

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

2                               Marie Delacre<sup>1</sup> & Christophe Leys<sup>1</sup>

3   <sup>1</sup> Université Libre de Bruxelles, Service of Analysis of the Data (SAD), Bruxelles, Belgium

4                               Author Note

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

7

Abstract

8

9

*Keywords:* keywords

10

Word count: X

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

## Intro

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

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

Facing this argument, it has become an advised practice to report the *p*-value assorted by a measure of the effect size, that is, a quantitative measure of the magnitude of the experimental effect (Cohen, 1965; Fan, 2001; Hays, 1963). This practice is also highly endorsed by the American Psychological Association (APA) and the American Educational Research Association (AERA) (American Educational Research Association, 2006; American Psychological Association, 2010). However, limited studies properly report effect size in the

---

<sup>1</sup> Tiny differences might be due to sampling error, or to other factors than the one of interest: even under the assumption of random assignment (which is a necessary but not sufficient condition), it is almost impossible to be sure that the only difference between two conditions is the one defined by the factor of interest. Other tiny factors of no theoretical interest might slightly influence results, making the probability of getting an actual zero effect very low. This is what Meehl (1990) calls 'systematic noise'.

last several decades.

Generally, there is a high confusion between the effect size and other related concepts such as the Applied significance<sup>2</sup>. Moreover, there are several situations that call for effect size measures and in the current literature, it's not always easy to know which measure using in a specific context. We will therefore introduce this paper with 3 sections in which we will:

1. Clearly define what is a measure of effect size;
2. Listing the different situations that call for effect sizes measure;
3. Define required properties as a function of the situations.

After these general adjustments, we will focus our attention on “between-subject” designs where individuals are randomly assigned into one of two independent groups and groups scores are compared based on their means<sup>3</sup>. Because it has been widely argued that there are many fields in psychology where the assumption of equal variances between groups is ecologically unlikely (Erceg-Hurn & Mirosevich, 2008; Grissom, 2000), it is becoming more common in statistical software to present a *t*-test that does not hold on this assumption by default, namely the Welch's *t*-test (e.g., R, Minitab). However, similar issues for the measures of effect sizes has received less attention (Shieh, 2013), and Cohen's  $d_s$  remains persistent<sup>4</sup>. One possible reason is that researchers cannot find a consensus on which alternative should be in use (Shieh, 2013). We will limit our study to the standardized mean difference, called the *d*-family, because it is the dominant family of estimators of effect size when comparing two groups based on their means (Peng, Chen, Chiang, & Chiang, 2013; Shieh, 2013), and we will see that even in this very specific context, there is little agreement between researchers as to which is the most suitable estimator. According to us, the main

---

<sup>2</sup> In our conception Applied significance" encompass all what refers to the relevance of an effect in real life, e.g. clinical, personal, social, professional.

<sup>3</sup> We made this choice because *t*-tests are still the most commonly used tests in the field of Psychology.

<sup>4</sup> For example, in Jamovi, Cohen's  $d_s$  is provided, whatever one performs Student's or Welch's *t*-test.

reason is that it is difficult, based on currently existing measures, to optimally serve all the purposes of an effect size measure. Throughout this section, we will:

1. Present the main measures of the  $d$ -family that are proposed in the literature, related to the purpose they serve, and introduce a new one, namely the “transformed Shieh’s  $d$ ” that should help at reaching all the purposes simultaneously;
2. Present and discuss the results of simulations we performed, in order to compare existing measures and the new introduced one;
3. Summarize our conclusions in practical recommendations.

### Measure of effect size: what it is, what it is not

The effect size is commonly referred to the practical significance of a test. Grissom & Kim (2005) define the effect size as the extent to which results differ from what is implied by the null hypothesis. In the context of the comparison of two groups based on their mean, depending on the defined null hypothesis (considering the absence of effect as the null hypothesis), we could define the effect size either as the magnitude of differences between parameters of two populations groups are extracted from (e.g. the mean; Peng & Chen, 2014) or as the magnitude of the relation between one dichotomous factor and one dependent variable (American Educational Research Association, 2006). Both definitions refer to as the most famous families of measures of effect sizes (Rosenthal, 1994): respectively the  $d$ -family and the  $r$ -family.

Very often, the contribution of the measures of effect size is overestimated. First, benchmarks about what should be a small, medium or large effect size might have contributed at seeing the effect size as a measure of the importance or the relevance of an effect in real life, but it is not (Stout & Ruble, 1995). The effect size is only a mathematical indicator of the magnitude of a difference, which depends on the way a variable is converted into numerical indicator. In order to assess the meaningfulness of an effect, we should be able to relate this effect with behaviors/meaningful consequences in the real world

(Andersen, McCullagh, & Wilson, 2007). For example, let us imagine a sample of students in serious school failure who are randomly divided into two groups: an experimental group following a training program and a control group. At the end of the training, students in the experimental group have on average significantly higher scores on a test than students in the control group, and the difference is large (e.g. 30 percents). Does it mean that students in the experimental condition will be able to pass to the next grade and to continue normal schooling? Whether the computed magnitude of difference is an important, meaningful change in everyday life refers to another construct: the *Applied significance* (Bothe & Richardson, 2011). It refers to the interpretation of treatment outcomes and is neither statistical nor mathematical, it is related to underlying theory that posits an empirical hypothesis. In other words, the relation between *practical* and *Applied* significance is more a theoretical argument than a statistical one.

Second, in the context of the comparison of two groups based on their means, it should not replace the null hypothesis testing. Statistical testing allows the researcher to determine whether the observed departure from  $H_0$  occurred by chance or not (Stout & Ruble, 1995) while effect size estimators allow to assess the practical significance of an effect, and as reminds Fan (2001): “*a practically meaningful outcome may also have occurred by chance, and consequently, is not trustworthy*” (p.278). For this reason, the use of confidence intervals around the effect size estimate is highly recommended (Bothe & Richardson, 2011).

### Different purposes of effect size measures

Effect size measures can be used in an *inferential* perspective:

- The effect sizes from previous studies can be used in 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);
- We can also compute confidence limits around the point estimator (Shieh, 2013), in order to replace conventional hypothesis testing : if the null hypothesis area is out of the

confidence interval, we can conclude that the null hypothesis is false.

Measures of effect size can also be used in a *comparative* perspective, that is to assess the stability of results across designs, analysis, samples sizes (Wilkinson & the Task Force on Statistical Inference, 1999). It includes:

- To compare results of 2 or more studies (Prentice & Miller, 1990);
- To incorporate results in meta-analysis (Lakens, 2013; Li, 2016; Nakagawa & Cuthill, 2007; Stout & Ruble, 1995; Wilkinson & the Task Force on Statistical Inference, 1999).

Finally, effect size measures can be used for *interpretative* purposes: in order to assess the practical significance of a result (beyond statistical significance; Lakens, 2013; American Psychological Association, 2010; Prentice & Miller, 1990).

### Properties of a good effect size estimator

The estimate of an estimator depends on the sampling. That is to say, based on different samples extracted from the same population, one would obtain different estimates of the same estimator. 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 very useful, as it allows to assess the goodness of an effect size estimator and more specifically, three desirable properties of a good estimator for inferential purposes: **unbiasedness**, **consistency** and **efficiency**.

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 smaller) than the true population parameter (see Figure 1). In other words, the bias tells us if estimates are good, on average. The *bias* of a point estimator  $\hat{\delta}$  can be computed as follows:

$$\hat{\delta}_{bias} = E(\hat{\delta}) - \delta \quad (1)$$

Where  $E(\hat{\delta})$  is the expectancy of the sampling distribution of the estimator (i.e. the average estimate) and  $\delta$  is the true parameter.

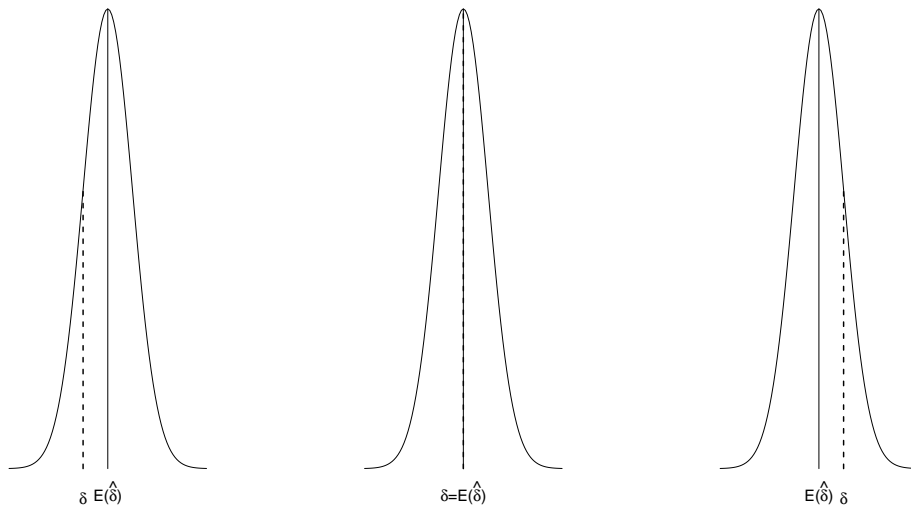


Figure 1. Sampling distribution for a positively biased (left), an unbiased (center) and a negatively biased estimator (right)

Bias informs us about the goodness of estimates averages, but 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. It would be problematic, as long as we have only one estimate, the one based on our sample, and we don't know how far is this estimate from the center of the sampling distribution. We hope that *all* possible estimates are close enough of the true population parameter, in order to be sure that for *any* estimate, one has a correct estimation of the real parameter. In other words, we expect the variability of estimates around the true population parameter to be as small as possible. It refers to the **efficiency** of the point estimator ( $\hat{\delta}$ ) and can be computed as follows:



$$\hat{\delta}_{efficiency} = Var(\hat{\delta}) \quad (2)$$

Among all unbiased estimators, the more efficient will be the one with the smallest variance. Note that both unbiasedness and efficiency are very important. An unbiased estimator with such a large variance that some estimates are extremely far from the real parameter is as undesirable as a parameter which is highly biased. In some situations, it is better to have a very slightly biased estimator with a tight shape around the biased value, so each estimate remains relatively close to the true parameter, than an unbiased estimator with a large variance (Raviv, 2014). Because both *unbiasedness* and *efficiency* must be considered, it is interesting to compute an indicator that take simultaneously both properties into account (Wackerly, Mendenhall, & Scheaffer, 2008). The *mean square error* of a point estimator  $\hat{\delta}$  is defined as follows:

$$MSE(\hat{\delta}) = E[(\hat{\delta} - \delta)^2] \quad (3)$$

It can be proven that the *mean square error* is a function of the bias and the variance of  $\hat{\delta}$ :

$$MSE(\hat{\delta}) = \hat{\delta}_{efficiency} + \hat{\delta}_{bias}^2 \quad (4)$$

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

Beyond the inferential properties, Cumming (2013) reminds that an effect size estimator need to have a constant value across designs in order to be easily interpretable and to be included in meta-analysis. In other word, it should achieve the property of **generality**.

## Different measures of effect sizes TOUCHER UN MOT SUR L'IMPACT DE LA DIFFERENCE DE MOYENNE.

The  $d$ -family effect sizes are commonly used with “between-subject” designs where individuals are randomly assigned into one of two independent groups and groups scores means are compared. The population effect size is defined as follows:

$$\delta = \frac{\mu_1 - \mu_2}{\sigma} \quad (5)$$

They exist different estimators of this effect size measure varying as a function of the chosen standardizer ( $\delta$ ). For all estimators, the mean difference is estimated by the difference of both sample means ( $\bar{X}_1 - \bar{X}_2$ ). When used for inference, some of them rely on both assumptions of normally distributed residuals and equality of variances, while others rely solely on the normally distributed residuals assumption.

### Alternatives when variances are equal between groups

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):

$$Cohen's d_s = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{\frac{(n_1-1) \times SD_1 + (n_2-1) \times SD_2}{n_1 + n_2 - 2}}} \quad (6)$$

The reasoning behind this measure is that considering both samples as extracted from a common population variance (n.d.), we achieve a more accurate estimation of the population variance by pooling both estimates of this parameter (i.e  $SD_1$  and  $SD_2$ ) and because the larger the sample size, the more accurate the estimate, we give more weight to the estimate based on the larger sample size ( $\max(n_j)$ ). Unfortunately, even under the assumptions that residuals are independent and identically normally distributed with the same variance across

groups, Cohen's  $d_s$  is known to be positively biased (Lakens, 2013) and for this reason, Hedges & Olkin (1985) has defined a bias-corrected version, which is referred to:

$$\text{Hedge's } g_s = \text{Cohen's } d_s \times \left(1 - \frac{3}{4 \times (n_1 + n_2) - 9}\right) \quad (7)$$

The pooled SD is the best choice when variances are equal between groups (Grissom & Kim, 2001) but they may not be well advised for use with data that violates this assumption (Cumming, 2013; Grissom & Kim, 2001, 2005; Kelley, 2005, 2005; Shieh, 2013). In case of a positive pairing (i.e. the group with the larger sample size also has the larger variance), the variance will be over-estimated and therefore, the estimator will be lower as it should be. On the other side, in case of negative pairing (i.e. the group with the larger sample size has the smaller variance), the estimator will be larger as it should be. Because the assumption of equal variances across groups is very rare in practice (Cain, Zhang, & Yuan, 2017; Delacre, Lakens, & Leys, 2017; Delacre, Leys, Mora, & Lakens, 2019; Erceg-Hurn & Mirosevich, 2008; Glass, Peckham, & Sanders, 1972; Grissom, 2000; Micceri, 1989; Yuan, Bentler, & Chan, 2004), both Cohen's  $d_s$  and Hedge's  $g_s$  should be abandoned in favor of a robust alternative to unequal variances.

### Alternatives when variances are unequal between groups

In his review, Shieh (2013) mention three alternative available in the literature: the sample mean difference, divided by the non pooled average of both variance estimates (A), the Glass's  $d_s$  (B) and the Shieh's  $d_s$  (C).

The **sample mean difference, divided by the non pooled average of both variance estimates** (A) was suggested by (Cohen, 1988). We immediately exclude this alternative because it suffers of many limitations:

- it results in a variance term of an artificial population and is therefore very difficult to interpret (Grissom & Kim, 2001);

- unless both sample sizes are equal, the variance term does not correspond to the variance of the mean difference (Shieh, 2013);

- unless the mean difference is null, the measure is biased. Moreover, the bigger the sample size, the larger the variance around the estimate.

When comparing one control group with one experimental group, Glass, McGav, & Smith (2005) recommend using the SD of the control group as standardizer. It is also advocated by Cumming (2013), because according to him, it is what makes the most sense, conceptually speaking.

$$Glass's\ d_s = \frac{\bar{X}_{experimental} - \bar{X}_{control}}{SD_{control}} \quad (8)$$

One argument in favour of using the standard deviation of the control group 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. However, defining this group is not always obvious (Coe, 2002). This could induce large ambiguity because depending of the chosen standard deviation, as standardizer, measures could be substantially different (Shieh, 2013). Moreover, glass  $d_s$  also have limitations when used for inference. While it is a consistant measure, it can be shown that it can be highly positively biased when there are less than 300 participants (Hedges, 1981; Olejnik & Hess, 2001), especially for small effect sizes.

Kulinskaya and Staudte (2007) adviced the use of a standardizer that take the sample sizes allocation ratios into account, in addition to the variance of both groups. It results in a modification of the exact  $SD$  of the sample mean difference:

$$Shieh's d_s = \frac{\bar{X}_1 - \bar{X}_2}{\sqrt{SD_1^2/q_1 + SD_2^2/q_2}} \quad (9)$$

According to the statistical properties of Welch's statistic under heteroscedasticity, it does not appear possible to define a proper standardised effect size without accounting for the relative group size of subpopulations in a sampling scheme. At the same time, 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 keeping constant the mean difference as well as  $SD_1$  and  $SD_2$ , Shieh's  $d_s$  will vary as a function of the sample sizes allocation ratio (dependency of Shieh's  $d_s$  value on the sample sizes allocation ratio is detailed and illustrated in Appendix 1, and also in the following shiny application: <http://127.0.0.1:3461/>).

Fortunately, this paradox can be resolved. It is possible to find a modified measure of Shieh's  $d_s$  that does not depend on sample sizes ratio, in answering the following question: “whatever the real sample sizes ratio, what value of Shieh's  $d_s$  would have been computed if design were balanced (i.e.  $n_1 = n_2$ ), keeping all other parameters constant?”

It can be shown that the relationship between Shieh's  $\delta_s$  when samples sizes are equal between groups (i.e.  $\delta_{Shieh, n_1=n_2}$ ) and Shieh's  $\delta_s$  for any other sample sizes allocation ratios can be expressed as follows:

$$\delta_{Shieh, n_1=n_2} = \delta_{Shieh} \times \frac{(nratio + 1) \times \sigma_{n_1 \neq n_2}}{2 \times \sigma_{n_1=n_2} \times \sqrt{nratio}} \quad (10)$$

With

$$\sigma_{n_1=n_2} = \sqrt{\frac{\sigma_1^2 + \sigma_2^2}{2}}$$

and

$$SD_{n_1 \neq n_2} = \sqrt{(1 - \frac{n_1}{N}) \times \sigma_1^2 + (1 - \frac{n_2}{N}) \times \sigma_2^2}$$

$\delta_{Shieh, n_1=n_2}$  can therefore be estimated using this equation:

$$Shieh's\ ds_{n_1=n_2} = Shieh's\ ds \times \frac{(nratio + 1) \times SD_{n_1 \neq n_2}}{2 \times SD_{n_1=n_2} \times \sqrt{nratio}} \quad (11)$$

With

$$SD_{n_1=n_2} = \sqrt{\frac{SD_1^2 + SD_2^2}{2}}$$

and

$$SD_{n_1 \neq n_2} = \sqrt{(1 - \frac{n_1}{N}) \times SD_1^2 + (1 - \frac{n_2}{N}) \times SD_2^2}$$

$Shieh's\ ds_{n_1=n_2}$  can be compared across two different studies using different sample sizes allocation ratio and could be included in meta-analysis.

## Monte Carlo Simulations

### Simulation 1: 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  $d_s$ , Hedge's  $g_s$ , Glass's  $d_s$  (using respectively the sample standard deviation of the first or second group as a standardizer), Shieh's  $d_s$  and our transformed measure of Shieh's  $d_s$ , that we will note later  $d_s^*$ .

100,000 datasets were generated for 1,260 scenarios. In 315 scenarios, samples were extracted from a normal population distribution and in 945 scenarios, samples were extracted from non normal population distributions. In order to assess the goodness of estimators under realistic deviations from the normality assumption, we referred to the review of Cain et al. (2017). Based on their investigation<sup>5</sup>, Cain et al. (2017) found values of

---

<sup>5</sup> 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 the Fisher's skewness (G1) and kurtosis (G2).

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, i.e.  $G2=95.75$ . Regarding skewness, we simulated population parameter values which correspond to the 1st and 99th percentile of their distribution, i.e. respectively  $G1 = -2.08$  and  $G1 = 6.32$ . We also simulated null population parameter values (i.e.  $G1 = 0$ ), in order to assess the main effect of high kurtosis on the goodness of estimators. All possible combinations of skewness and kurtosis and the number of scenarios for each combination are summarized in Table 1.

Table 1. *Number of Combinations of skewness and kurtosis in our simulations*

		<b>Kurtosis</b>		
		0	95.75	<b>TOTAL</b>
<b>Skewness</b>	0	315	315	<b>630</b>
	-2.08	/	315	<b>315</b>
	6.32	/	315	<b>315</b>
	<b>TOTAL</b>	<b>315</b>	<b>945</b>	<b>1260</b>

*Note.* Fisher's skewness ( $G1$ ) and kurtosis ( $G2$ ) are presented in Table 1. The 315 combinations where both  $G1$  and  $G2$  equal 0 correspond to the normal case.

For the 4 resulting combinations of skewness and kurtosis (see Table 1), all other population parameter values were chosen in order to illustrate the consequences of factors known to play a key role on goodness of 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  $SD$ -ratio ( $SD\text{-ratio} = \frac{\sigma_1}{\sigma_2}$ ), and the sample size and variance pairing. In our scenarios,  $\mu_2$  was always 0 and  $\mu_1$  varied from 0 to 4, in step of 1 (so does  $\mu_1 - \mu_2$ ). Moreover,  $\sigma_1$  always equals 1, and  $\sigma_2$  equals .1, .25, .5, 1, 2, 4 or 10 (so does the  $SD$ -ratio). The simulations for which both  $\sigma_1$  and  $\sigma_2$  equal 1 are the particular case of homoscedasticity (i.e. equal variances across groups). Sample size of both groups ( $n_1$  and  $n_2$ ) were 20, 50 or 100. When sample sizes of both groups are equal, the  $n$ -ratio equals 1 (it is known as a balanced design). All possible combinations of  $n$ -ratio and  $SD$ -ratio were performed in order to distinguish positive pairings (the group with the largest sample size is extracted from the population with the largest  $SD$ ), negative pairings (the group with the smallest sample size is extracted from the population with the smallest  $SD$ ), and no pairing (sample sizes and/or population  $SD$  are equal across all groups). In sum, the simulations grouped over different sample sizes yield 5 conditions based on the  $n$ -ratio,  $SD$ -ratio, and sample size and variance pairing, as summarized in Table 2.

Table 2. 5 conditions based on the  $n$ -ratio,  $SD$ -ratio, and sample size and variance pairing

		<b><math>n</math>-ratio</b>		
		<b>1</b>	<b>&gt;1</b>	<b>&lt;1</b>
<b><math>SD</math>-ratio</b>	<b>1</b>	a	b	b
	<b>&gt;1</b>	c	d	e
	<b>&lt;1</b>	c	e	d

*Note.* The  $n$ -ratio is the sample size of the first group ( $n_1$ ) divided by the sample size of the second group ( $n_2$ ). When all sample sizes are equal across groups, the  $n$ -ratio equals 1.



When  $n_1 > n_2$ ,  $n\text{-ratio} > 1$ , and when  $n_1 < n_2$ ,  $n\text{-ratio} < 1$ .  $SD\text{-ratio}$  is the population  $SD$  of the first group ( $\sigma_1$ ) divided by the population  $SD$  of the second group ( $\sigma_2$ ). When  $\sigma_1 = \sigma_2$ ,  $SD\text{-ratio} = 1$ . When  $\sigma_1 > \sigma_2$ ,  $SD\text{-ratio} > 1$ . Finally, when  $\sigma_1 < \sigma_2$ ,  $SD\text{-ratio} < 1$ .

**Results.** In all Figures presented below, averaged results for each sub-condition are presented under five different configurations of distributions, using the legend described in Figure 2.

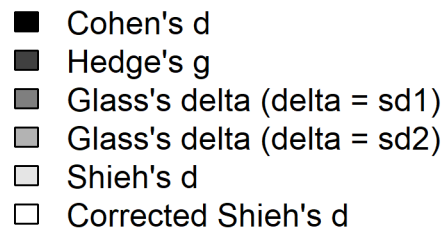


Figure 2. Legend

Figure 3 shows that when variances and sample sizes are equal between groups, estimator bias tends to decrease and precision is also improved with increasing sample sizes, meaning that all estimators are consistent.

Among all effect size indicators, glass's  $d_s$  shows least precision and highest bias rates. Shieh's  $d_s$  and Shieh's  $d_s^*$  are identical (because  $n_1 = n_2$ ) and perform the lowest bias rates and precision. We can easily demonstrate that their bias is exactly half the size of Cohen's  $d_s$ , and that their variance is exactly four times smaller than Cohen's  $d_s$ . Due to the relation described in equation 12, when sample sizes are equal between groups, one can deduce from

these proportions that relative to their respective true effect size, Cohen's  $d_s$ , Shieh's  $d_s$  and  $d_s^*$  perform all as well (see the second row in Figure 3).

$$\delta_{Shieh, n_1=n_2} = \frac{\delta_{Cohen, n_1=n_2}}{2} \quad (12)$$

Moreover, details in Supplemental Material reveals that in our simulations, all effect size estimators tend to be more biased and variable when the difference between two groups means enlarges. However, the relative bias (i.e. the size of the bias relative to the true effect size) is quite stable, except for the glass's  $d_s$  when distribution are highly skewed ( $G_1 = 6.32$ ): in that case, keeping all other parameters constant, the relative bias of glass's  $d_s$  using  $sd_1$  as standardizer decreases when the mean difference enlarges, and the relative bias of glass's  $d_s$  using  $sd_2$  as standardizer increases when the mean difference enlarges (when  $\mu_1 - \mu_2$  varies from 1 to 4, one can observe variations as big as approximately 13% with 20 subjects per group).

Figure 4 shows that when variances are equal but sample sizes are unequal between groups, as when sample sizes were equal, glass's  $d_s$  is more biased and variable than all other estimators. As previously, the bias of Shieh's  $d_s$  is smaller than the Cohen's  $d_s$  one. However, the difference is much smaller than previously. Remember that when sample sizes differ between groups, Shieh's  $d_s$  is always more than twice smaller than Cohen's  $d_s$  (see Appendix 1 for more details). As a consequence, if both Cohen's  $d_s$  and Shieh's  $d_s$  performed as well, the bias of Shieh's  $d_s$  should be more than twice smaller than Cohen's  $d_s$  bias, but it's not. It's confirmed by the second row in Figure 4 where we can see that the relative bias of Shieh's  $d_s$  is larger than the relative bias of Cohen's  $d_s$ . (cf. notes en vue de la réunion avec Christophe à ce sujet).

Moreover, our transformed Shieh's  $d_s^*$  is on average less biased than original Shieh's  $d_s$ , both if raw and relative terms, but slightly more variable in order that the mean square error

of both estimators are quite similar.

Details in Supplemental Material reveals that in our simulations, all effect size estimators tend to be more biased and variable when the difference between two groups means enlarges (this is especially true for the Glass  $d_s$  when one chooses the  $sd$  of the smallest group as a standardizer). However, the relative bias of all estimators is quite stable, except for the glass's  $d_s$ , and in a lesser extent for Shieh's  $d_s$  when distribution are highly skewed ( $G_1 \neq 0$ ). The relative bias of glass's  $d_s$  using  $sd_1$  as standardizer decreases when the mean difference enlarges, and the relative bias of glass's  $d_s$  using  $sd_2$  as standardizer increases when the mean difference enlarges: when  $\mu_1 - \mu_2$  varies from 1 to 4, one can observe variations as big as approximately 13% when the standardizer is computed based on 20 subjects (i.e. when  $n_{stdizer} = 20$ ). About Shieh's  $d_s$ , one observes variations until approximately 6%.

Figure 5 shows that when variances are unequal but sample sizes are equal between groups, ...

Again, because  $n_1 = n_2$  the bias and variance of Shieh's  $d_s$  and  $d_s^*$  are respectively two times and for times smaller than bias and variance of Cohen's  $d_s$

Figure 6 shows that when variances are unequal, and the largest group is associated with largest variance, ...

Figure 7 shows that when variances are unequal, and the largest group is associated with smallest variance, ...

### ***Conclusion.***

### **Simulation 2: confidence intervals.**

### ***Method.***

### ***Results.***

**Conclusion.** American Educational Research Association. (2006). Standards for reporting on empirical social science research in aera publications. *Educational Researcher*, 35, 33–40. doi:10.3102/0013189X035006033

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

Andersen, M. B., McCullagh, P., & Wilson, G. J. (2007). But what do the numbers really tell us? Arbitrary metrics and effect size reporting in sport psychology research. *Journal of Sport & Exercise Psychology*, 29, 664–672.

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

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. doi: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.

Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (Routledge Academic.). New York, NY.

Cumming, G. (2013). Cohen's d needs to be readily interpretable: Comment on shieh (2013). *Behavior Research Methods*, 45, 968–971. doi: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. doi: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. doi:http://doi.org/10.5334/irsp.198

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. doi:10.1037/0003-066X.63.7.591

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

Glass, G. V., 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. doi:10.3102/00346543042003237

Grissom, R. J. (2000). Heterogeneity of variance in clinical data. *Journal of Consulting and Clinical Psychology*, 68(1), 155–165. doi: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. doi:10.1037/1082-989X.6.2.135

Grissom, R. R., & Kim, J. J. (2005). *Effect size for research: A broad practical*

approach. (Lawrence Erlbaum Associates, Mahwah, N.J.). London.

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

Hedges, L. V. (1981). Distribution theory for glass's estimator of effect size and related estimators. *Journal of Educational Statistics*, 6(2), 107–128.

Hedges, L. V., & Olkin, I. (1985). *Statistical methods for meta-analysis* (Academic Press.). Cambridge, Massachusetts. doi:10.1016/C2009-0-03396-0

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

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. doi:10.1177/0013164404264850

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

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. doi:10.3389/fpsyg.2013.00863

Li, J. (2016). Effect size measures in a two-independent-samples case with nonnormal and nonhomogeneous data. *Behavior Research Methods*, 48(4), 1560–1574. doi:10.3758/s13428-015-0667-z

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

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

Micceri, T. (1989). The unicorn, the normal curve, and other improbable creatures. *Psychological Bulletin*, 105(1), 156–166. doi:10.1037/0033-2909.105.1.156

Nakagawa, S., & Cuthill, I. C. (2007). Effect size, confidence interval and statistical significance: A practical guide for biologists. *Biological Reviews*, 82, 591–605. doi:10.1111/j.1469-185X.2007.00027.x

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

Olejnik, S., & Hess, B. (2001). *Revisiting the efficacy of glass's estimator of effect size for program impact analysis*. Retrieved from <https://eric.ed.gov/?id=ED452210>

Peng, C.-Y., & 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. doi:10.1080/00220973.2012.745471

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

Prentice, D., & 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/>

Rosenthal, R. (1994). Parametric measures of effect size. In H. Cooper & L. V. Hedges

(Eds.), *The hand-book of research synthesis* (pp. 231–244). New-York: Sage.

Shieh, G. (2013). Confidence intervals and sample size calculations for the weighted eta-squared effect sizes in one-way heteroscedastic ANOVA. *Behavior Research Methods*, 45(1), 2–37. doi:10.3758/s13428-012-0228-7

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

Stout, D. D., & Ruble, T. L. (1995). Assessing the practical significance of empirical results in accounting education research: The use of effect size information. *Journal of Accounting Education*, 13(3), 281–298.

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

Wackerly, D. D., Mendenhall, W., & Scheaffer, R. L. (2008). *Mathematical statistics with applications (7th edition)* (Brooks/Cole, Cengage Learning.). Belmont, USA.

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. doi:10.1007/bf02295644

(n.d.).

American Educational Research Association. (2006). Standards for reporting on empirical social science research in aera publications. *Educational Researcher*, 35, 33–40. doi:10.3102/0013189X035006033

American Psychological Association. (2010). *Publication manual of the american*



psychological association [apa] (6 ed.) (American Psychological Association.). Washington, DC:

Andersen, M. B., McCullagh, P., & Wilson, G. J. (2007). But what do the numbers really tell us? Arbitrary metrics and effect size reporting in sport psychology research. *Journal of Sport & Exercise Psychology*, 29, 664–672.

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

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. doi: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.

Cohen, J. (1988). *Statistical power analysis for the behavioral sciences* (Routledge Academic.). New York, NY.

Cumming, G. (2013). Cohen's d needs to be readily interpretable: Comment on shieh (2013). *Behavior Research Methods*, 45, 968–971. doi: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. doi: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.  
doi:http://doi.org/10.5334/irsp.198

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. doi:10.1037/0003-066X.63.7.591

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

Glass, G. V., 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. doi:10.3102/00346543042003237

Grissom, R. J. (2000). Heterogeneity of variance in clinical data. *Journal of Consulting*  
*and Clinical Psychology*, 68(1), 155–165. doi: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.  
doi:10.1037/1082-989X.6.2.135

Grissom, R. R., & Kim, J. J. (2005). *Effect size for research: A broad practical*  
*approach*. (Lawrence Erlbaum Associates, Mahwah, N.J.). London.

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

Hedges, L. V. (1981). Distribution theory for glass's estimator of effect size and related

estimators. *Journal of Educational Statistics*, 6(2), 107–128.

Hedges, L. V., & Olkin, I. (1985). *Statistical methods for meta-analysis* (Academic Press.). Cambridge, Massachusetts. doi:10.1016/C2009-0-03396-0

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

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. doi:10.1177/0013164404264850

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

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. doi:10.3389/fpsyg.2013.00863

Li, J. (2016). Effect size measures in a two-independent-samples case with nonnormal and nonhomogeneous data. *Behavior Research Methods*, 48(4), 1560–1574. doi:10.3758/s13428-015-0667-z

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

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

Micceri, T. (1989). The unicorn, the normal curve, and other improbable creatures. *Psychological Bulletin*, 105(1), 156–166. doi:10.1037/0033-2909.105.1.156

Nakagawa, S., & Cuthill, I. C. (2007). Effect size, confidence interval and statistical significance: A practical guide for biologists. *Biological Reviews*, 82, 591–605. doi:10.1111/j.1469-185X.2007.00027.x

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

Olejnik, S., & Hess, B. (2001). *Revisiting the efficacy of glass's estimator of effect size for program impact analysis*. Retrieved from <https://eric.ed.gov/?id=ED452210>

Peng, C.-Y., & 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. doi:10.1080/00220973.2012.745471

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

Prentice, D., & 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/>

Rosenthal, R. (1994). Parametric measures of effect size. In H. Cooper & L. V. Hedges (Eds.), *The hand-book of research synthesis* (pp. 231–244). New-York: Sage.

Shieh, G. (2013). Confidence intervals and sample size calculations for the weighted eta-squared effect sizes in one-way heteroscedastic ANOVA. *Behavior Research Methods*, 45(1), 2–37. doi:10.3758/s13428-012-0228-7

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

Stout, D. D., & Ruble, T. L. (1995). Assessing the practical significance of empirical results in accounting education research: The use of effect size information. *Journal of Accounting Education*, 13(3), 281–298.

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

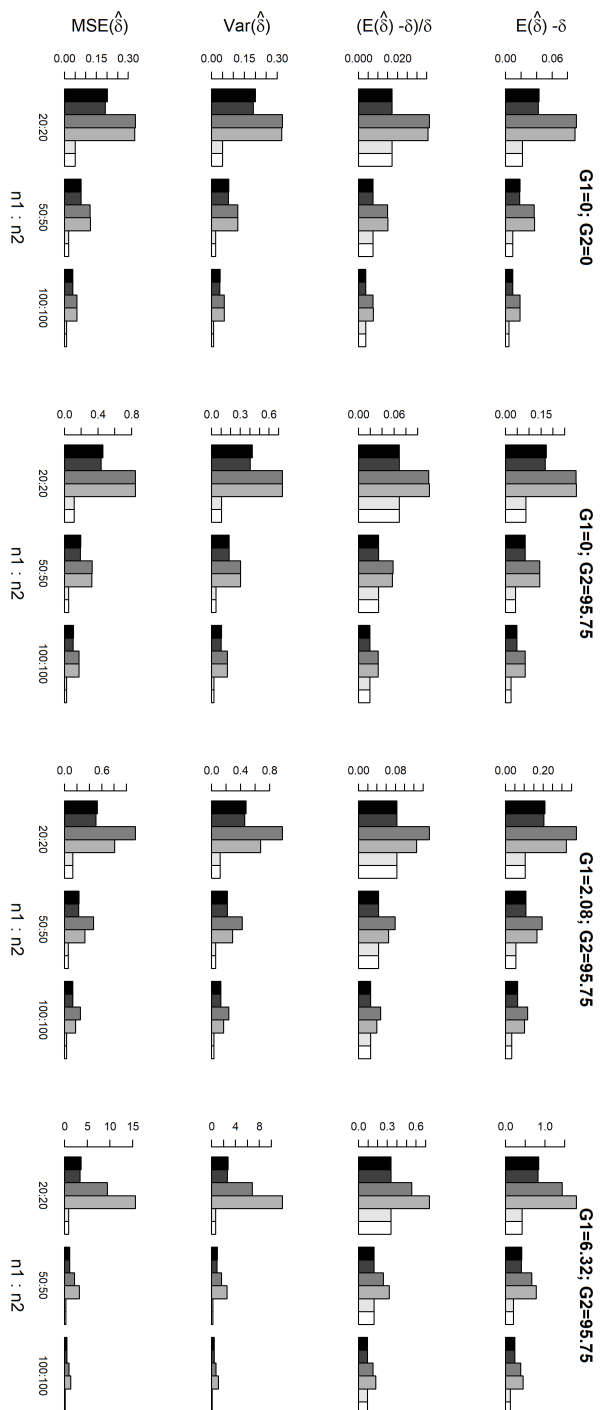
Wackerly, D. D., Mendenhall, W., & Scheaffer, R. L. (2008). *Mathematical statistics with applications (7th edition)* (Brooks/Cole, Cengage Learning.). Belmont, USA.

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. doi:10.1007/bf02295644

(n.d.).

Figure 3. Bias and efficiency of five estimator of standardized mean difference, when variances and sample sizes are equal across groups



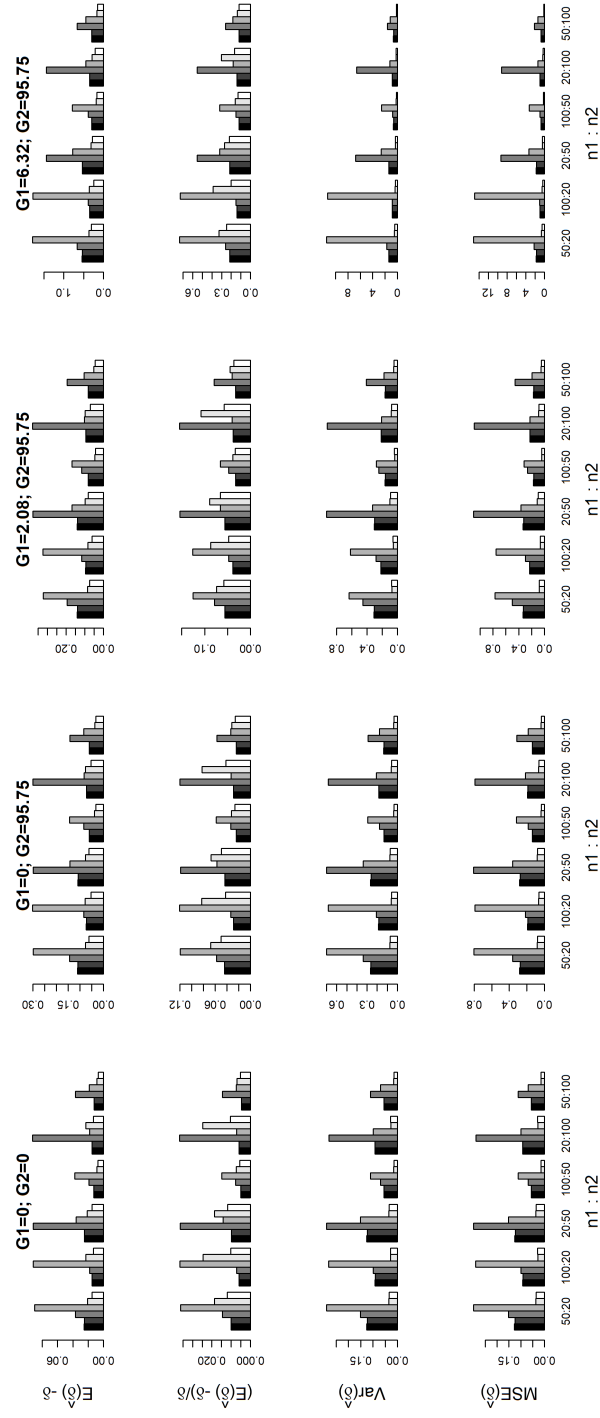
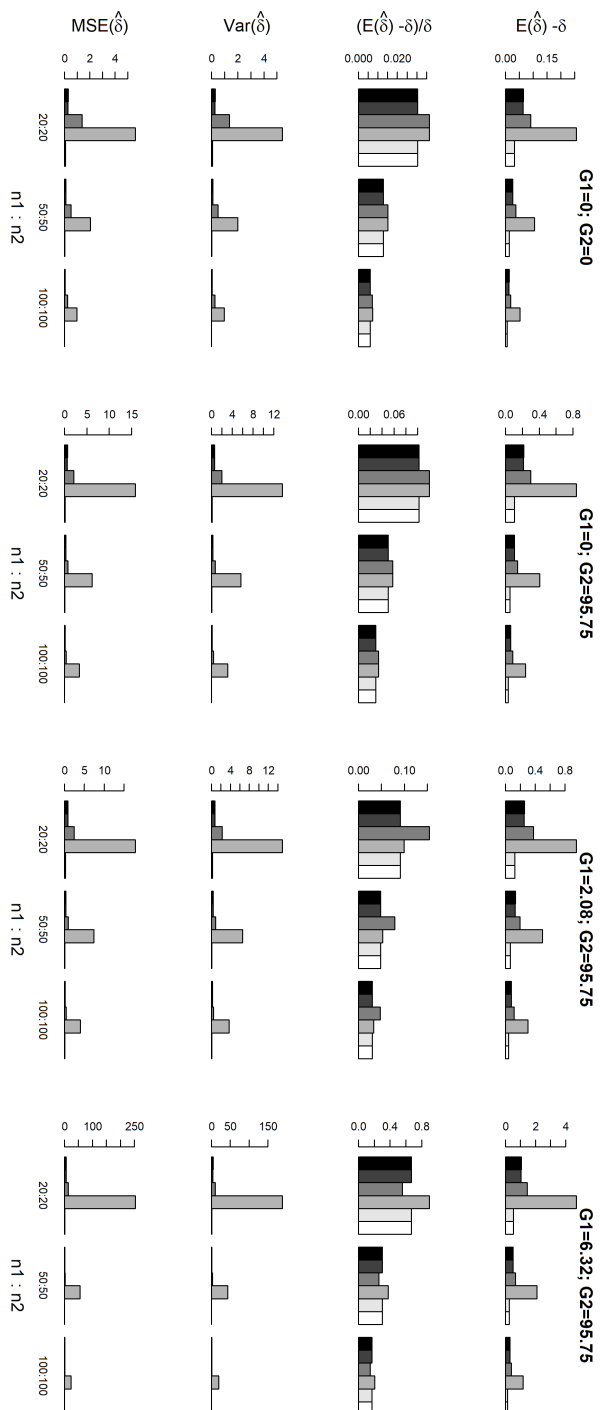


Figure 4. Bias and efficiency of five estimator of standardized mean difference, when variances are equal across groups and sample sizes are unequal

Figure 5. Bias and efficiency of five estimator of standardized mean difference, when variances are unequal across groups and sample sizes are equal





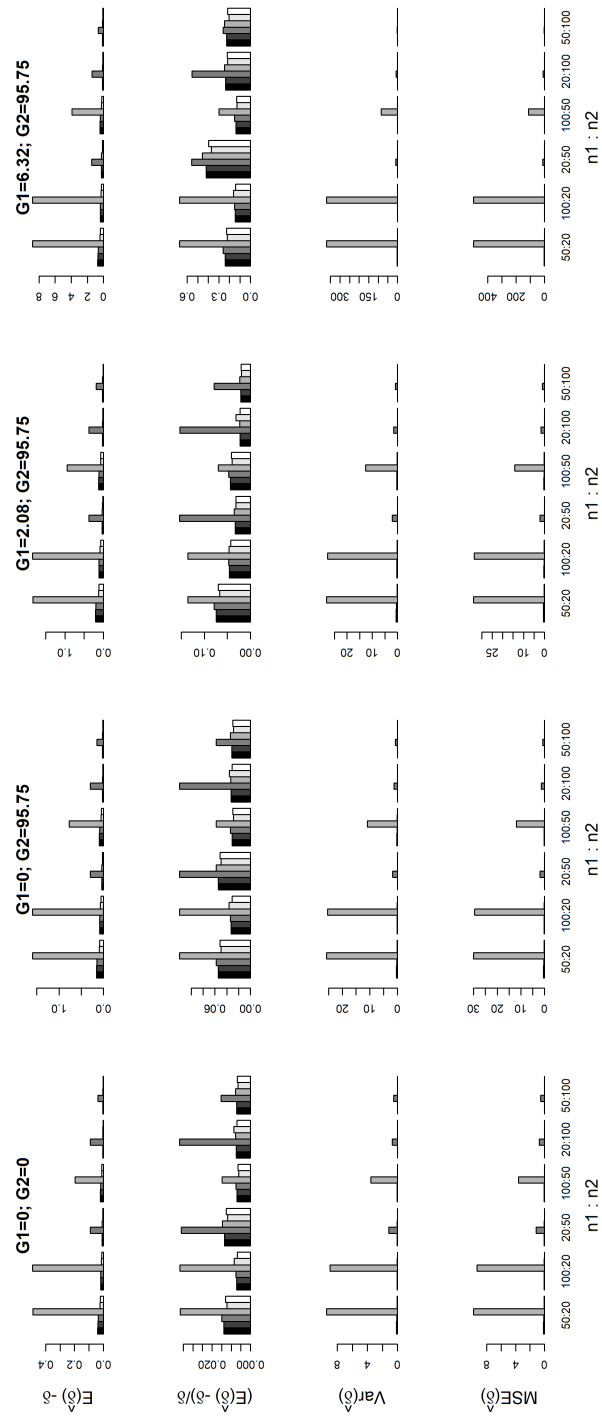


Figure 6. Bias and efficiency of five estimator of standardized mean difference, when variances and sample sizes are unequal across groups, with positive correlation between them

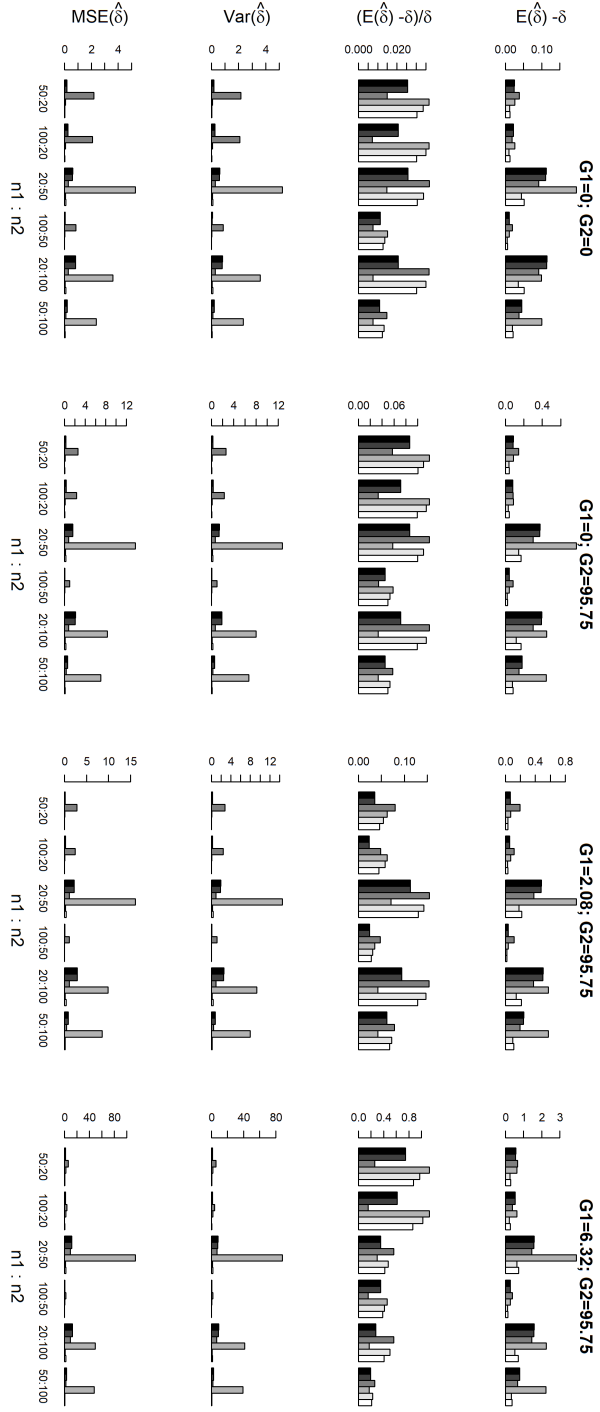


Figure 7. Bias and efficiency of five estimator of standardized mean difference, when variances and sample sizes are unequal across groups, with negative correlation between them

## Appendix

568 **Appendix 1: The mathematical study of Shieh's  $\delta$**

569 Paste Appendix 1 when it will be finished

570 **Appendix 2: Confidence intervals**

571 Paste Appendix 2 when it will be finished

572 **Appendix 3: a priori power analyses**

573 Paste Appendix 3 when it will be finished (Cumming & Finch, 2001)

574 Cumming, G., & Finch, S. (2001). A primer on the understanding, use, and calculation of  
575 confidence intervals that are based on central and noncentral distributions. *Educational and*  
576 *Psychological Measurement*, 61(532), 532–574.