

# Data Analyses When Sample Sizes Are Small: Modern Advances for Dealing With Outliers, Skewed Distributions, and Heteroscedasticity

Rand Wilcox,<sup>1</sup> Travis J. Peterson,<sup>2</sup> and Jill L. McNitt-Gray<sup>1</sup>

<sup>1</sup>University of Southern California; <sup>2</sup>California Lutheran University

The paper reviews advances and insights relevant to comparing groups when the sample sizes are small. There are conditions under which conventional, routinely used techniques are satisfactory. But major insights regarding outliers, skewed distributions, and unequal variances (heteroscedasticity) make it clear that under general conditions they provide poor control over the type I error probability and can have relatively poor power. In practical terms, important differences among groups can be missed and poorly characterized. Many new and improved methods have been derived that are aimed at dealing with the shortcomings of classic methods. To provide a conceptual basis for understanding the practical importance of modern methods, the paper reviews some modern insights related to why methods based on means can perform poorly. Then some strategies for dealing with nonnormal distributions and unequal variances are described. For brevity, the focus is on comparing 2 independent groups or 2 dependent groups based on the usual difference scores. The paper concludes with comments on issues to consider when choosing from among the methods reviewed in the paper.

**Keywords:** robust methods, distribution-free techniques, nonparametric methods, nonnormality, ANOVA

Classic, routinely taught and used statistical methods for comparing groups are based on means. From basic principles, there are 2 fundamental goals: (1) to control the probability of a type I error, meaning rejecting the null hypothesis when it is true, and (2) to have relatively high power, meaning the probability of detecting a true difference. A positive feature of standard methods for comparing means is that they control the type I error probability reasonably well when observations are sampled from populations of individuals that have identical means, variances, skewness, and so forth. That is, sampling is from identical distributions. However, recent insights indicate that under general conditions when group distributions differ these 2 fundamental goals are not achieved,<sup>1-7</sup> resulting in poor power and a highly misleading summary of the data. If distributions differ, even with relatively large sample sizes ( $n > 200$ ), classic methods can yield inaccurate confidence intervals and can have relatively poor power.<sup>6,7</sup> In practical terms, if Student's  $t$  test fails to reject, one possibility is that there is little or no difference between the groups, or it might be because power is poor due to nonnormality (see [Supplementary Material](#) [available online]; Welch's 2-sample  $t$  test and Yuen's method with trimmed means). One strategy for dealing with nonnormal distributions is to use classic rank-based methods. They test the hypothesis that groups have identical distributions. But in terms of comparing means or medians, they are unsatisfactory under general conditions.<sup>8,9</sup> This remains the case when using more modern rank-based methods described by Brunner et al.<sup>10</sup> Another strategy is to rely on the

generalized linear model. But when using means, this does not address the theoretical issues described by Huber and Ronchetti,<sup>2</sup> Hampel et al.,<sup>11</sup> and Staudte and Sheather.<sup>5</sup> Other concerns are described by Keselman et al.<sup>12</sup>

Recent advances in modern statistical methods provide robust hypothesis testing by dealing with outliers and skewed distributions in ways that maintain power and control the probability of a type I error. No single method dominates in that each approach is sensitive to different features of the data and each tells something different about how groups compare. To provide a conceptual basis for understanding the practical importance of modern methods, this paper reviews some modern insights related to why methods based on means can perform poorly. Strategies for dealing with nonnormal distributions and unequal variances are also described. Modern advances for dealing with outliers, skewed distributions, limitations associated with the central limit theorem, and unequal variances (heteroscedasticity) in the context of small sample sizes ( $<10$  participants) will be discussed. For brevity, examples will be limited to the comparison of 2 independent groups or 2 dependent groups based on the usual difference scores. The paper concludes with comments on issues to consider when choosing between the methods reviewed in the paper.

## Analysis of Small Sample Sizes

Small sample sizes may be due to a very limited subject pool such as when studying elite athletes or data collection that might be difficult due to the time and expense that is required. For such situations, there is the practical issue of choosing an appropriate statistical technique. A related issue is interpreting the results in an appropriate manner. This entails having a precise understanding of what a statistical method reveals about the data and just as important, understanding what it does not reveal.

Space limitations preclude a detailed explanation regarding when and why modern methods offer a practical advantage.

Wilcox is with the Department of Psychology, University of Southern California, Los Angeles, CA, USA. Peterson is with the Department of Exercise Science, California Lutheran University, Thousand Oaks, CA, USA. McNitt-Gray is with the Department of Biomedical Engineering, University of Southern California, Los Angeles, CA, USA. McNitt-Gray is also with the Department of Biological Sciences, University of Southern California, Los Angeles, CA, USA. McNitt-Gray ([mcnitt@usc.edu](mailto:mcnitt@usc.edu)) is corresponding author.

The goal here is to provide some indications regarding their relative merits. The paper begins by briefly outlining the nature of 3 major advances related to outliers, skewed distributions, and heteroscedasticity (unequal variances). Figures are included to provide context. Some unsatisfactory strategies for salvaging classic techniques are also discussed. This is followed by a section on robust methods for testing hypotheses that deal effectively with small sample sizes in a manner that addresses the issues previously described and a description of some of the more basic methods that are now available for consideration (see [Supplementary Material](#) [available online] and corresponding references). Illustrative examples from sports biomechanics research are also provided.

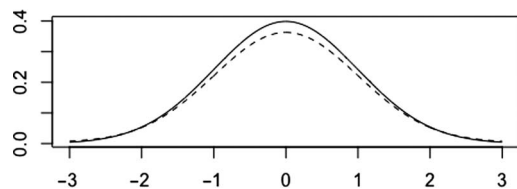
## Outliers

The first major advance has to do with outliers and what are generally known as heavy-tailed distributions. Here it is assumed that values are flagged as outliers based on a boxplot or the so-called median absolute deviation-median rule. Methods based on the mean and variance are known to be unsatisfactory.<sup>6</sup> Roughly, when sampling from a heavy-tailed distribution, outliers tend to occur more frequently versus sampling from a normal distribution. This can result in poor power because outliers can inflate the standard deviation. Even a small departure from normality can result in poor power.<sup>5-7,13</sup> This result follows almost immediately from Tukey's<sup>14</sup> discussion of a mixed normal distribution.

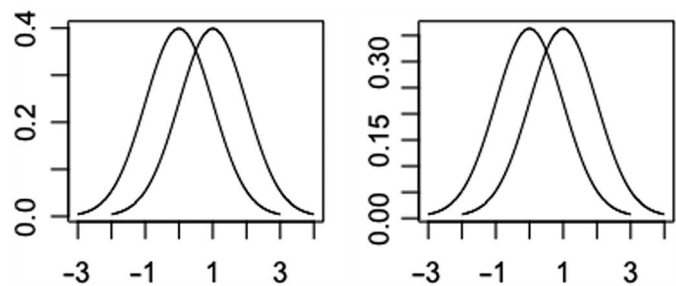
The population variance is sensitive to slight changes in the tails of a distribution. For example, Figure 1 shows a standard normal distribution and a mixed normal distribution, which has heavy tails, meaning that the tails are thicker than the normal distribution. Despite the obvious similarity, the standard normal distribution has variance of 1, but the variance of the mixed normal is 10.9. This comparison illustrates the basic principle that the population variance is very sensitive to slight changes in the tails of a distribution.

Even when a distribution appears to be approximately normal, power can be relatively low. Figure 2 illustrates this point. The left panel shows 2 normal distributions, both of which have variances equal to 1. With sample sizes  $n_1 = n_2 = 25$ , power is .96 when testing at the .05 level with Student's  $t$ . By contrast, the distributions in the right panel are mixed normals. Now power is only .28 and is considered relatively low.

Put another way, outliers are more likely to occur when dealing with heavy-tailed distributions such as the mixed normal. Outliers can inflate the standard deviation to the point that power is relatively low when using means.



**Figure 1** — The solid line is a standard normal distribution having variance 1. The dashed line is a mixed normal distribution, which is characterized as having heavy tails, meaning that the tails are thicker than a normal distribution. Despite the evident similarity, the mixed normal has variance 10.9. An implication is that even a slight departure from a normal distribution can result in relatively low power.



**Figure 2** — In the left panel, power is .96 based on Student's  $t$ ,  $\alpha = .05$  and sample sizes  $n_1 = n_2 = 25$ . But in the right panel, the distributions are not normal, they are mixed normal, and power is only .28.

## Skewed Distributions

The second advance has to do with skewed distributions, which are a much more serious concern, when using means, than is generally recognized. Suppose that for a nominal .05 type I error probability, Student's  $t$  is judged to be reasonably accurate if the actual type I error probability is between .03 and .08. When observations are randomly sampled from a skewed distribution that has relatively light tails (outliers are relatively rare), approximately 130 observations can be required. When dealing with a skewed distribution where outliers are relatively common, 300 observations can be required.<sup>7</sup>

Skewed distributions have implications regarding the 2-sample  $t$  test. When testing at the .05 level, if the distributions being compared have the same amount of skewness, then Student's  $t$  is satisfactory in the sense that the actual type I error probability will not exceed .05. But otherwise, the actual type I error probability can be substantially larger than .05. Concerns persist regardless of how large the sample sizes might be.<sup>15</sup> For a more detailed summary regarding why skewed distributions are a concern, see Wilcox and Rousselet.<sup>16</sup>

## Heteroscedasticity (Unequal Variances)

Conventional methods for comparing independent groups assume homoscedasticity (equal population variances). Heteroscedasticity (unequal population variances) can negatively impact both power and the ability to control the type I error probability. For example, when testing at the .05 level, the actual type I error probability can exceed .3 when using the analysis of variance  $F$  test.<sup>6</sup>

## Some Unsatisfactory Strategies for Salvaging Classic Techniques

Simple transformations are sometimes suggested for salvaging methods based on means. For example, taking logarithms sometimes results in a distribution that is approximately normal. But in general, simple transformations are relatively unsatisfactory. Typically, distributions remain skewed and the deleterious impact of outliers remains.<sup>17-19</sup>

Another strategy is to test assumptions. For example, when comparing 2 independent groups, test the hypothesis of equal variances and if a nonsignificant result is obtained, use Student's  $t$ . However, extant publications do not support this approach.<sup>20-24</sup> The basic problem is that tests of the hypothesis of equal variances do not have enough power to detect situations where the equal variance assumption should be abandoned.

## Dealing Effectively with Small Sample Sizes

There are robust methods for testing hypotheses that deal effectively with small sample sizes in a manner that addresses the issues previously described. Moreover, they are readily applied via the software R.<sup>25</sup> For the 1-sample case or when comparing dependent groups using difference scores, a few methodologies are mentioned here. The first method is based on the median, which assumes random sampling only (Supplementary Material [available online]).<sup>6,7,9</sup> A possible criticism is that in some situations the sample median trims too many observations. In this case, one might use a compromised amount of trimming. A 20% trimmed mean is often a good choice and has been studied extensively.<sup>6,7</sup> This means that after putting the observations in ascending order, the lower 20% are trimmed as well as the upper 20%. Power is nearly the same as Student's  $t$  when sampling from a normal distribution and power can be relatively high when outliers are likely to occur, as will be illustrated. It is stressed that technically sound methods for making inferences based on a trimmed mean are not obvious based on standard training. Simply trimming and using Student's  $t$  results in using an incorrect estimate of the standard error in that the observations left after trimming are dependent, which invalidates the derivation of Student's  $t$ .

In terms of type I errors, the best method uses a percentile bootstrap technique.<sup>6,7</sup> If the sample size is not too small, say greater than 20, the (nonbootstrap) method derived by Tukey and McLaughlin<sup>26</sup> is another way of dealing with trimmed means; it uses a correct estimate of the standard error. When dealing with 2 dependent groups, a third approach is to use improved versions of the sign test.<sup>6,27</sup>

As for comparing 2 independent groups, a method for comparing medians, via a percentile bootstrap method, performs very well in terms of controlling the type I error probability (see Supplementary Material [available online]).<sup>28,29</sup> Currently, it is the only known method that effectively deals with tied (duplicated) values.<sup>7</sup> If the likelihood of outliers is relatively low, using a 20% trimmed mean might result in better power. In addition, both theoretical results and simulation studies indicate that problems with controlling the type I error probability, when dealing with skewed distributions, diminish as the amount of trimming increases.<sup>30</sup>

To provide some perspective on power, consider again the 2 normal distributions in the left panel of Figure 2. As previously noted, power is .96 when comparing means at the .05 level and when both sample sizes equal to 25. Using instead a 20% trimmed mean, power is .91, and it is .78 when comparing the medians. Now consider the right panel and recall that power is .28 when comparing means. By contrast, when using a 20% trimmed mean or median, power is .75 and .70, respectively. This illustrates the general principle that when testing hypotheses based on means, power can be relatively low when outliers tend to occur.

Let  $P$  be the probability that a randomly sampled observation from the first group is less than a randomly sampled observation from the second. The classic Wilcoxon–Mann–Whitney (WMW) test is based on a direct estimate of  $P$ . But under general conditions it is not a satisfactory test of the hypothesis that  $P = .5$  and it does not provide a satisfactory confidence interval for  $P$ .<sup>8</sup> The reason is that it uses the wrong standard error when distributions differ. In addition, under general conditions, the WMW does not test the hypothesis of equal medians as is sometimes claimed.<sup>31</sup> There are several methods for testing the hypothesis that  $P = .5$  that continue to perform well when distributions differ in shape. Based on results

in Neuhaus et al,<sup>32</sup> the method derived by Cliff<sup>8</sup> performs particularly well when the sample sizes are small. Unlike the WMW test, Cliff's method is designed to perform well when there is heteroscedasticity.

## Choosing a Method: Different Methods Provide Different Information

Several methods have been mentioned regarding how to compare 2 independent groups when the sample size is relatively small. This raises a natural question: which one is best? Because different methods are sensitive to different features of the data, the only certainty is that no single method dominates in terms of power. In addition, different methods provide different perspectives. If, for example, distributions are skewed, comparing the medians is not the same as comparing 20% trimmed means or testing the hypothesis that  $P = .5$ . In some situations, multiple perspectives might provide a deeper and more nuanced understanding of how groups compare.

As for comparing dependent groups, again different methods provide different perspectives. Modern variations of the sign test might seem relatively uninteresting because in some sense they ignore information that is utilized by methods based on a measure of central tendency. But situations are encountered where the sign test makes a practical difference as will be illustrated.

The first illustration deals with the comparison of vertical velocity generated by volleyball players ( $n = 12$ ) when performing a block maneuver initiated with and without horizontal momentum. Five trials were performed under each condition. Here, to control the probability of one or more type I errors, an improvement on the Bonferroni method is used that was derived by Hochberg.<sup>33</sup> Testing the hypothesis that the difference scores have a mean of 0, using the usual paired  $t$  test, no significant differences are found. Using instead a 20% trimmed mean and the Tukey–McLaughlin method, 2 significant differences are found. Again, using a 20% trimmed mean but with the percentile bootstrap method, a significant difference is found on all 5 occasions. Using the median, a significant difference is found on 4 of the 5 occasions. The sign test returns 2 significant differences. The median absolute deviation–median rule<sup>7</sup> indicates that among the difference scores, the number of outliers for the 5 occasions is 1, 0, 2, 4, and 0. The very presence of outliers does not necessarily mean more power when using a trimmed mean or median, the only point is that the choice of method can make a practical difference. Russell et al<sup>34</sup> provide another example dealing with the sign test in a clinical context.

The second illustration stems from a study comparing the kinematics and reaction forces generated during a golf swing under 2 conditions.<sup>35</sup> Two types of shots were attempted on 4 different occasions resulting in 4 sets of difference scores. Again, the probability of one or more type I errors was controlled using Hochberg's method. On the third occasion, the median absolute deviation–median rule finds a kinematic variable with 6 outliers, 4 of which occur among the higher measures. The point is that with a 20% trimmed mean, only 3 of these 4 outliers are trimmed. Moreover, if the goal is to have the probability of one or more type I errors equal to .05, the hypothesis that the trimmed mean is 0 is not rejected. However, the sign test, as well as the hypothesis of a median difference equal to 0 is rejected as well. That is, despite the extreme amount of trimming used by the median, significant results can occur in contrast to results when using less trimming.

A similar result was obtained when comparing the peak force measures between the 2 shot conditions. For completeness, the paired  $t$  test also rejects in this particular case, despite the outliers. But why it



Currently, we have the technology to get a deeper and more accurate understanding of our data. In particular, substantially improved statistical methods are available for dealing with small sample sizes. Moreover, free R software for applying modern methods is available.<sup>6,7,36,37</sup> Three of the books cited are aimed at an introductory statistics course and include both classic and modern robust methods. For SAS and SPSS users interested in using R, the book by Muenchen<sup>38</sup> might help. Although space limitations make it impossible to foster a deep appreciation of the many advances that have occurred, our hope is that this brief review will stimulate the use of modern techniques.

## References

- (Ahead of Print)**