

StatNews #17

Bootstrap Estimation of Variability

February 1997

Updated May 2012

The availability of more powerful computers has enabled us to use statistical techniques that were not feasible in the past. One such category of techniques is resampling, which uses an empirical process to obtain estimates and their precision. This newsletter will introduce bootstrapping, which is probably the most common of these resampling techniques.

The usual approach to learning about a population is to draw a sample from the population, and then calculate a statistic, for example the mean, from the sample that provides information about the population. If we were to draw repeated samples from the population and calculate the statistic for each one, we would get a different answer every time. How varied these answers are across the different samples would give us an idea of how much confidence to place in the calculated statistic. This variation is conveyed by the standard error of the statistic.

In practice, it is often impossible to draw repeated samples from the population in order to find the standard error of an estimator. For basic statistics such as the mean and regression parameters, the standard error can be readily calculated using the standard deviation of the sample. Other standard errors may be calculated from a single sample if the sample size is quite large, or if we assume that the distribution of the population has certain characteristics such as normality.

For many statistics, however, there is either no algebraic solution to obtain a standard error, or the distributional assumptions necessary to obtain one are not tenable and, if assumed, would yield very distorted results. Several different resampling techniques have been developed that enable us to obtain a good measure of the precision of a statistic.

The technique that is probably the best known is the bootstrap. The idea behind this technique is to treat the sample as if it were the whole population. This idea rests on the assumption that the sample captures the essential aspects of the population. We can then take repeated samples with replacement from the original sample. The statistic of interest can then be calculated from each of these samples. The average of the statistic over all of these samples is known as the bootstrap estimator. Now that we have repeated estimates of the same statistic, we can calculate its standard error, known as the bootstrap standard error. The accuracy of the bootstrap estimator increases with the number of samples drawn.

One example where the bootstrap has proven useful is for calculating the precision of a mean after removing some of the extreme values in the sample. Since this is not exactly a mean, the usual standard error formula would not be accurate. Bootstrapping can also be used, for example, to obtain precision estimates for structural equation models or for regression trees.

It is usually rather easy to implement the bootstrap algorithm. Many statistical packages such as Stata, SAS, R and SPSS offer bootstrap estimators for a variety of statistics. If you would like assistance in obtaining a bootstrap estimator and its associated standard error, do not hesitate to contact us.

References:

Efron, B. and Tibshirani, R. (1991), Statistical Data Analysis in the Computer Age, Science, vol. 253, p.390-395.

Efron, B. and Tibshirani, R. (1993), An Introduction to the Bootstrap, New York: Chapman & Hall.

Author: Françoise Vermeylen (fmv1@cornell.edu)

(This newsletter was distributed by the Cornell Statistical Consulting Unit. Please forward it to any interested colleagues, students, and research staff. Anyone not receiving this newsletter who would like to be added to the mailing list for future newsletters should contact us at cscu@cornell.edu. Information about the Cornell Statistical Consulting Unit and copies of previous newsletters can be obtained at <http://www.cscu.cornell.edu>).