

[Page One](#)[Campus  
Computing  
News](#)[New Software  
Available](#)[Important  
Summer  
Reading](#)[Free Virus  
Protection for  
Home PCs](#)[Today's Cartoon](#)[RSS Matters](#)[SAS Corner](#)[The Network  
Connection](#)[Link of the  
Month](#)[WWW@UNT.EDU](#)[Short Courses](#)[IRC News](#)[Staff Activities](#)[Subscribe to  
Benchmarks  
Online](#)

# Research and Statistical Support

## University of North Texas

### RSS Matters

The previous article in this series can be found in a previous [issue](#) of Benchmarks Online: *Controlling the False Discovery Rate in Multiple Hypothesis Testing* - Ed.

## Using Robust Mean and Robust Variance Estimates to Calculate Robust Effect Size

By [Dr. Rich Herrington](#), Research and Statistical Support Consultant

This month we demonstrate the calculation of robust effect sizes. The GNU S language, "R" is used to implement this procedure. R is a statistical programming environment that is a clone of the S and S-Plus language developed at Lucent Technologies. In the following document we illustrate the use of a GNU Web interface to the R engine on the "rss" server, <http://rss.acs.unt.edu/cgi-bin/R/Rprog>. This GNU Web interface is a derivative of the "Rcgi" Perl scripts available for download from the CRAN website, <http://www.cran.r-project.org> (the main "R" website). Scripts can be submitted interactively, edited, and be re-submitted with changed parameters by selecting the hypertext link buttons that appear below the figures. For example, clicking the "Run Program" button below creates a vector of 100 random normal deviates; displays the results; sorts and displays the results; then creates a histogram and a density plot of the random numbers. To view any text output, scroll to the bottom of the browser window. To view any graphical output, select the "Display Graphic" link. The script can be edited and resubmitted by changing the script in the form window and then selecting "Run the R Program". Selecting the browser "back page" button will return the reader to this document.

### Introduction - Calculating Power and Effect Size

Power analysis involves the relationships between four variables involved in statistical inference: sample size (N), a significance criterion ( $\alpha$ ), the population effect size ( $d_{\text{cohen}}$ ), and statistical power. For any statistical inference, these relationships are a function of the other three (Cohen, 1988). For research planning, it is most useful to determine the N necessary to have a specified power for a given  $\alpha$  and  $d_{\text{cohen}}$ . The statistical power of a test is the long term probability of rejecting  $H_0$  (null hypothesis) given a specified  $\alpha$  criterion and sample size N. When the effect size is not equal to zero,  $H_0$  is false, so a failure to reject  $H_0$  is a decision error on the part of the researcher. This is called a type II error ( $\beta$ ) and is related mathematically to power. The probability of rejecting the null if it needs to be rejected (power) is one minus the type II error ( $1-\beta$ ). Figure 1. below is a graphical representation of the relationship between the null distribution, the alternate distribution, and the critical scores under the null distribution. The area underneath the  $H_1$  distribution (the alternate distribution), past the critical score of the left tail of  $H_0$ , and past the critical score of the right tail of  $H_0$ , represents the power of the statistical test being performed (the shaded area).

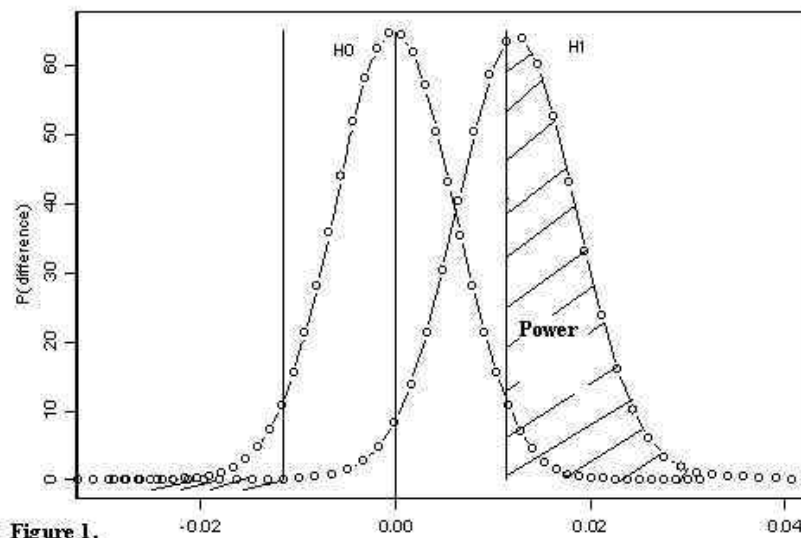


Figure 1.

Effect size is the degree to which  $H_0$  (null hypothesis) is false and is indexed by the discrepancy between the null hypothesis and the alternate hypothesis. Power analysis specifies a non-centrality parameter to quantify this discrepancy. The noncentrality parameter for the difference between means is:

$$\delta = \frac{|\hat{\mu}_1 - \hat{\mu}_2|}{\hat{\sigma}_{diff}}$$

where the difference between estimated population means is scaled in  $\hat{\sigma}_{diff}$  units (known as the estimated standard error of the difference between means):

$$\hat{\sigma}_{diff} = \sqrt{\hat{\sigma}_{pooled}^2 \cdot \left( \frac{n_1 \cdot n_2}{n_1 + n_2} \right)}$$

where,

$$\hat{\sigma}_{pooled}^2 = \frac{(n_1 - 1) \cdot \hat{\sigma}_1^2 + (n_2 - 1) \cdot \hat{\sigma}_2^2}{n_1 + n_2 - 2}$$

and  $\hat{\sigma}$  is the sample estimate of the population standard deviation. The denominator of the non-centrality parameter represents the estimated standard deviation of the sampling distribution for the null hypothesis for differences between means. Usually,  $\hat{\sigma}_{diff}$  is calculated on the basis of a formula that assumes normality in the population since the standard deviation of the null sampling distribution cannot be calculated directly on the basis of the observed data without normality assumptions. For robust measures of location (i.e. M-estimate), the numerator would be the difference between two M-estimates, and the denominator would represent the standard deviation of the null hypothesis re-sampling distribution for the difference between M-estimates. For robust estimates (as well as the sample mean), the standard error can be estimated directly by calculating the standard deviation of the bootstrap estimates of the differences between the robust estimates of location (see [September 2001 issue of Benchmarks](#)). An alternative effect size for group differences has been advocated by Cohen (1988). Cohen's  $d_{cohen}$  measure is based on the pooled estimated population standard deviation:

$$d_{cohen} = \frac{|\hat{\mu}_1 - \hat{\mu}_2|}{\hat{\sigma}_{pooled}}$$

Cohen provides guidelines for interpreting the practical importance of an effect size based on  $d_{cohen}$  when no prior research is available to anchor  $d_{cohen}$  meaningfully. Cohen's rule of thumb for a small, medium and large effect size are based on a wide examination of the typical difference found in psychological data. A small effect size for  $d_{cohen}$  is .20; a medium effect size for  $d_{cohen}$  is .50; and a large effect size for  $d_{cohen}$  is .80 (Cohen, 1992).

Equating  $\delta$  and  $d_{cohen}$  using algebra, the expression for  $\delta$  is:

$$\delta = d_{\text{cohen}} \cdot \sqrt{\frac{n_1 \cdot n_2}{n_1 + n_2}}$$

It is noted that  $d_{\text{cohen}}$  is not a robust measure of effect size. The pooled sample standard deviation, which is used to estimate the population standard deviation ( $\hat{\sigma}$ ) will be inflated in the presence of outliers thereby biasing the effect size measure. Furthermore,  $d_{\text{cohen}}$  assumes a normal distribution in the calculation of power estimates.

## Measures of Robust Effect Size

Several problems exist with the  $d_{\text{cohen}}$  measure of effect size. The assumption of equal variances in the population is often dealt with by substituting a pooled variance estimate for  $\sigma$ . With data that appear to have unequal variances, questions arise about how to interpret  $H_0$ . Another criticism of  $d_{\text{cohen}}$  is that both the location and scale (mean difference and sample standard deviation) of the sample are non-resistant measures. One strategy would be to replace the means and standard deviation with more resistant measures of location and scale. For example, one variation might be a difference of medians divided by MAD (median absolute deviation):

$$d_{\text{median}} = \frac{\mu_{M1} - \mu_{M2}}{MAD}$$

where  $MAD = MED[|X_1 - M|, \dots, |X_n - M|]$  and  $M$  is the median of the scores in the control group. This effect size estimator does not seem like a good candidate since both the median and MAD are both known to be inefficient for Normal distributions compared to the mean and standard deviation.

## Robust Effect Size based on M-estimators

Lax (1985) examined the performance of 17 different estimators of scale with heavy tailed distributions. Lax examined the performance of these scale estimators with the Normal distribution; a distribution with Cauchy tails (large kurtosis relative to the Normal The Slash dist.); and a mixture distribution of  $N(0,1)$  and  $N(0,100)$  for samples of size 20. The mixture distribution had 19 points sampled from  $N(0,1)$  and 1 point sampled from  $N(0,100)$  (One-Wild dist.). Lax combined the efficiencies (see [July 2001 issue of Benchmarks](#)) of the estimators for the three distributions into what was defined as triefficiency. The biweight midvariance (with  $c=9$ ) estimator performed best, with favorable efficiencies across all scenarios: Normal (86.7%), One-Wild (85.8), and Slash (86.1). Following Wilcox (1997) the biweight midvariance can be calculated as follows. Setting (with  $c=9$ ,  $M$ =sample median):

$$Y_i = \frac{X_i - M}{c \cdot MAD} \quad \text{and} \quad a_i = \begin{cases} 1, & \text{if } |Y_i| < 1 \\ 0, & \text{if } |Y_i| \geq 1 \end{cases}$$

the following is calculated:

$$\hat{s}_{bi} = \frac{\sqrt{n} \cdot \sqrt{\sum a_i (X_i - M)^2 (1 - Y_i^2)^4}}{\left| \sum a_i (1 - Y_i^2) (1 - 5Y_i^2) \right|}$$

The square of  $\hat{s}_{bi}$  is called the biweight midvariance. It appears to have a breakdown point of approximately .5 (Hoaglin, Moesteller, & Tukey, 1983). Based on this robust variance, the following robust effect size can be calculated:

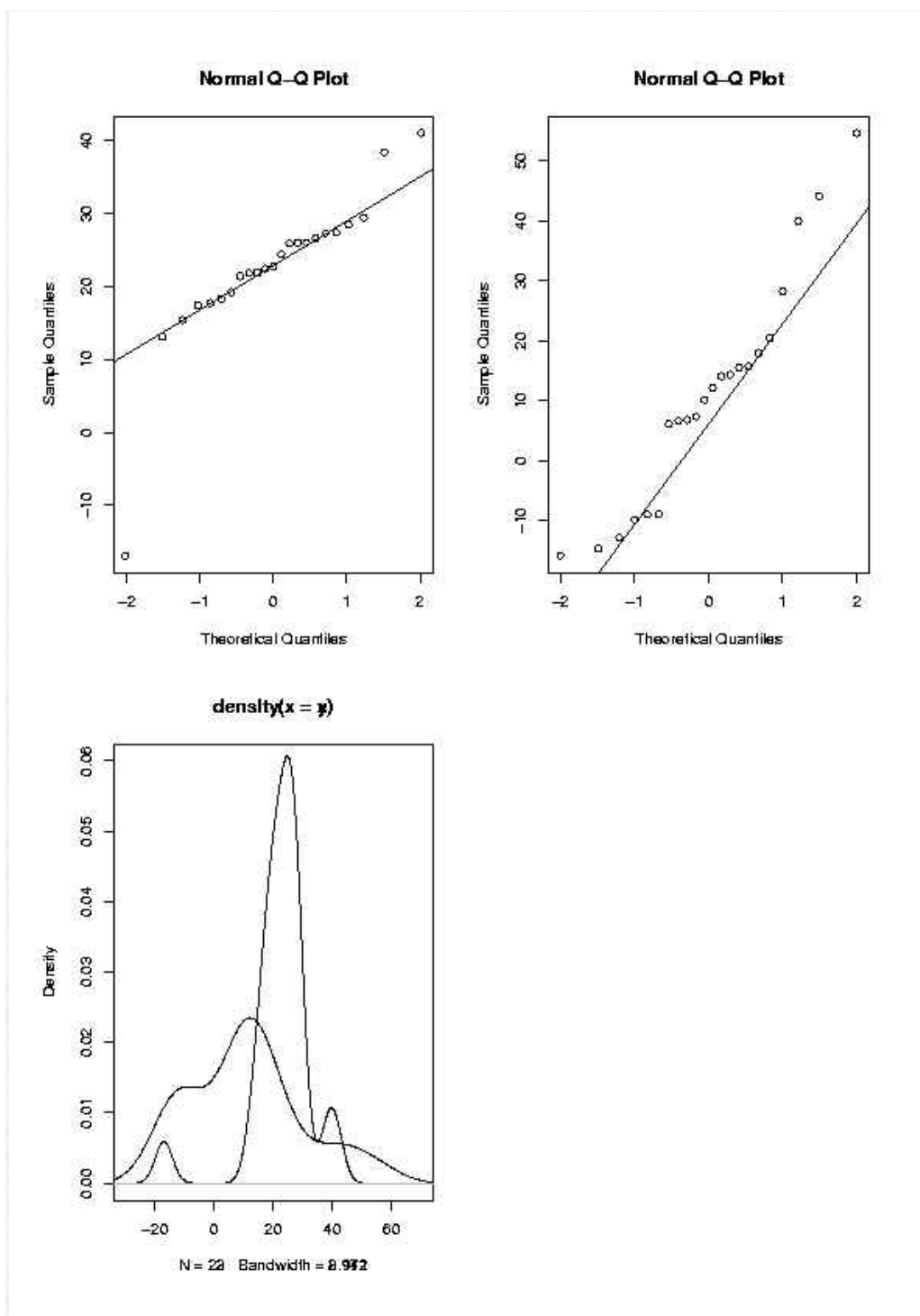
$$d_{\text{robust}} = \frac{\hat{\mu}_{\text{Mest}1} - \hat{\mu}_{\text{Mest}2}}{\hat{s}_{bi1}}$$

where,  $\hat{\mu}_{\text{Mest}1}$  is the robust M-estimator for group 1 (using Huber objective function, with  $k=1.28$  for both groups),  $\hat{\mu}_{\text{Mest}2}$  is the robust M-estimator for group 2, and  $\hat{s}_{bi1}$  is the square root of the biweight midvariance for group 1 (control group). The robust effects size  $d_{\text{robust}}$  does not assume equal variances among groups since only the robust variance for the control group is used (alternatively, a pooled estimated of both the control and experimental group biweight midvariances could be used, assuming equal variances).

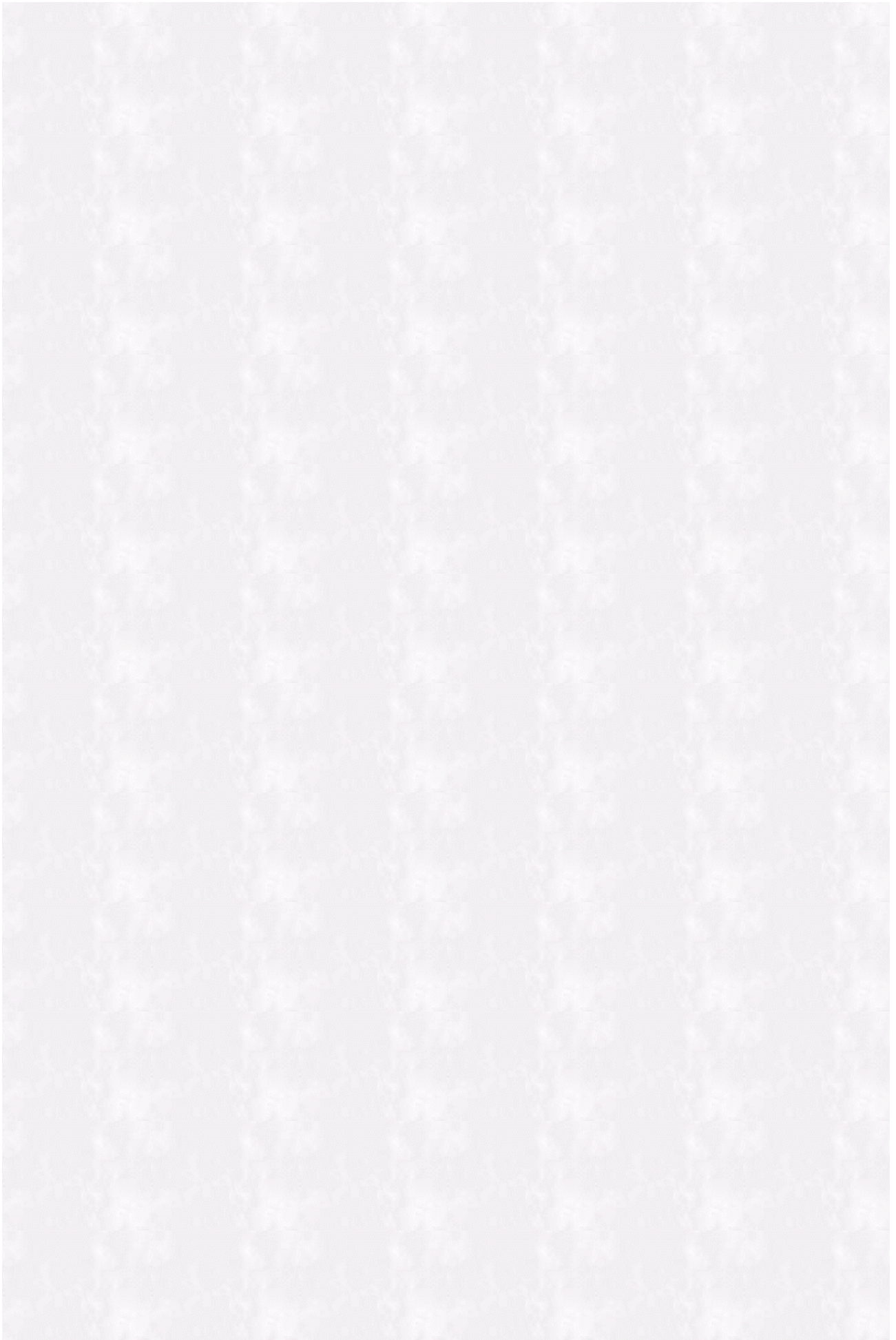
## An Example Using GNU-S ("R")

Doksum & Sievers (1976) report data on a study designed to assess the effects of ozone on weight gain in rats. The experimental group consisted of 22 seventy-day old rats kept in an ozone environment for 7 days (group y). The control group consisted of 23 rats of the same age (group x), and were kept in an ozone-free environment. Weight gain is measured in grams. The following R code produces quantile-quantile plots and non-parametric density plots of the two groups of data:

Resulting qqnorm plots and density plots from R code above:



Both groups appear to have right and left tails which are "heavy". It appears as if the classical mean difference between the groups will be underestimated (smaller). The R code below estimates both the classical means, and robust means; classical estimated pooled standard deviation, and estimated pooled robust root biweight midvariance.



## Results and Conclusion

```

> ## Classical Mean
> mean(x)
[1] 22.40435
> mean(y)
[1] 11.00909
>
> ## Robust Mean (M-estimator)
> mest(x)
[1] 23.14211
> mest(y)
[1] 9.687215
>
> ## Classical Pooled Variance
> std.dev.pooled.x.y<-std.dev.pooled(var(x)^.5, var(y)^.5, length(x), length(y))
> std.dev.pooled.x.y
[1] 15.36189
>
> ## Pooled Robust variance
> robust.std.dev.pooled.x.y<-std.dev.pooled(bivar(x)^.5, bivar(y)^.5, length(x), length(y))
> robust.std.dev.pooled.x.y
[1] 14.3583
>

```

The resulting M-estimators suggest that the population control group mean is downwardly biased (23.24 - robust; 22.40 - classical) and the experimental population group mean is biased upwardly (9.69 - robust; 11.01 - classical). Additionally, the robust pooled scale estimate is smaller than the classical pooled scale estimate (14.36 - robust; 15.36 - classical). Using these estimates to calculate Cohen's d measure indicates that the effect size is downwardly biased. Cohen's d based on classical estimators suggests a medium to large effect size (.74), whereas Cohen's d based on robust estimators suggests a very large effect size (.94). In terms of sample size planning for future experiments, the robust Cohen's d would suggest that a much smaller sample size would be needed to achieve the same power for a smaller effect size using non-robust estimators of location and scale - a considerable savings in terms of data that needs to be collected.

```
> ## Cohens d (non robust)
> cohen.d<- (mean(x)-mean(y))/std.dev.pooled.x.y
> cohen.d
[1] 0.7417875
>
> ## Robust Cohen's d
> robust.d<- (mest(x)-mest(y))/robust.std.dev.pooled.x.y
> robust.d
[1] 0.937081
>
```

## References

Cohen J. (1992) A power primer. *Psychological Bulletin*, 112, 155-159.

Cohen, J. 1988. *Statistical Power Analysis for the Behavioral Sciences*, 2nd Edition. Lawrence Erlbaum Associates, Inc., Hillsdale, New Jersey.

Doksum, K.A. & Sievers, G.L. (1976). Plotting with confidence: graphical comparisons of two populations. *Biometrika* 63, 421-434.

Hoaglin, D. C., Mosteller, F., & Tukey, J.W. (1983). *Understanding robust and exploratory data analysis*. New York: Wiley.

Lax, D.A. (1985). Robust estimators of scale: finite sample performance in long-tailed symmetric distributions. *Journal of the American Statistical Association*, 80(391), 736-741.

Wilcox, Rand R. (1997). *Introduction to Robust Estimation and Hypothesis Testing*. Academic Press, New York.