larger than the bias of Shieh's  $g_s$  and Hedges'  $g_s$ . This is mainly due to the fact that when two samples of same sizes are extracted from two skewed distributions with unequal variances (the two last columns in Figure 5), there will be non-null correlations of opposite sign between the mean difference  $(\bar{X}_1 - \bar{X}_2)$  and the standardizer of all estimators, depending on which population variance is larger  $^{12}$ .

Figures 6 to 9 are dedicated to scenarios where both sample sizes and population variances differ across groups. Due to a high number of combinations between the sample sizes-ratio and the SD-ratio in our simulations, we decided to present only some conditions. Because equations in Table 3 revealed an interaction effect between the sample sizes ratio and the SD-ratio on the bias and variance of Hedges'  $g_s^*$  and Shieh's  $g_s$  (see Supplemental Material 1), we chose to present all configurations where the larger SD is 10 times larger than the smaller SD (Figures 6 and 7), and configurations where the larger SD is twice larger than the smaller SD (Figures 8 and 9), in order to compare the effect of the sample sizes ratio on the bias and variance of all estimators when the SD-ratio is large  $(\frac{\sigma_1}{\sigma_2} = 10 \ or \ .1$ ) or medium  $(\frac{\sigma_1}{\sigma_2} = 2 \ or \ .5)$ .

When distributions are symmetric, the bias of Glass's  $g_s$  only depends on the size of the control group and is therefore not impacted by either the sample sizes ratio or the total sample size. When comparing Figures 6 to 9, one can also notice that the bias of Glass's  $g_s$  does not depend on the SD-ratio either. Unlike the bias of Glass's  $g_s$ , its variance depends on both sample sizes, but not in an equivalent way. In most scenarios it is more efficient, in order to reduce the variance of Glass's  $g_s$ , to add subjects in the control group. Regarding

When population variances are unequal, a non-null correlation occurs between standardizer estimates and  $\bar{X}_1 - \bar{X}_2$ . For standardizers computed based on both  $S_1$  and  $S_2$ , the sign of the correlation between the standardizer and the mean difference will be the same as the sign of the correlation between the mean difference and the estimate of the larger population variance. For interested readers, this is detailed in Supplemental Material 3.

Hedges'  $g_s^*$  and Shieh's  $g_s$ , their respective biases and variances depend on an interaction effect between the sample sizes ratio and the SD-ratio  $\left(\frac{n_1}{n_2} \times \frac{\sigma_1}{\sigma_2}\right)$ : the sample sizes ratio associated with the smallest bias and variance is not the same when the more variable group is 10 times more variable than the other group (Figures 6 and 7) than when it is only twice more variable (Figures 8 and 9). However, the respective biases and variances of Hedges'  $g_s^*$  and Shieh's  $g_s$  are always smaller when there is a positive pairing between sample sizes and variances. When samples are extracted from skewed distributions, the bias and variance of Glass's  $g_s$  are sometimes smaller and sometimes larger than the bias of Shieh's  $g_s$  and Hedges'  $g_s^*$ , due to a combination of three factors: (1) which group is larger, (2) which group has the smallest standard deviation and (3) what is the correlation between the standardizer and the mean difference.

### Recommendations

We recommend using Hedges'  $g_s^*$  in order to assess the magnitude of the effect when comparing two independent means, because a) it does not rely on the equality of population variances assumption (unlike Hedges'  $g_s$ ), b) it is always consistent (unlike Shieh's  $g_s$ ), c) it is easy to interpret (Hedges'  $g_s^*$  can be interpreted in the same way as Hedges'  $g_s$ ) and d) it remains constant for any sample sizes ratio, even when population variances are unequal across groups, as shown in the Shiny App at https://effectsize.shinyapps.io/ShiehvsCohen/.

Effect sizes estimates such as Hedges'  $g_s^*$  should always be reported with a confidence interval. To help researchers compute Hedges'  $g_s^*$  and it's confidence interval we created the R package deffectsize (see https://anonymous.4open.science/r/deffectsize-B1D2/). The datacohen\_CI function was built in order to compute point estimators and confidence intervals based on raw data and the cohen\_CI function was built in order to compute point estimators and confidence intervals based on descriptive statistics (sample means, sample variances and sample sizes). By default, unbiased Cohen's  $g_s^*$  is computed but it is also possible to compute biased estimators (e.g. Cohen's  $g_s^*$ ) and/or to use a pooled error

term as standardizer by assuming that the equality of population variances is met (e.g. Hedges'  $g_s$  or Cohen's  $d_s$ , depending on whether we choose to compute unbiased or biased estimator). Other functions ( $datashieh\_CI$ ,  $shieh\_CI$ ,  $dataglass\_CI$  and  $glass\_CI$ ) are available in order to compute Shieh's  $g_s$  (or Shieh's  $d_s$ ) and Glass's  $g_s$  (or Glass's  $d_s$ ) as well as their respective confidence intervals, even though we don't recommend to use these effect sizes by default. Researchers who do not use R can use a Shiny app to compute point estimators and confidence intervals based on descriptive statistics: https://effectsize.shinyapps.io/deffsize/.

### References

- Algina, J., Keselman, H. J., & Penfield, R. D. (2006). Confidence intervals for an effect size when variances are not equal. *Journal of Modern Applied Statistical Methods*, 5(1), 1–13. https://doi.org/10.22237/jmasm/1146456060
- Altman, G. D. (2005). Why we need confidence intervals. World Journal of Surgery, 29, 554–556. https://doi.org/10.1007/s00268-005-7911-0
- American Educational Research Association. (2006). Standards for reporting on empirical social science research in AERA publications. *Educational Researcher*, 35, 33–40. https://doi.org/10.3102/0013189X035006033
- Cain, M. K., Zhang, Z., & Yuan, K.-H. (2017). Univariate and multivariate skewness and kurtosis for measuring nonnormality: Prevalence, influence and estimation. *Behavior Research Methods*, 49(5), 1716–1735. https://doi.org/10.3758/s13428-016-0814-1
- Coe, R. (2002). It's the effect size, stupid. What effect size is and why it is important. Retrieved from https://www.leeds.ac.uk/educol/documents/00002182.htm
- Cohen, J. (1965). Some statistical issues in psychological research. In *Handbook of Clinical Psychology* (B. B. Wolman, pp. 95–121). New York: McGraw-Hill.

Cumming, G. (2013). Cohen's d needs to be readily interpretable: Comment on Shieh (2013). Behavior Research Methods, 45, 968–971. https://doi.org/10.3758/s13428-013-0392-4

- Delacre, M., Lakens, D., & Leys, C. (2017). Why psychologists should by default use Welch's t-test instead of Student's t-test. *International Review of Social Psychology*, 30(1), 92–101. https://doi.org/10.5334/irsp.82
- Delacre, M., Leys, C., Mora, Y. L., & Lakens, D. (2019). Taking parametric assumptions seriously: Arguments for the use of Welch's F-test instead of the classical F-test in one-way ANOVA. *International Review of Social Psychology*, 32(1), 1–12. https://doi.org/http://doi.org/10.5334/irsp.198
- Ellis, P. D. (2015). The Essential Guide to Effect Sizes: Statistical Power, Meta-Analysis, and the Interpretation of Research Results (Cambridge University Press). Cambridge, UK.
- Erceg-Hurn, D. M., & Mirosevich, V. M. (2008). Modern robust statistical methods: An easy way to maximize the accuracy and power of your research. *American Psychologist*, 63(7), 591–601. https://doi.org/10.1037/0003-066X.63.7.591
- Glass, G. V., McGav, B., & Smith, M. L. (2005). *Meta-analysis in Social Research* (Sage). Beverly Hills, CA.
- 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. https://doi.org/10.3102/00346543042003237
- Goulet-Pelletier, J.-C., & Cousineau, D. (2018). A review of effect sizes and their confidence intervals, Part I: The Cohen's d family. *The Quantitative Methods for Psychology*, 14(4), 242–265. https://doi.org/10.20982/tqmp.14.4.p242

Grissom, R. J. (2000). Heterogeneity of variance in clinical data. *Journal of Consulting* and Clinical Psychology, 68(1), 155–165. https://doi.org/10.1037//0022-006x.68.1.155

- Grissom, R. J., & Kim, J. J. (2001). Review of assumptions and problems in the appropriate conceptualization of effect size. *Psychological Methods*, 6(2), 135–146. https://doi.org/10.1037/1082-989X.6.2.135
- Grissom, R. J., & Kim, J. J. (2005). Effect size for research: A broad practical approach. (Lawrence Erlbaum Associates, Mahwah, N.J.). London.
- Hedges, L. V., & Olkin, I. (1985). Statistical Methods for Meta-analysis (Academic Press). Cambridge, Massachusetts. https://doi.org/10.1016/C2009-0-03396-0
- Huynh, C.-L. (1989). A unified approach to the estimation of effect size in meta-analysis.
  San Francisco: Paper presented at the Annual Meeting of the American Educational Research Association.
- Kelley, K. (2005). The effects of nonnormal distributions on confidence intervals around the standardized mean difference: Bootstrap and parametric confidence intervals.

  Educational and Psychological Measurement, 65(1), 51–69.

  https://doi.org/10.1177/0013164404264850
- Keselman, H. J., Algina, J., Lix, L. M., Deering, K. N., & Wilcox, R. R. (2008). A generally robust approach for testing hypotheses and setting confidence intervals for effect sizes.
  Psychological Methods, 13(2), 110–129. https://doi.org/10.1037/1082-989X.13.2.110
- Kulinskaya, E., & Staudte, R. G. (2007). Confidence intervals for the standardized effect arising in the comparison of two normal populations. *Statistics in Medicine*, 26, 2853–2871. https://doi.org/10.1002/sim.2751
- Lakens, D. (2013). Calculating and reporting effect sizes to facilitate cumulative science: A practical primer for t-tests and ANOVAs. Frontiers in Psychology, 4 (863), 1–12. https://doi.org/10.3389/fpsyg.2013.00863

Micceri, T. (1989). The unicorn, the normal curve, and other improbable creatures.

Psychological Bulletin, 105(1), 156–166. https://doi.org/10.1037/0033-2909.105.1.156

- Peng, C.-Y. J., & Chen, L.-T. (2014). Beyond Cohen's d: Alternative effect size measures for between-subject designs. *The Journal of Experimental Education*, 82(1), 22–50. https://doi.org/10.1080/00220973.2012.745471
- Peng, C.-Y. J., 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. https://doi.org/10.1080/00220973.2012.745471
- Prentice, D. A., & Miller, D. T. (1990). When small effects are impressive. *Psychological Bulletin*, 112(1), 160–164.
- Raviv, E. (2014). Bias vs. consistency. Retrieved March 25, 2020, from https://eranraviv.com/bias-vs-consistency/
- Shieh, G. (2013). Confidence intervals and sample size calculations for the standardized mean difference effect size between two normal populations under heteroscedasticity.

  Behavior Research Methods, 45, 955–967. https://doi.org/10.3758/s13428-012-0228-7
- Stout, D. D., & Ruble, T. L. (1995). Assessing the practical signficance of empirical results in accounting education research: The use of effect size information. *Journal of Accounting Education*, 13(3), 281–298.
- Sullivan, G. M., & Feinn, R. (2012). Using effect size—or why the p value is not enough.
  Journal of Graduate Medical Education, 279–282.
  https://doi.org/10.4300/JGME-D-12-00156.1
- Wackerly, D. D., Mendenhall, W., & Scheaffer, R. L. (2008). *Mathematical Statistics with Applications (7th edition)* (Brooks/Cole, Cengage Learning). Belmont, USA.
- Welch, B. L. (1938). The significance of the difference between two means when the population variances are unequal. *Biometrika*, 29, 350–362.

Wilkinson, L., & the Task Force on Statistical Inference. (1999). Statistical methods in psychology journals: Guidelines and explanations. *American Psychologist*, 54(8), 594–604.

Yuan, K.-H., Bentler, P. M., & Chan, W. (2004). Structural equation modeling with heavy tailed distributions. *Psychometrika*, 69(3), 421–436. https://doi.org/10.1007/bf02295644 groups (condition a)

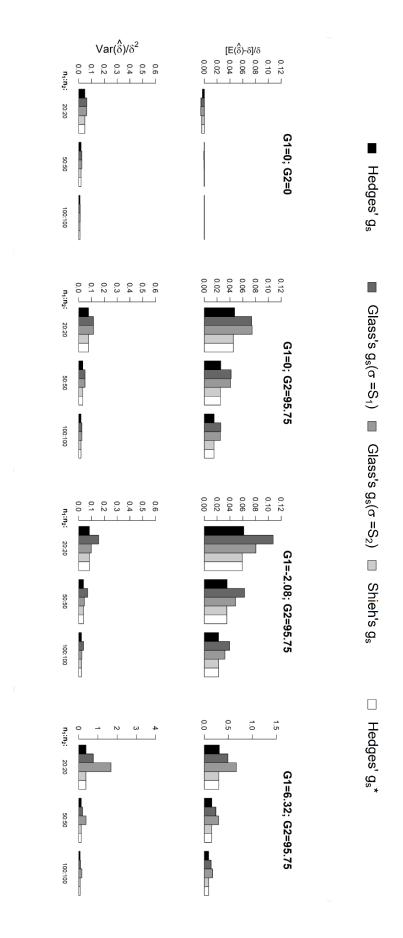


Figure 2. Bias and efficiency of estimators of standardized mean difference, when variances and sample sizes are equal across

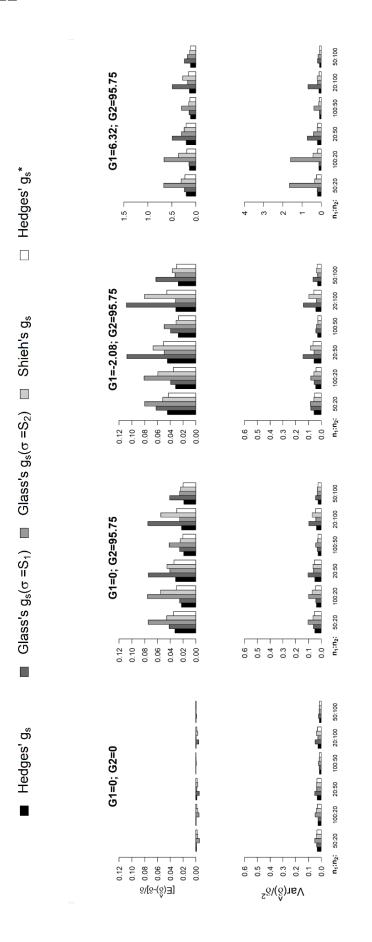
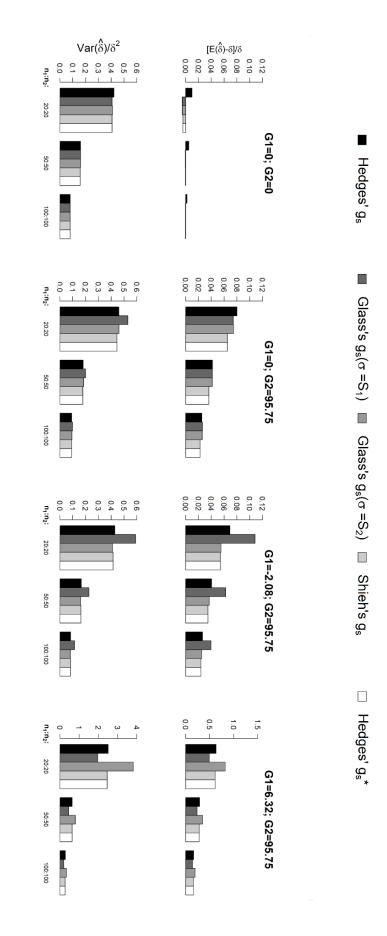


Figure 3. Bias and efficiency of estimators of standardized mean difference, when variances are equal across groups and sample sizes are unequal (condition b)



sample sizes are equal (condition c), as a function of n-ratio Figure 4. Bias and efficiency of estimators of standardized mean difference, when variances are unequal across groups and

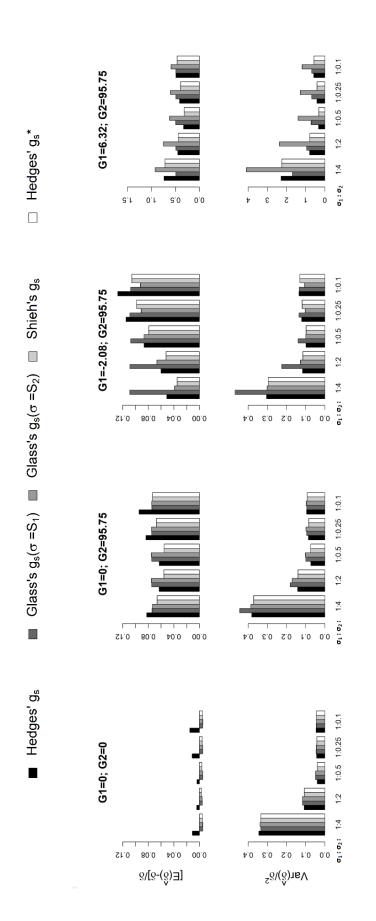


Figure 5. Bias and efficiency of estimators of standardized mean difference, when variances are unequal across groups and sample sizes are equal (condition c) as a function of the SD-ratio (when  $n_1 = n_2 = 20$ )

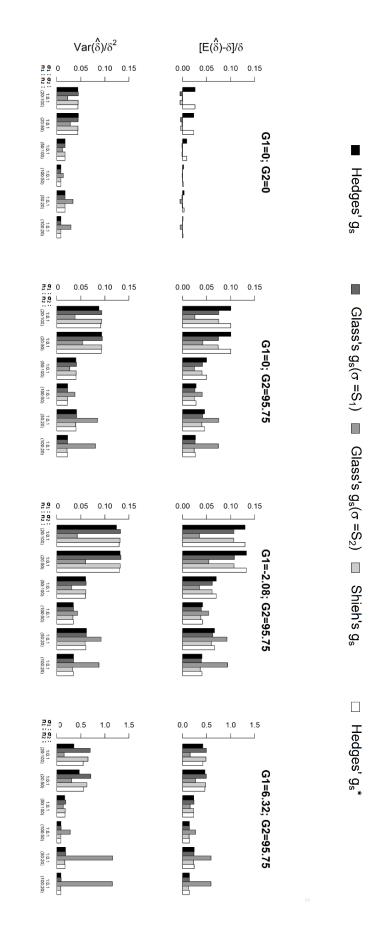


Figure 6. Bias and efficiency of estimators of standardized mean difference, when variances and sample sizes are unequal across

groups (condition d),  $\frac{\sigma_1}{\sigma_2} = 10$  and  $\sigma_1 > \sigma_2$ 

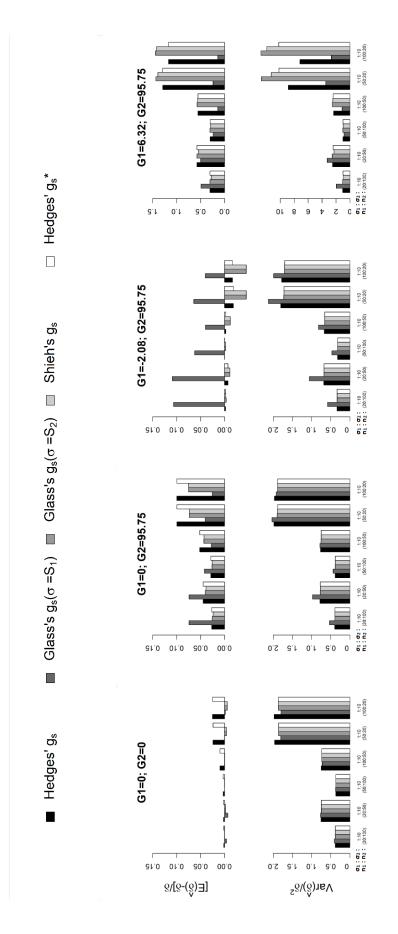


Figure 7. Bias and efficiency of estimators of standardized mean difference, when variances and sample sizes are unequal across groups (condition d),  $\frac{\sigma_1}{\sigma_2} = 10$  and  $\sigma_1 < \sigma_2$ 

groups (condition d),  $\frac{\sigma_1}{\sigma_2} = 2$  and  $\sigma_1 > \sigma_2$ 

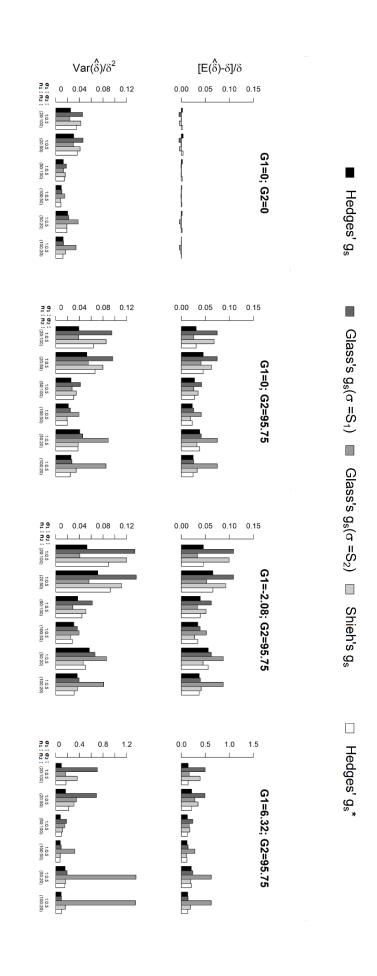


Figure 8. Bias and efficiency of estimators of standardized mean difference, when variances and sample sizes are unequal across

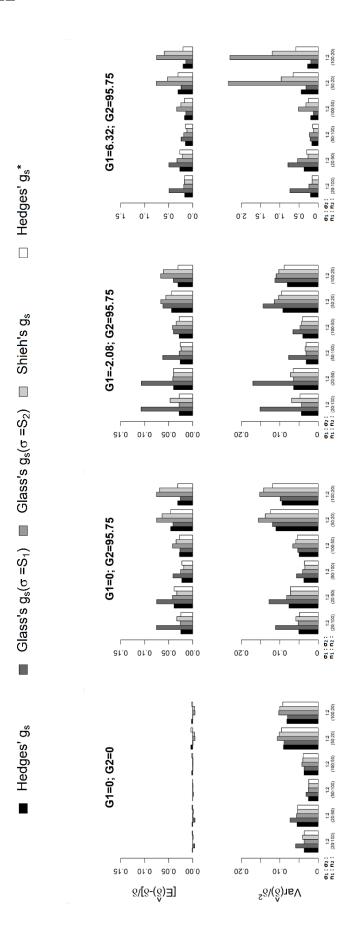


Figure 9. Bias and efficiency of estimators of standardized mean difference, when variances and sample sizes are unequal across groups (condition d),  $\frac{\sigma_1}{\sigma_2}=2$  and  $\sigma_1<\sigma_2$ 

# Appendix

The bias of Cohen's  $d_s$  is twice as large as the bias of Shieh's  $d_s$  when population variances and sample sizes are equal across groups: mathematical demonstration.

As mentioned in Table 1, the bias of Cohen's  $d_s$  is defined as

$$Bias_{Cohen's d_s} = \delta_{Cohen} \times \left( \frac{\sqrt{\frac{df_{Student}}{2}} \times \Gamma\left(\frac{df_{Student}-1}{2}\right)}{\Gamma\left(\frac{df_{Student}}{2}\right)} - 1 \right)$$
(17)

with

$$\delta_{Cohen} = \frac{\mu_1 - \mu_2}{\sqrt{\frac{(n_1 - 1) \times \sigma_1^2 + (n_2 - 1) \times \sigma_2^2}{n_1 + n_2 - 2}}}$$

and

$$df_{Student} = n_1 + n_2 - 2$$

As mentioned in Table 2, the bias of Shieh's  $d_s$  is defined as

$$Bias_{Shieh's d_s} = \delta_{Shieh} \times \left( \frac{\sqrt{\frac{df_{Welch}}{2}} \times \Gamma\left(\frac{df_{Welch}-1}{2}\right)}{\Gamma\left(\frac{df_{Welch}}{2}\right)} - 1 \right)$$
(18)

with

$$\delta_{Shieh} = \frac{\mu_1 - \mu_2}{\sqrt{\frac{\sigma_1^2}{n_1/N} + \frac{\sigma_2^2}{n_2/N}}} \quad (N = n_1 + n_2)$$

and

$$df_{Welch} = \frac{\left(\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}\right)^2}{\frac{(\sigma_1^2/n_1)^2}{n_1 - 1} + \frac{(\sigma_2^2/n_2)^2}{n_2 - 1}}$$

When  $n_1 = n_2 = n$  and  $\sigma_1 = \sigma_2 = \sigma$ ,  $\delta_{Cohen}$  is twice larger than  $\delta_{Shieh}$ , as shown below in equations 19 and 20:

$$\delta_{Cohen} = \frac{\mu_1 - \mu_2}{\sqrt{\frac{2(n-1)\sigma^2}{2(n-1)}}} = \frac{\mu_1 - \mu_2}{\sigma}$$
 (19)

$$\delta_{Shieh} = \frac{\mu_1 - \mu_2}{\sqrt{2\left(\frac{\sigma^2}{n/(2n)}\right)}} = \frac{\mu_1 - \mu_2}{2\sigma}$$
(20)

Moreover, degrees of freedom associated with Student's t-test and Welch's t-test are identical, as shown below in equations 21 and 22:

$$df_{Student} = 2(n-1) \tag{21}$$

$$df_{Welch} = \frac{[2(\sigma^2/n)]^2}{\frac{2(\sigma^2/n)^2}{n-1}} = \mathbf{2(n-1)}$$
 (22)

Equations 17 and 18 can therefore be redefined as follows:

$$Bias_{Cohen's d_s} = \frac{\mu_1 - \mu_2}{\sigma} \times \left( \frac{\sqrt{n-1} \times \Gamma\left(\frac{2n-3}{2}\right)}{\Gamma(n-1)} - 1 \right)$$
 (23)

$$Bias_{Shieh's d_s} = \frac{\mu_1 - \mu_2}{2\sigma} \times \left(\frac{\sqrt{n-1} \times \Gamma\left(\frac{2n-3}{2}\right)}{\Gamma(n-1)} - 1\right)$$
(24)

We can therefore conclude that the bias of Cohen's  $d_s$  is twice larger than the bias of Shieh's  $d_s$ .

The variance of Cohen's  $d_s$  is four times larger than the bias of Shieh's  $d_s$  when population variances and sample sizes are equal across groups: mathematical demonstration.

The variance of Cohen's  $d_s$  is defined in Table 1 as

$$Var_{Cohen's\ d_s} = \frac{N \times df_{Student}}{n_1 n_2 \times (df_{Student} - 2)} + \delta_{Cohen}^2 \left[ \frac{df_{Student}}{df_{Student} - 2} - \left( \frac{\sqrt{\frac{df_{Student}}{2}} \times \Gamma\left(\frac{df_{Student} - 1}{2}\right)}{\Gamma\left(\frac{df_{Student}}{2}\right)} \right)^2 \right]$$
(25)

and the variance of Shieh's  $d_s$  is defined in Table 2 as

$$Var_{Shieh's\ ds} = \frac{df_{Welch}}{(df_{Welch} - 2)N} + \delta_{Shieh}^2 \left[ \frac{df_{Welch}}{df_{Welch} - 2} - \left( \frac{\sqrt{\frac{df_{Welch}}{2}} \times \Gamma\left(\frac{df_{Welch} - 1}{2}\right)}{\Gamma\left(\frac{df_{Welch}}{2}\right)} \right)^2 \right]$$
(26)

We have previously shown in equations 21 and 22 that degrees of freedom associated with Student's t-test and Welch's t-test equal 2(n-1), when  $n_1 = n_2 = n$  and  $\sigma_1 = \sigma_2 = \sigma$ . As a consequence, the first term of the addition in equation 25 is 4 times larger than the first term of the addition in equation 26:

$$\frac{N \times df_{Student}}{n_1 n_2 \times (df_{Student} - 2)} = \frac{2n \times 2(n-1)}{n^2 \times (2n-4)} = \frac{\mathbf{4(n-1)}}{\mathbf{n(2n-4)}}$$
$$\frac{df_{Welch}}{(df_{Welch} - 2)N} = \frac{2(n-1)}{2n(2n-4)} = \frac{\mathbf{n-1}}{\mathbf{n(2n-4)}}$$

We have also previously shown in equations 19 and 20 that  $\delta_{Cohen}$  is twice larger than  $\delta_{Shieh}$  when  $n_1 = n_2 = n$  and  $\sigma_1 = \sigma_2 = \sigma$  and, therefore,  $\delta_{Cohen}^2$  is four times larger than  $\delta_{Shieh}^2$ . As a consequence, the second term of the addition in equation 25 is also 4 times larger than the second term of the addition in equation 26. Because both terms of the addition in equation 25 are four times larger than those in equation 26, we can conclude that the variance of Cohen's  $d_s$  is four times larger than the variance of Shieh's  $d_s$ .