



Figure 1: Estimated power of Levene's test as a function of sample size, SDR and centering parameter.

are slightly above power curves of the Levene's test based on the median, meaning that it leads to slightly higher power than Levene's test based on the median. This can be due to the fact that data is extracted from normal distributions. With asymmetric data, the median would perform better. When $SDR = 2$, approximately 50 subjects are needed to have 80 percent power to detect differences, while approximately 70 subjects are needed to have 95 percent power to detect differences (for both versions of Levene's test). To detect an SDR of 1.5 with Levene's test, approximately 120 subjects are needed to reach a power of 0.80 and about 160 to reach a power of 0.95. Since such an SDR is already very problematic in terms of the type 1 error rate for the Student's t -test (Bradley, 1978), needing such a large sample size to detect it is a serious hurdle. This issue becomes even worse for lower SDR , since an SDR as small as 1.1 already calls for the use of Welch's t -test (See table A3.1 to A3.9 in the additional file). Detecting such a small SDR calls for a huge sample size (a sample size of 160 provides a power rate of 0.16).

Since Welch's t -test has practically the same power as Student's t -test, even when $SDR = 1$, as explained using simulations later, we should seriously consider using Welch's t -test by default.

The problems in using a two-step procedure (first testing for equality of variances, then deciding upon which test to use) have already been discussed in the field of statistics (see e.g., Rasch, Kubinger, & Moder, 2011; Ruxton, 2006; Wilcox, Granger, & Clark, 2013; Zimmerman, 2004), but these insights have not changed the current practices in psychology, as of yet. More importantly, researchers do not even seem to take the assumptions of Student's t -test into consideration before performing the test, or at least rarely discuss assumption checks.

We surveyed statistical tests reported in the journal *SPPS* (*Social Psychological and Personality Science*) between April 2015 and April 2016. From the total of 282 studies, 97 used a t -test (34.4%), and the homogeneity of variance was explicitly discussed in only 2 of them. Moreover, based on the reported degrees of freedom in the results section, it seems that Student's t -test is used most often and that alternatives are considerably less popular. For 7 studies,

there were decimals in the values of the degrees of freedom, which suggests Welch's t -test might have been used, although the use of Welch's t -test might be higher but not identifiable because some statisticians recommend rounding the degrees of freedom to round numbers.

To explain this lack of attention to assumption checks, some authors have argued that researchers might have a lack of knowledge (or a misunderstanding) of the parametric assumptions and consequences of their violations or that they might not know how to check assumptions or what to do when assumptions are violated (Hoekstra, Kiers, & Johnson, 2012).⁶ Finally, many researchers don't even know there are options other than the Student's t -test for comparing two groups (Erceg-Hurn & Miroseovich, 2008). How problematic this is depends on how plausible the assumption of equal variances is in psychological research. We will discuss circumstances under which the equality of variances assumption is especially improbable and provide real-life examples where the assumption of equal variances is violated.

Homogeneity of Variance Assumptions

The homogeneity of variances assumption is rarely true in real life and cannot be taken for granted when performing a statistical test (Erceg-Hurn & Miroseovich, 2008; Zumbo & Coulombe, 1997). Many authors have examined real data and noted that SDR is often different from the 1:1 ratio (see, e.g., Grissom, 2000; Erceg-Hurn & Miroseovich, 2008). This shows that the presence of unequal variances is a realistic assumption in psychological research.⁷ We will discuss three different origins of unequal standard deviations across two groups of observations.

A first reason for unequal variances across groups is that psychologists often use *measured variables* (such as age, gender, educational level, ethnic origin, depression level, etc.) instead of random assignment to condition. In their review of comparing psychological findings from all fields of the behavioral sciences across cultures, Henrich, Heine, and Norenzayan (2010) suggest that parameters vary largely from one population to another. In other words, variance is not systematically the same in every pre-existing group. For example, Feingold (1992) has