

*Department Copy*

# **THE BOOTSTRAP METHOD FOR ASSESSING STATISTICAL ACCURACY**

*Bradley Efron*

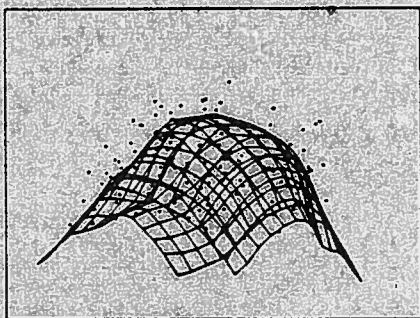
*and*

*Robert Tibshirani*

**Technical Report No. 19**

**October 1985**

**Laboratory for  
Computational  
Statistics**



**Department of Statistics  
Stanford University**

# THE BOOTSTRAP METHOD FOR ASSESSING STATISTICAL ACCURACY

by

B. Efron and R. Tibshirani

L.C.S. Technical Report No. 19

October 1985

Keywords: Bootstrap method, estimated standard errors, approximate confidence intervals, non-parametric methods.

## Abstract

This is an invited review of bootstrap methods. It begins with an exposition of the bootstrap estimate of standard error for one-sample situations. Several examples, some involving quite complicated statistical procedures, are given. The bootstrap is then extended to other measures of statistical accuracy, like bias and prediction error, and to complicated data structures such as time series, censored data, and regression models. Several more examples are presented illustrating these ideas. The last third of the paper deals mainly with bootstrap confidence intervals. The paper ends with a FORTRAN program for bootstrap standard errors.

This work was supported by an Office of Naval Research contract N00014-83-K-0472 and Public Health Service Grant 5 RO1 GM21215-10.

# The Bootstrap Method for Assessing Statistical Accuracy

B. Efron and R. Tibshirani

Stanford University

## 1. Introduction.

A typical problem in applied statistics is the estimation of an unknown parameter  $\theta$ . The two main questions asked are (1) what estimator  $\hat{\theta}$  should be used? And (2) having chosen to use a particular  $\hat{\theta}$ , how accurate is it as an estimator of  $\theta$ ? The bootstrap is a general methodology for answering the second question. It is a computer-based method, which substitutes considerable amounts of computation in place of theoretical analysis. As we shall see, the bootstrap can routinely answer questions which are far too complicated for traditional statistical analysis. Even for relatively simple problems computer-intensive methods like the bootstrap are an increasingly good data-analytic bargain in an era of exponentially declining computational costs.

This paper describes the basis of the bootstrap theory, which is very simple, gives several examples of its use, and ends with a bootstrap computer program, also very simple. Related ideas like the jackknife, the delta method, and Fisher's information bound are also discussed. Most of the proofs and technical details are omitted. These can be found in the references given, particularly Efron (1982). Some of the discussion here is abridged from Efron and Gong (1983), and also from Efron (1984b).

Before beginning the main exposition, we will describe how the bootstrap works in terms of a problem where it is not needed, assessing the accuracy of the sample mean. Suppose that our data consists of a random sample from an unknown probability distribution  $F$  on the real line,

$$X_1, X_2, \dots, X_n \sim F. \quad (1.1)$$

Having observed  $X_1 = x_1, X_2 = x_2, \dots, X_n = x_n$ , we compute the sample mean  $\bar{x} = \sum_1^n x_n/n$ , and wonder how accurate it is as an estimate of the true mean  $\theta = E_F\{X\}$ .

If the second central moment of  $F$  is  $\mu_2(F) \equiv E_F X^2 - (E_F X)^2$ , then the standard error  $\sigma(F; n, \bar{x})$ , that is the standard deviation of  $\bar{x}$  for a sample of size  $n$  from distribution  $F$ , is

$$\sigma(F) = [\mu_2(F)/n]^{1/2}. \quad (1.2)$$

(The shortened notation  $\sigma(F) \equiv \sigma(F; n, \bar{x})$  is allowable because the sample size  $n$  and statistic of interest  $\bar{x}$  are known, only  $F$  being unknown.) This is the traditional measure of  $\bar{x}$ 's accuracy. Unfortunately we can't actually use (1.2) to assess the accuracy of  $\bar{x}$ , since we don't know  $\mu_2(F)$ , but we can use the *estimated standard error*

$$\hat{\sigma} = [\hat{\mu}_2/n]^{1/2}, \quad (1.3)$$

where  $\hat{\mu}_2 = \sum_{i=1}^n (x_i - \bar{x})^2 / (n - 1)$ , the unbiased estimate of  $\mu_2(F)$ .

There is a more obvious way to estimate  $\sigma(F)$ . Let  $\hat{F}$  indicate the empirical probability distribution,

$$\hat{F} : \text{probability mass } 1/n \text{ on } x_1, x_2, \dots, x_n. \quad (1.4)$$

Then we can simply replace  $F$  by  $\hat{F}$  in (1.2), obtaining

$$\hat{\sigma} \equiv \sigma(\hat{F}) = [\mu_2(\hat{F})/n]^{1/2}, \quad (1.5)$$

as the estimated standard error for  $\bar{x}$ . This is the *bootstrap estimate*. The reason for the name "bootstrap" will be apparent in Section 2, when we evaluate  $\sigma(\hat{F})$  for statistics more complicated than  $\bar{x}$ . Since

$$\hat{\mu}_2 \equiv \mu_2(\hat{F}) = \sum_{i=1}^n \frac{(x_i - \bar{x})^2}{n}, \quad (1.6)$$

$\hat{\sigma}$  is not quite the same as  $\hat{\sigma}$ , but the difference is too small to be important in most applications.

Of course we don't really need an alternative formula to (1.3) in this case. The trouble begins when we want a standard error for estimators more complicated than  $\bar{x}$ , for example a median or a correlation or a slope coefficient from a robust regression. In most cases there is no equivalent to formula (1.2), which expresses the standard error  $\sigma(F)$  as a simple function of the sampling distribution  $F$ . As a result, formulas like (1.3) do not exist for most statistics.

This is where the computer comes in. It turns out that we can always numerically evaluate the bootstrap estimate  $\hat{\sigma} = \sigma(\hat{F})$ , even without knowing a simple expression for  $\sigma(F)$ . The

evaluation of  $\hat{\sigma}$  is a straightforward Monte Carlo exercise, described in the next section.

Standard errors are crude but useful measures of statistical accuracy. They are frequently used to give approximate confidence intervals for an unknown parameter  $\theta$ ,

$$\theta \in \hat{\theta} \pm \hat{\sigma} z^{(\alpha)}, \quad (1.7)$$

where  $z^{(\alpha)}$  is the  $100 \cdot \alpha$  percentile point of a standard normal variate, e.g.  $z^{(.95)} = 1.645$ . Interval (1.7) is sometimes good, and sometimes not so good. Sections 7 and 8 discuss a more sophisticated use of the bootstrap, which gives better approximate confidence intervals than (1.7).

The standard interval (1.7) is based on taking literally the large-sample normal approximation  $(\hat{\theta} - \theta)/\hat{\sigma} \sim N(0, 1)$ . Applied statisticians use a variety of tricks to improve this approximation. For instance if  $\theta$  is the correlation coefficient, and  $\hat{\theta}$  the sample correlation, then the transformation  $\phi = \tanh^{-1}(\theta)$ ,  $\hat{\phi} = \tanh^{-1}(\hat{\theta})$  greatly improves the normal approximation, at least in those cases where the underlying sampling distribution is bivariate normal. The correct tactic then is to transform, compute the interval (1.7) for  $\phi$ , and transform this interval back to the  $\theta$  scale.

We will see that bootstrap confidence intervals can automatically incorporate tricks like this, without requiring the data analyst to produce special techniques, like the  $\tanh^{-1}$  transformation, for each new situation. An important theme of what follows is the substitution of raw computing power for theoretical analysis. This is not an argument against theory, of course, only against unnecessary theory. Most common statistical methods were developed in the 1920's and 1930's, when computation was slow and expensive. Now that computation is fast and cheap we can hope for and expect changes in statistical methodology. This paper discusses one such potential change, Efron (1979b) discusses several others.

## 2. The Bootstrap Estimate of Standard Error.

This section presents a more careful description of the bootstrap estimate of standard error. For now we will assume that the observed data  $y = (x_1, x_2, \dots, x_n)$  consists of independent and identically distributed (i.i.d.) observations  $X_1, X_2, \dots, X_n \stackrel{\text{iid}}{\sim} F$ , as in (1.1). Here

$F$  represents an unknown probability distribution on  $X$ , the common sample space of the observations. We have a statistic of interest, say  $\hat{\theta}(y)$ , to which we wish to assign an estimated standard error.

Figure 1 shows an example. The sample space  $X$  is  $\mathbb{R}^{2+}$ , the positive quadrant of the plane. We have observed  $n = 15$  bivariate data points, each corresponding to an American law school. Each point  $x_i$  consists of two summary statistics for the 1973 entering class at law school  $i$ ,

$$x_i = (\text{LSAT}_i, \text{GPA}_i); \quad (2.1)$$

$\text{LSAT}_i$  is the class' average score on a nationwide exam called "LSAT";  $\text{GPA}_i$  is the class' average undergraduate grades. The observed Pearson correlation coefficient for these 15 points is  $\hat{\theta} = .776$ . We wish to assign a standard error to this estimate.

Let  $\sigma(F)$  indicate the standard error of  $\hat{\theta}$ , as a function of the unknown sampling distribution  $F$ ,

$$\sigma(F) = [\text{Var}_F\{\hat{\theta}(y)\}]^{1/2}. \quad (2.2)$$

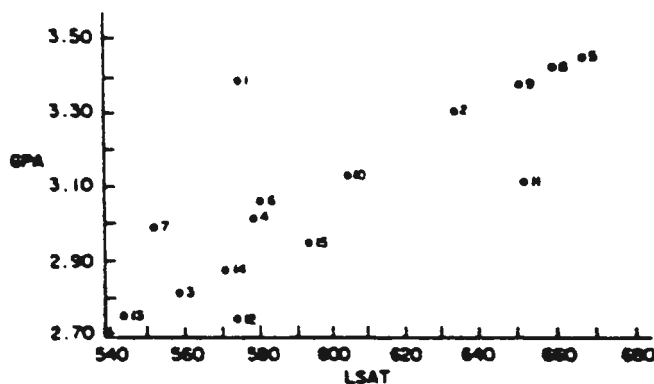
Of course  $\sigma(F)$  is also a function of the sample size  $n$  and the form of the statistic  $\hat{\theta}(y)$ , but since both of these are known they needn't be indicated in the notation. The bootstrap estimate of standard error is

$$\hat{\sigma} = \sigma(\hat{F}), \quad (2.3)$$

where  $\hat{F}$  is the empirical distribution (1.4), putting probability  $1/n$  on each observed data point  $x_i$ . In the law school example,  $\hat{F}$  is the distribution putting mass  $1/15$  on each point in Figure 1, and  $\hat{\sigma}$  is the standard deviation of the correlation coefficient for 15 i.i.d. points drawn from  $\hat{F}$ .

In most cases, including that of the correlation coefficient, there is no simple expression for the function  $\sigma(F)$  in (2.2). Nevertheless it is easy to numerically evaluate  $\hat{\sigma} = \sigma(\hat{F})$  by means of a Monte Carlo algorithm which depends on the following notation:  $y^* = (x_1^*, x_2^*, \dots, x_n^*)$  indicates  $n$  independent draws from  $\hat{F}$ , called a *bootstrap sample*. Because  $\hat{F}$  is the empirical distribution of the data, a bootstrap sample turns out to be the same as a random sample of size  $n$  drawn with replacement from the actual sample  $\{x_1, x_2, \dots, x_n\}$ .

The Monte Carlo algorithm proceeds in three steps.



**Figure 1.** The law school data (Efron 1979b). The data points, beginning with School No. 1, are (576, 3.39), (635, 3.30), (558, 2.81), (578, 3.03), (666, 3.44), (580, 3.07), (555, 3.00), (661, 3.43), (651, 3.36), (605, 3.13), (653, 3.12), (575, 2.74), (545, 2.76), (572, 2.88), (594, 2.96).

- (i) Using a random number generator, independently draw a large number of bootstrap samples, say  $y^*(1), y^*(2), \dots, y^*(B)$ .
- (ii) For each bootstrap sample  $y(b)$ , evaluate the statistic of interest, say  $\hat{\theta}^*(b) = \hat{\theta}(y^*(b))$ ,  $b = 1, 2, \dots, B$ .
- (iii) Calculate the sample standard deviation of the  $\hat{\theta}^*(b)$  values,

$$\hat{\sigma}_B = \left[ \frac{\sum_{b=1}^B \{\hat{\theta}^*(b) - \hat{\theta}^*(\cdot)\}^2}{B-1} \right]^{1/2} \quad \hat{\theta}^*(\cdot) = \frac{\sum_{b=1}^B \hat{\theta}^*(b)}{B} \quad (2.4)$$

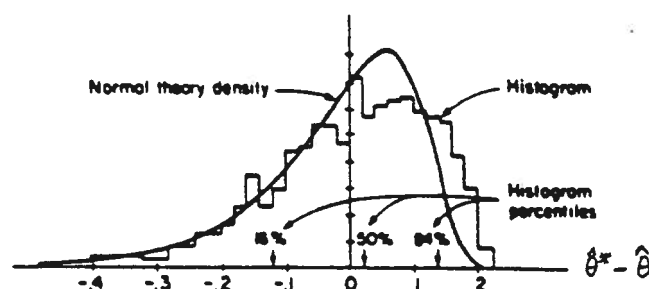
It is easy to see that as  $B \rightarrow \infty$ ,  $\hat{\sigma}_B$  will approach  $\hat{\sigma} = \sigma(\hat{F})$ , the bootstrap estimate of standard error. All we are doing is evaluating a standard deviation by Monte Carlo sampling. Later, in Section 9, we will discuss how large  $B$  need be taken. For most situations  $B$  in the range 50 to 200 is quite adequate. In what follows we will usually ignore the difference between  $\hat{\sigma}_B$  and  $\hat{\sigma}$ , calling both simply " $\hat{\sigma}$ ".

Figure 2 shows the histogram of  $B = 1000$  bootstrap replications of the correlation coefficient, from the law school data. For convenient reference the abscissa is plotted in terms of  $\hat{\theta}^* - \hat{\theta} = \hat{\theta}^* - .776$ . Formula (2.4) gives  $\hat{\sigma} = .127$  as the bootstrap estimate of standard error.

This can be compared with the usual normal theory estimate of standard error for  $\hat{\theta}$ ,

$$\hat{\sigma}_{\text{NORM}} = \frac{1 - \hat{\theta}^2}{(n - 3)^{1/2}} = .115, \quad (2.5)$$

Johnson and Kotz (1970), p. 229.



**Figure 2.** Histogram of  $B = 1000$  bootstrap replications of  $\hat{\theta}^*$  for the law school data. The normal theory density curve has a similar shape, but falls off more quickly at the upper tail.

There is another way to describe the bootstrap standard error:  $\hat{F}$  is the nonparametric maximum likelihood estimate (MLE) of the unknown distribution  $F$ , Kiefer and Wolfowitz (1956). This means that the bootstrap estimate  $\hat{\sigma} = \sigma(\hat{F})$  is the nonparametric MLE of  $\sigma(F)$ , the true standard error.

In fact there is nothing which says that the bootstrap must be carried out nonparametrically. Suppose for instance that in the law school example we believed the true sampling distribution  $F$  must be bivariate normal. Then we could estimate  $F$  with its *parametric* MLE  $\hat{F}_{\text{NORM}}$ , the bivariate normal distribution having the same mean vector and covariance matrix as the data. The bootstrap samples at step (i) of the algorithm could then be drawn from  $\hat{F}_{\text{NORM}}$  instead of  $\hat{F}$ , and steps (ii) and (iii) carried out as before.

The smooth curve in Figure 2 shows the results of carrying out this “normal theory bootstrap” on the law school data. Actually there is no need to do the bootstrap sampling in



this case, because of Fisher's formula for the sampling density of a correlation coefficient in the bivariate normal situation, see Chapter 32 of Johnson and Kotz (1970). This density is a close approximation to  $\hat{\sigma}_{\text{NORM}} = \sigma(\hat{F}_{\text{NORM}})$ , the parametric bootstrap estimate of standard error.

In considering the merits or demerits of the bootstrap, it is worth remembering that all of the usual formulas for estimating standard errors, like one over the square root of the observed Fisher information, are essentially bootstrap estimates carried out in a parametric framework. This point is carefully explained in Section 5 of Efron (1981b). The straightforward nonparametric algorithm (i)–(iii) has the virtues of avoiding all parametric assumptions, all approximations (such as those involved with the Fisher information expression for the standard error of an MLE), and in fact all analytic difficulties of any kind. The data analyst is free to obtain standard errors for enormously complicated estimators, subject only to the constraints of computer time. Sections 3 and 6 discuss some interesting applied problems which are far too complicated for standard analyses.

How well does the bootstrap work? Table 1 shows the answer in one situation. Here  $X$  is the real line,  $n = 15$ , and the statistic  $\hat{\theta}$  of interest is the 25% trimmed mean. If the true sampling distribution  $F$  is  $N(0, 1)$ , then the true standard error is  $\sigma(F) = .286$ . The bootstrap estimate  $\hat{\sigma}$  is nearly unbiased, averaging .287 in a large sampling experiment. The standard deviation of the bootstrap estimate  $\hat{\sigma}$  is itself .071 in this case, with coefficient of variation  $.071/.287 = .25$ . [Notice that there are two levels of Monte Carlo involved in Table 1: first drawing the actual samples  $y = (x_1, x_2, \dots, x_{15})$  from  $F$ , and then drawing bootstrap samples  $(x_1^*, x_2^*, \dots, x_{15}^*)$  with  $y$  held fixed. The bootstrap samples evaluate  $\hat{\sigma}$  for a fixed value of  $y$ . The standard deviation .071 refers to the variability of  $\hat{\sigma}$  due to the random choice of  $y$ .]

The jackknife is another common method of assigning nonparametric standard errors, discussed in Section 10. The jackknife estimate  $\hat{\sigma}_J$  is also nearly unbiased for  $\sigma(F)$ , but has higher coefficient of variation ( $CV$ ). The minimum possible  $CV$  for a scale-invariant estimate of  $\sigma(F)$ , assuming full knowledge of the parametric model, is shown in brackets. The nonparametric bootstrap is seen to be moderately efficient in both cases considered in Table 1.

|   | <i>F</i> Standard Normal |      |              | <i>F</i> Negative Exponential |      |              |
|---|--------------------------|------|--------------|-------------------------------|------|--------------|
|   | Ave                      | Sd   | Coeff<br>Var | Ave                           | Sd   | Coeff<br>Var |
| Bootstrap $\hat{\sigma}$ :<br>( $B = 200$ ) | .287                     | .071 | .25          | .242                          | .078 | .32          |
| Jackknife $\hat{\sigma}_J$                  | .280                     | .084 | .30          | .224                          | .085 | .38          |
| True :<br>[Minimum C.V.]                    | .286                     |      | [.19]        | .232                          |      | [.27]        |

**Table 1.** A sampling experiment comparing the bootstrap and jackknife estimates of standard error for the 25% trimmed mean, sample size  $n = 15$ .

Table 2 returns to the case of  $\hat{\theta}$  the correlation coefficient. Instead of real data we have a sampling experiment in which the true  $F$  is bivariate normal, true correlation  $\theta = .50$ , sample size  $n = 14$ . Table 2 is abstracted from a larger table in Efron (1981c), in which some of the methods for estimating a standard error required the sample size to be even.

The left side of Table 2 refers to  $\hat{\theta}$ , while the right side refers to  $\hat{\phi} = \tanh^{-1}(\hat{\theta}) = .5 \log(1 + \hat{\theta})/(1 - \hat{\theta})$ . For each estimator of standard error, the root mean squared error of estimation  $[E(\hat{\sigma} - \sigma)^2]^{1/2}$  is given in the column headed  $\sqrt{\text{MSE}}$ .

The bootstrap was run with  $B = 128$  and also with  $B = 512$ , the latter value yielding only slightly better estimates in accordance with the results of Section 9. Further increasing  $B$  would be pointless. It can be shown that  $B = \infty$  gives  $\sqrt{\text{MSE}} = .063$  for  $\hat{\theta}$ , only .001 less than  $B = 152$ . The normal theory estimate (2.5), which we know to be ideal for this sampling experiment, has  $\sqrt{\text{MSE}} = .056$ .

We can compromise between the totally nonparametric bootstrap estimate  $\hat{\sigma}$  and the totally parametric bootstrap estimate  $\hat{\sigma}_{\text{NORM}}$ . This is done in lines 3, 4, and 5 of Table 2. Let  $\hat{\Sigma} = \sum_{i=1}^n (z_i - \bar{z})(z_i - \bar{z})'/n$  be the sample covariance matrix of the observed data. The *normal smoothed bootstrap* draws the bootstrap sample from  $\hat{F} \otimes N_2(0, .25\hat{\Sigma})$ ,  $\otimes$  indicating convolution. This amounts to estimating  $F$  by an equal mixture of the  $n$  distributions  $N_2(z_i, .25\hat{\Sigma})$ , that is by a normal window estimate. Each point  $z_j^*$  in a smoothed bootstrap sample is the sum of a randomly selected original data point  $z_j$ , plus an independent bivariate normal point  $z_j \sim N_2(0, .25\hat{\Sigma})$ . Smoothing makes little difference on the left side of the table, but is

spectacularly effective in the  $\hat{\phi}$  case. The latter result is suspect since the true sampling distribution is bivariate normal, and the function  $\hat{\phi} = \tanh^{-1}\hat{\theta}$  is specifically chosen to have nearly constant standard error in the bivariate-normal family. The *uniform smoothed bootstrap* samples from  $\hat{F} \otimes U(0, .25\hat{\Sigma})$ , where  $U(0, .25\hat{\Sigma})$  is the uniform distribution on a rhombus selected so  $U$  has mean vector 0 and covariance matrix  $.25\hat{\Sigma}$ . It yields moderate reductions in  $\sqrt{MSE}$  for both sides of the table.

| Summary Statistics for 200 Trials            |   |         |     |              |   |         |     |              |
|--|---|---------|-----|--------------|---|---------|-----|--------------|
|  | Standard Error Estimates for $\hat{\theta}$ |         |     |              | Standard Error Estimates for $\hat{\phi}$ |         |     |              |
|  | Ave   | Std Dev | CV  | $\sqrt{MSE}$ | Ave                                       | Std Dev | CV  | $\sqrt{MSE}$ |
| 1. Bootstrap B = 128                         | .206  | .066    | .32 | .067         | .301                                      | .065    | .22 | .065         |
| 2. Bootstrap B = 512                         | .206  | .063    | .31 | .064         | .301                                      | .062    | .21 | .062         |
| 3. Normal Smoothed Bootstrap B = 128         | .200  | .060    | .30 | .063         | .296                                      | .041    | .14 | .041         |
| 4. Uniform Smoothed Bootstrap B = 128        | .205  | .061    | .30 | .062         | .298                                      | .058    | .19 | .058         |
| 5. Uniform Smoothed Bootstrap B = 512        | .205  | .059    | .29 | .060         | .296                                      | .052    | .18 | .052         |
| 6. Jackknife                                 | .223  | .085    | .38 | .085         | .314                                      | .090    | .29 | .091         |
| 7. Delta Method<br>(Infinitesimal Jackknife) | .175  | .058    | .33 | .072         | .244                                      | .052    | .21 | .076         |
| 8. Normal Theory                             | .217  | .056    | .26 | .056         | .302                                      | 0       | 0   | .003         |
| True Standard Error                          | .218  |         |     |              | .299                                      |         |     |              |

**Table 2.** Estimates of standard error for the correlation coefficient  $\hat{\theta}$  and for  $\hat{\phi} = \tanh^{-1}\hat{\theta}$ , sample size  $n = 14$ , distribution  $F$  bivariate normal with true correlation  $\rho = .5$ . From a larger table in Efron (1981c).

Line 6 of Table 2 refers to the *delta method*, which is the most common method of assigning nonparametric standard error. Surprisingly enough, it is badly biased downwards on both sides of the table. The delta method, also known as the method of statistical differentials, the Taylor series method, and the infinitesimal jackknife, are discussed in Section 10.

### 3. Examples.

#### Example 1: Cox's proportional hazards model

In this section we apply bootstrap standard error estimation to some complicated statistics.

The data for this example come from a study of leukemia remission times in mice, taken from Cox (1972). They consist of measurements of remission time ( $y$ ) in weeks for two groups,

treatment ( $x = 0$ ) and control ( $x = 1$ ), and a 0-1 variable ( $\delta_i$ ) indicating whether or not the remission time is censored (0) or complete (1). There are 21 mice in each group.

The standard regression model for censored data is Cox's proportional hazards model (Cox 1972). It assumes that the hazard function  $h(t | x)$ , the probability of going into remission in next instant given no remission up to time  $t$  for a mouse with covariate  $x$ , is of the form

$$h(t | x) = h_0(t)e^{\beta x}. \quad (3.1)$$

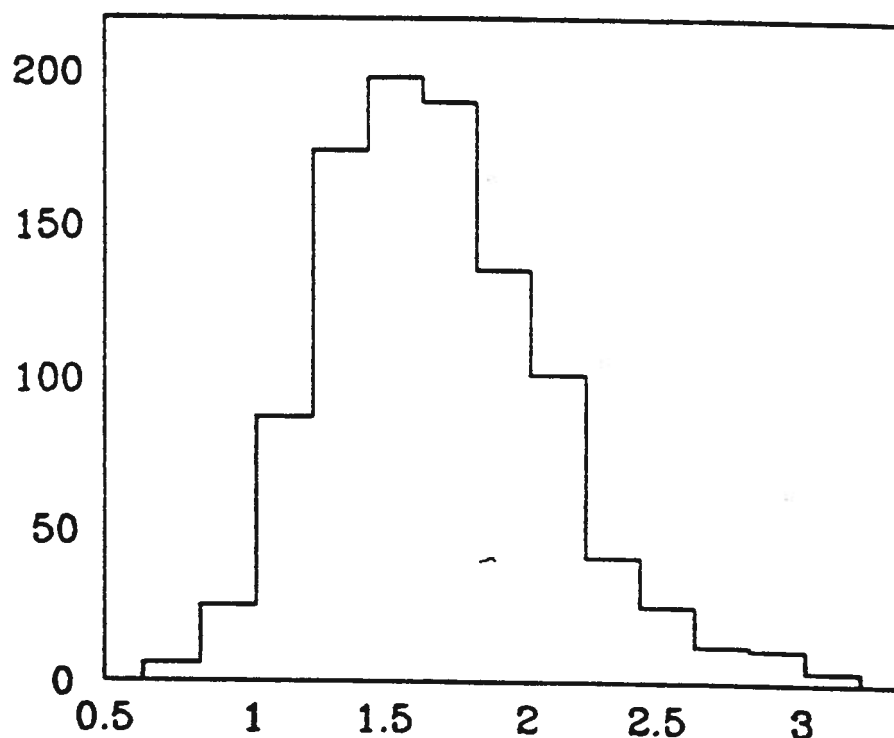
Here  $h_0(t)$  is an arbitrary unspecified function. Since  $x$  here is a group indicator, this means simply that the hazard for the control group is  $e^\beta$  times the hazard for the treatment group. The regression parameter  $\beta$  is estimated independently of  $h_0(t)$  through maximization of the so called "partial likelihood"

$$PL = \prod_{i \in D} \frac{e^{\beta x_i}}{\sum_{j \in R_i} e^{\beta x_j}} \quad (3.2)$$

where  $D$  is the set of indices of the failure times and  $R_i$  is the set of indices of those at risk at time  $y_i$ . This maximization requires an iterative computer search.

The estimate  $\hat{\beta}$  for these data turns out to be 1.51. Taken literally, this says that the hazard rate is  $e^{1.51} = 4.33$  times higher in the control group than in the treatment group, so the treatment is very effective. What's the standard error of  $\hat{\beta}$ ? The usual asymptotic maximum likelihood theory, one over the square root of the observed Fisher information, gives an estimate of .41. Despite the complicated nature of the estimation procedure, we can also estimate the standard error using the bootstrap. We sample with replacement from the triples  $\{(y_1, x_1, \delta_1), \dots, (y_{42}, x_{42}, \delta_{42})\}$ . For each bootstrap sample  $\{(y_1^*, x_1^*, \delta_1^*), \dots, (y_{42}^*, x_{42}^*, \delta_{42}^*)\}$  we form the partial likelihood and numerically maximize it to produce the bootstrap estimate  $\hat{\beta}^*$ . A histogram of 1000 bootstrap values is shown in Figure 3.

The bootstrap estimate of the standard error of  $\hat{\beta}$  based on these 1000 numbers is .42. Although that the bootstrap and standard estimates agree, it is interesting to note that the bootstrap distribution is skewed to the right. This leads us to ask: is there other information that we can extract from the bootstrap distribution other than a standard error estimate? The answer is yes—in particular, the bootstrap distribution can be used to form a confidence interval for  $\beta$ , as we will see in Section 9. The shape of the bootstrap distribution will help



**Figure 3.** Histogram of 1000 bootstrap replications for the mouse leukemia data

determine the shape of the confidence interval.

In this example our resampling unit was the triple  $(y_i, z_i, \delta_i)$ , and we ignored the unique elements of the problem, i.e. the censoring, and the particular model being used. In fact, there are other ways to bootstrap this problem. We'll see this when we discuss bootstrapping censored data in Section 5.

### **Example 2: Linear and Projection Pursuit Regression**

We illustrate an application of the bootstrap to standard linear least squares regression as well as to a non-parametric regression technique.

Consider the standard regression setup. We have  $n$  observations on a response  $Y$  and covariates  $(X_1, X_2, \dots, X_p)$ . Denote the  $i$ th observed vector of covariates by  $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ip})'$ .

The usual linear regression model assumes

$$E(Y_i) = \alpha + \sum_{j=1}^p \beta_j x_{ij} \quad (3.3)$$

Friedman and Stuetzle (1981) introduced a more general model, the *projection pursuit* regression model

$$E(Y_i) = \sum_{j=1}^m s_j(\mathbf{a}_j \cdot \mathbf{z}_i) \quad (3.4)$$

The  $p$ -vectors  $\mathbf{a}_j$  are unit vectors ("directions"), and the functions  $s_j(\cdot)$  are unspecified.

Estimation of  $\{\mathbf{a}_1, s_1(\cdot)\}, \dots, \{\mathbf{a}_m, s_m(\cdot)\}$  is performed in a forward stepwise manner as follows. Consider  $\{\mathbf{a}_1, s_1(\cdot)\}$ . Given a direction  $\mathbf{a}_1$ ,  $s_1(\cdot)$  is estimated by a non-parametric smoother (e.g. running mean) of  $y$  on  $\mathbf{a}_1 \cdot \mathbf{z}$ . The projection pursuit regression algorithm searches over all unit directions to find the direction  $\hat{\mathbf{a}}_1$  and associated function  $\hat{s}_1(\cdot)$  that minimize  $\sum_{i=1}^n (y_i - \hat{s}_1(\hat{\mathbf{a}}_1 \cdot \mathbf{z}_i))^2$ . Then residuals are taken and the next direction and function are determined. This process is continued until no additional term significantly reduces the residual sum of squares.

Notice the relation of the projection pursuit regression model to the standard linear regression model. When the function  $s_1(\cdot)$  is forced to be linear, and is estimated by the usual least squares method, a one term projection pursuit model is exactly the same as the standard linear regression model. That is to say, the fitted model  $\hat{s}_1(\hat{\mathbf{a}}_1 \cdot \mathbf{z}_i)$  exactly equals the least squares fit  $\hat{\alpha} + \sum_{j=1}^p \hat{\beta}_j x_{ij}$ . This is because the least squares fit, by definition, finds the best direction and the best linear function of that direction. Note also that adding another linear term  $\hat{s}_2(\hat{\mathbf{a}}_2 \cdot \mathbf{z}_2)$  would not change the fitted model since the sum of two linear functions is another linear function.

Hastie and Tibshirani (1984) applied the bootstrap to the linear and projection pursuit regression models to assess the variability of the coefficients in each. The data they considered are taken from Breiman and Friedman (1984). The response  $Y$  is Upland atmospheric ozone concentration (ppm); the covariates  $X_1$ - Sandburg Air Force base temperature ( $C^\circ$ ),  $X_2$ - inversion base height (ft.) ,  $X_3$ - Daggot pressure gradient (mmhg),  $X_4$ - visibility (miles), and  $X_5$ - day of the year. There are 330 observations. The number of terms ( $m$ ) in the model (3.4) is taken to be two. The projection pursuit algorithm chose directions  $\hat{\mathbf{a}}_1 = (.80, -.38, .37, -.24, -.14)'$

and  $\hat{a}_2 = (.07, .16, .04, -.05, -.98)'$ . These directions consist mostly of Sandburg Air Force temperature and day of the year respectively. (We don't show graphs of the estimated functions  $\hat{s}_1(\cdot)$  and  $\hat{s}_2(\cdot)$  although in a full analysis of the data they would also be of interest.) Forcing  $\hat{s}_1(\cdot)$  to be linear results the direction  $\hat{a}_1 = (.90, -.37, .03, -.14, -.19)'$ . These are just the usual least squares estimates  $\hat{\beta}_1, \dots, \hat{\beta}_p$  scaled so that  $\sum_1^p \beta_j^2 = 1$ .

To assess the variability of the directions, a bootstrap sample is drawn with replacement from  $(y_1, x_{11}, \dots, x_{15}), \dots, (y_{3301}, x_{3301}, \dots, x_{3305})$  and the projection pursuit algorithm is applied. Figures 4 and 5 show histograms of the directions  $\hat{a}_1^*$  and  $\hat{a}_2^*$  for 200 bootstrap replications. Also shown in Figure 4 (broken histogram) are the bootstrap replications of  $\hat{a}_1$  with  $\hat{s}_1(\cdot)$  forced to be linear.

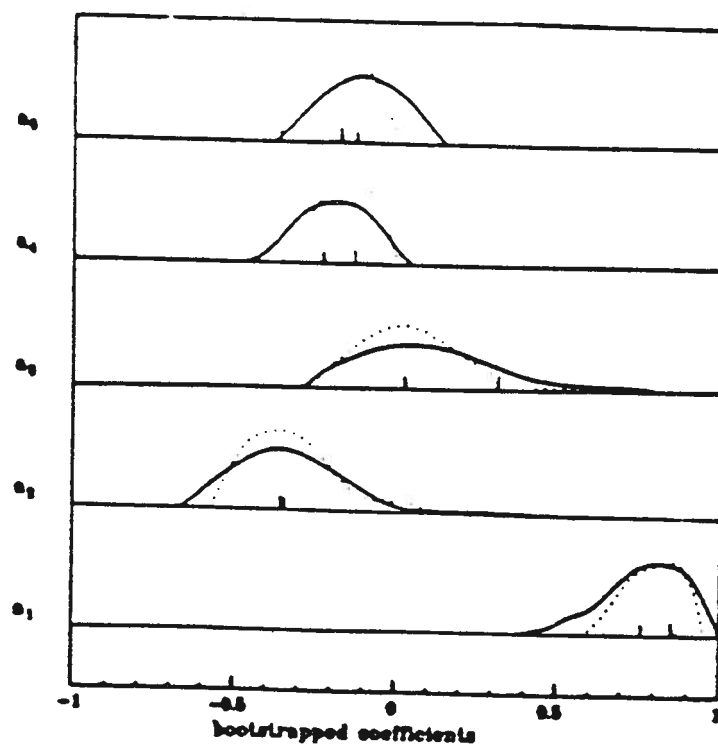
The first direction of the projection pursuit model is quite stable and only slightly more variable than the corresponding linear regression direction. But the second direction is extremely unstable! It is clearly unwise to put any faith in the second direction of the original projection pursuit model.

### Example 3: Cox's Model and Local Likelihood Estimation

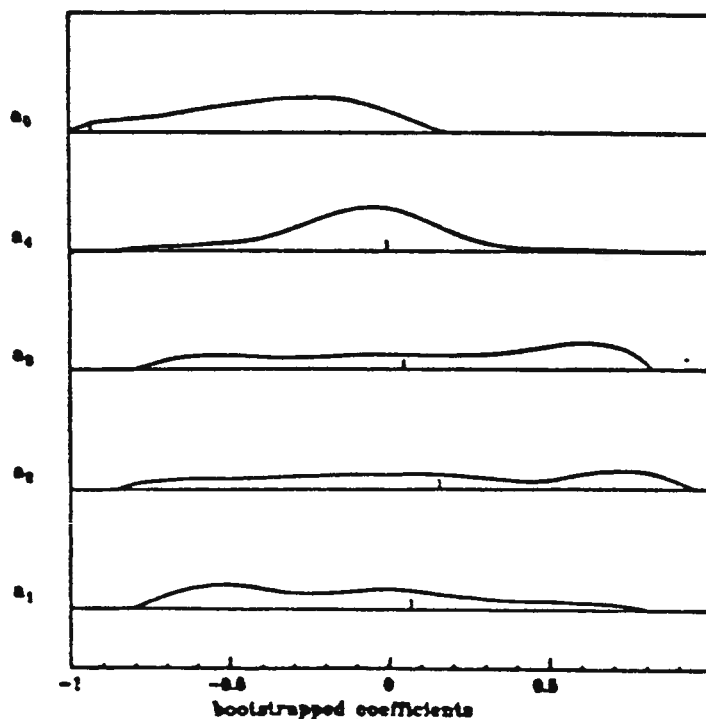
In this example, we return to Cox's proportional hazards model described in Example 1, but with a few added twists.

The data that we'll discuss come from the Stanford heart transplant program and are given in Miller and Halpern (1983). The response  $y$  is survival time in weeks after a heart transplant, the covariate  $x$  is age at transplant, and the 0-1 variable  $\delta$  indicates whether the survival time is censored (0) or complete (1). There are measurements on 157 patients. A proportional hazards model was fit to these data, with a quadratic term i.e.  $h(t|x) = h_0(t)e^{\beta_1 x + \beta_2 x^2}$ . Both  $\hat{\beta}_1$  and  $\hat{\beta}_2$  are highly significant; the broken curve in Figure 6 is  $\hat{\beta}_1 x + \hat{\beta}_2 x^2$  as a function of  $x$ .

For comparison, Figure 6 shows (solid line) another estimate. This was computed using *local likelihood estimation* (Tibshirani and Hastie 1984). Given a general proportional hazards model of the form  $h(t|x) = h_0(t)e^{s(x)}$ , the local likelihood technique assumes nothing about the parametric form of  $s(x)$ ; instead it estimates  $s(x)$  non-parametrically using a kind of local averaging. The algorithm is very computationally intensive, and standard maximum likelihood

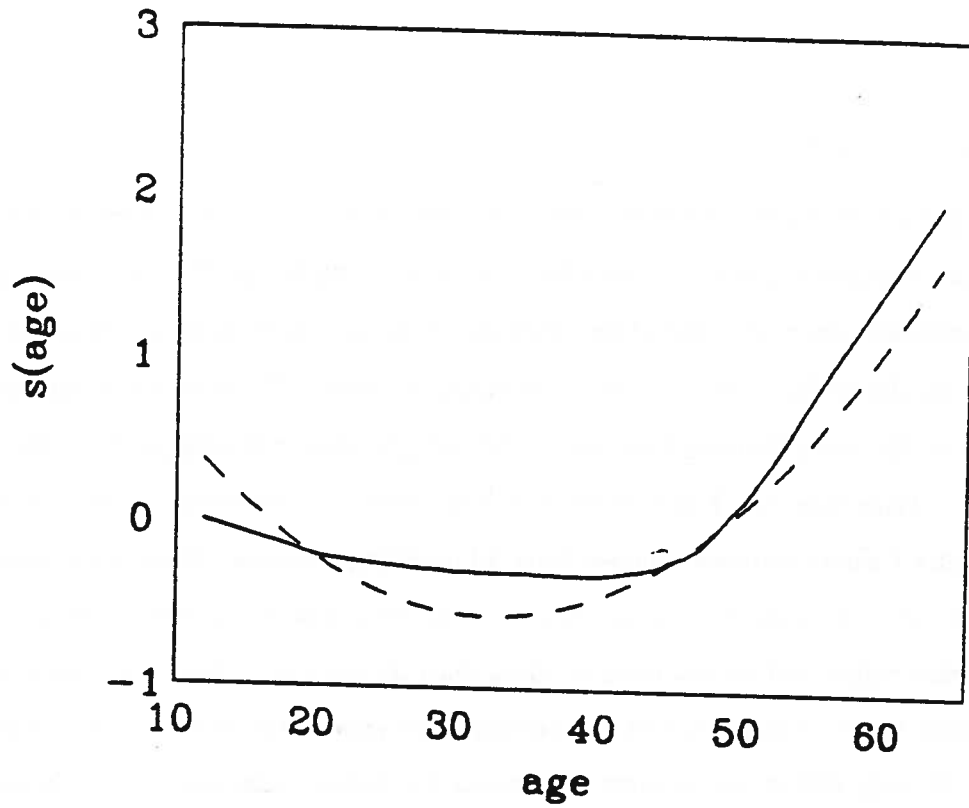


**Figure 4.** Smoothed histograms of the bootstrapped coefficients for the first term in the projection pursuit regression model. Solid histograms are for the usual projection pursuit model; the dotted histograms are for linear  $\hat{s}(\cdot)$ .

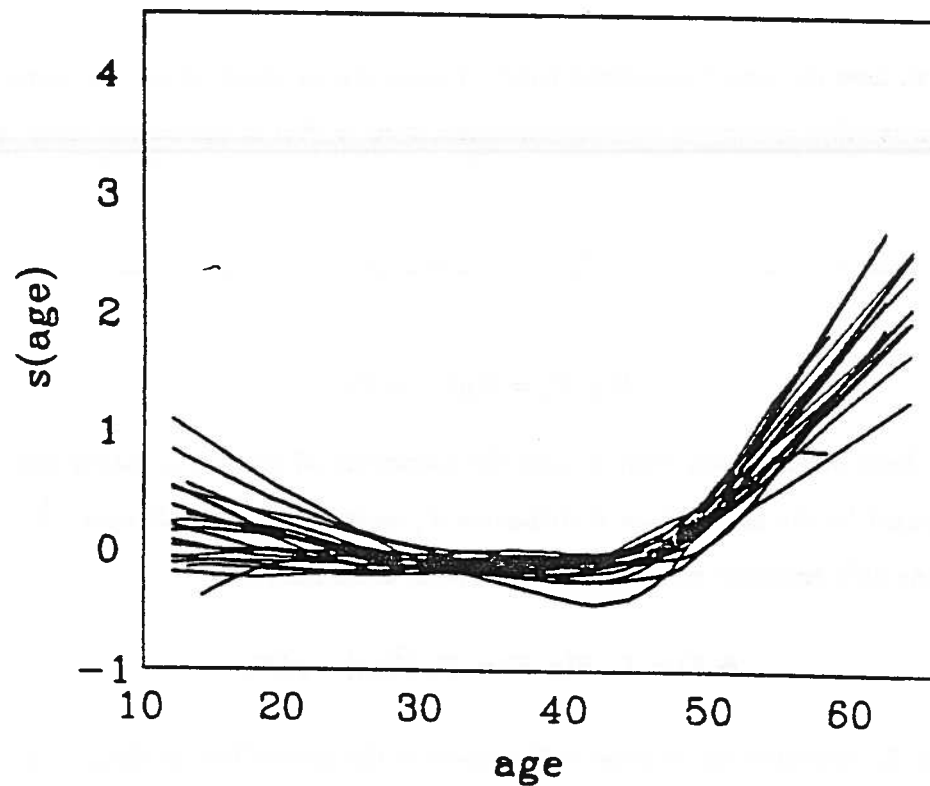


**Figure 5.** Smoothed histograms of the bootstrapped coefficients for the second term in the projection pursuit model.





**Figure 6.** Estimates of log relative risk for the Stanford heart transplant data. Broken curve: parametric estimate. Solid curve: local likelihood estimate



**Figure 7.** 20 Bootstraps of the local likelihood estimate, for the Stanford heart transplant data

theory cannot be applied.

A comparison of the two functions reveals an important qualitative difference: the parametric estimate suggests that the hazard decreases sharply up to age 34, then rises; the local likelihood estimate stays approximately constant up to age 45 then rises. Has the forced fitting of a quadratic function produced a misleading result? To answer this question, we can bootstrap the local likelihood estimate. We sample with replacement from the triples  $\{(y_1, x_1, \delta_1) \dots (y_{157}, x_{157}, \delta_{157})\}$  and apply the local likelihood algorithm to each bootstrap sample. Figure 7 shows estimated curves from 20 bootstrap samples. Some of the curves are flat up to age 45, others are decreasing. Hence the original local likelihood estimate is highly variable in this region and on the basis of these data we can't determine the true behaviour of the function there. A look back at the original data shows that while half of the patients were under 45, only 13% of the patients were under 30. Figure 7 also shows that the estimate is stable near the middle ages but unstable for the older patients.

#### 4. Other Measures of Statistical Error.

So far we have discussed statistical error, or accuracy, in terms of the standard error. It is easy to assess other measures of statistical error, such as bias or prediction error, using the bootstrap.

Consider the estimation of bias. For a given statistic  $\hat{\theta}(y)$ , and a given parameter  $\mu(F)$ , let

$$R(y, F) = \hat{\theta}(y) - \mu(F). \quad (4.1)$$

(It will help keep our notation clear to call the parameter of interest  $\mu$  rather than  $\theta$ .) For example  $\mu$  might be the mean of the distribution  $F$ , assuming the sample space  $X$  is the real line, and  $\hat{\theta}$  the 25% trimmed mean. The bias of  $\hat{\theta}$  for estimating  $\mu$  is

$$\beta(F) = E_F R(y, F) = E_F \{\hat{\theta}(y)\} - \mu(F). \quad (4.2)$$

The notation  $E_F$  indicates expectation with respect to the probability mechanism appropriate to  $F$ , in this case  $y = (x_1, x_2, \dots, x_n)$  a random sample from  $F$ .

The bootstrap estimate of bias is

$$\begin{aligned}\hat{\beta} &= \beta(\hat{F}) = E_{\hat{F}} R(y^*, \hat{F}) \\ &= E_{\hat{F}} \{\hat{\theta}(y^*)\} - \mu(\hat{F}).\end{aligned}\quad (4.3)$$

As in Section 2,  $y^*$  denotes a random sample  $(x_1^*, x_2^*, \dots, x_n^*)$  from  $\hat{F}$ , i.e. a bootstrap sample. To numerically evaluate  $\hat{\beta}$ , all we do is change step (iii) of the bootstrap algorithm in Section 2 to

$$\begin{aligned}\hat{\beta}_B &= \frac{1}{B} \sum_{b=1}^B R(y^*(b), \hat{F}) = \frac{\sum_{b=1}^B \hat{\theta}^*(b)}{B} - \mu(\hat{F}) \\ &= \hat{\theta}^*(\cdot) - \hat{\mu}(F).\end{aligned}\quad (4.4)$$

As  $B \rightarrow \infty$ ,  $\hat{\beta}_B$  goes to  $\hat{\beta}$ , as given in (4.3).

As an example consider the blood serum data of Table 3. Suppose we wish to estimate the true mean  $\mu = E_F\{X\}$  of this population using  $\hat{\theta}$ , the 25% trimmed mean. We calculate  $\hat{\mu} = \mu(\hat{F}) = 2.39$ , the sample mean of the 54 observations, and  $\hat{\theta} = 2.24$ , the trimmed mean. The trimmed mean is lower because it discounts the effect of the large observations 6.4 and 9.4. It looks like the trimmed mean might be more robust for this type of data, and as a matter of fact a bootstrap analysis,  $B = 1000$ , gave estimated standard error  $\hat{\sigma} = .16$  for  $\hat{\theta}$ , compared to .21 for the sample mean. But what about bias?

|      |      |      |      |      |      |      |      |      |      |      |      |
|------|------|------|------|------|------|------|------|------|------|------|------|
| 0.1, | 0.1, | 0.2, | 0.4, | 0.4, | 0.6, | 0.8, | 0.8, | 0.9, | 0.9, | 1.3, | 1.3, |
| 1.4, | 1.5, | 1.6, | 1.6, | 1.7, | 1.7, | 1.7, | 1.8, | 2.0, | 2.0, | 2.2, | 2.2, |
| 2.2, | 2.3, | 2.3, | 2.4, | 2.4, | 2.4, | 2.4, | 2.4, | 2.4, | 2.5, | 2.5, | 2.5, |
| 2.7, | 2.7, | 2.8, | 2.9, | 2.9, | 2.9, | 3.0, | 3.1, | 3.1, | 3.2, | 3.2, | 3.3, |
| 3.3, | 3.5, | 4.4, | 4.5, | 6.4, | 9.4  |      |      |      |      |      |      |

**Table 3.** BHCG blood serum levels for 54 patients having metastasized breast cancer, presented in ascending order.

The same 1000 bootstrap replications which gave  $\hat{\sigma} = .164$  also gave  $\hat{\theta}^*(\cdot) = 2.29$ , so

$$\hat{\beta} = 2.29 - 2.39 = -0.10 \quad (4.5)$$

according to (4.4). (The estimated standard deviation of  $\hat{\beta}_B - \hat{\beta}$  due to the limitations of having  $B = 1000$  bootstraps is only 0.005 in this case, so we can ignore the difference between  $\hat{\beta}_B$  and  $\hat{\beta}$ .) Whether or not a bias of magnitude  $-0.10$  is too large depends on the context of the

problem. If we attempt to remove the bias by subtraction, we get  $\hat{\theta} - \hat{\beta} = 2.24 - (0.10) = 2.34$ , which is close to the sample mean 2.39. Removing bias in this easy is frequently a bad idea, see Hinkley (1978), but at least the bootstrap analysis has given us a reasonable picture of the bias and standard error of  $\hat{\theta}$ .

Here is another measure of statistical accuracy, different than either bias or standard error. Let  $\hat{\theta}(y)$  be the 25% trimmed mean and  $\mu(F)$  be the mean of  $F$ , as in the serum example, and also let  $\hat{i}(y)$  be the interquartile range, the distance between the 25th and 75th percentiles of the sample  $y = (x_1, x_2, \dots, x_n)$ . Define

$$R(y, F) = \frac{\hat{\theta}(y) - \mu(F)}{\hat{i}(y)}. \quad (4.6)$$

$R$  is like a Student's  $t$  statistic, except that we have substituted the 25% trimmed mean for the sample mean, and the interquartile range for the standard deviation.

Suppose we know the 5th and 95th percentiles of  $R(y, F)$ , say  $\rho^{(.05)}(F)$  and  $\rho^{(.95)}(F)$ , where the definition of  $\rho^{(.05)}(F)$  is

$$\text{Prob}_F\{R(y, F) < \rho^{(.05)}(F)\} = .05, \quad (4.7)$$

and similarly for  $\rho^{(.95)}(F)$ . The relationship  $\text{Prob}_F\{\rho^{(.05)} \leq R < \rho^{(.95)}\} = .90$  combines with definition (4.6) to give a central 90% "t interval" for the mean  $\mu(F)$ ,

$$\mu \in [\hat{\theta} - \hat{i}\rho^{(.95)}, \hat{\theta} - \hat{i}\rho^{(.05)}]. \quad (4.8)$$

Of course we don't know  $\rho^{(.05)}(F)$  and  $\rho^{(.95)}(F)$ , but we can approximate them by their bootstrap estimates  $\rho^{(.05)}(\hat{F})$  and  $\rho^{(.95)}(\hat{F})$ . A bootstrap sample  $y^*$  gives a bootstrap value of (4.6),  $R(y^*, \hat{F}) = (\hat{\theta}(y^*) - \mu(\hat{F})/\hat{i}(y^*))$ , where  $\hat{i}(y^*)$  is the interquartile range of the bootstrap data  $x_1^*, x_2^*, \dots, x_n^*$ . For any fixed number  $\rho$ , the bootstrap estimate of  $\text{Prob}_F\{R < \rho\}$  based on  $B$  bootstrap samples is

$$\#\{R(y^*(b), \hat{F}) < \rho\}/B. \quad (4.9)$$

By keeping track of the empirical distribution of  $R(y^*(b), \hat{F})$ , we can pick off the values of  $\rho$  which make (4.9) equal .05 and .95. These approach  $\rho^{(.05)}(\hat{F})$  and  $\rho^{(.95)}(\hat{F})$  as  $B \rightarrow \infty$ .

For the serum data,  $B = 1000$  bootstrap replications gave  $\rho^{(.05)}(\hat{F}) = -.303$  and  $\rho^{(.95)}(\hat{F})$

= .078. Substituting these values into (4.9), and using the observed estimates  $\hat{\theta} = 2.24$ ,  $\hat{\sigma} = 1.40$ , gives

$$\mu \in [2.13, 2.66] \quad (4.10)$$

as a central 90% "bootstrap  $t$  interval" for the true mean  $\mu(F)$ . This compares with the standard  $t$  interval based on 53 degrees of freedom  $\bar{x} \pm 1.67\sigma = [2.04, 2.74]$ . Here  $\sigma = .21$  is the usual estimate of standard error (1.3).

It is interesting to notice that if we discard the 54th observation 9.4, then  $\sigma$  decreases to .16, and the Student's  $t$  interval  $\bar{x} \pm 1.67\sigma$  equals  $[2.12, 2.66]$  which is almost exactly the same as (4.10)! Bootstrap confidence intervals are discussed further in Sections 7 and 8. They require more bootstrap replications than does  $\hat{\sigma}$ , on the order of  $B = 1000$  rather than  $B = 50$  or 100. This point is discussed briefly in Section 9.

By now it should be clear that we can use any random variable  $R(y, F)$  to measure accuracy, not just (4.1) or (4.6), and then estimate  $E_F\{R(y, F)\}$  by its bootstrap value  $E_{\hat{F}}\{R(y^*, \hat{F})\} \doteq \sum_{b=1}^B R(y^*(b), \hat{F})/B$ . Similarly we can estimate  $E_F R(y, F)^2$  by  $E_{\hat{F}} R(y^*, \hat{F})^2$ , etc. Efron (1983) considers the prediction problem, in which a training set of data is used to construct a prediction rule. A naive estimate of the prediction rule's accuracy is the proportion of correct guesses it makes on its own training set, but this can be greatly overoptimistic since the prediction rule is explicitly constructed to minimize errors on the training set. In this case, a natural choice of  $R(y, F)$  is the overoptimism, the difference between the naive estimate and the actual success rate of the prediction rule for new data. Efron (1983) gives the bootstrap estimate of overoptimism, and shows that it is closely related to cross-validation, the usual method of estimating overoptimism. The paper goes on to show that some modifications of the bootstrap estimate greatly outperform both cross-validation and the bootstrap.

## 5. More Complicated Data Sets.

The bootstrap is not restricted to situations where the data is a simple random sample from a single distribution. Suppose for instance that the data consists of two independent random samples,

$$U_1, U_2, \dots, U_m \sim F \quad \text{and} \quad V_1, V_2, \dots, V_n \sim G, \quad (5.1)$$

|                 | Summary Statistics for $\hat{\sigma}_B$ |          |      |
|-----------------|---|----------|------|
|                 | Average                                 | St. Dev. | C.V. |
| B=100:          | .165                                    | .030     | .18  |
| B=200:          | .166                                    | .031     | .19  |
| True $\sigma$ : | .167                                    |          |      |

**Table 4.** Bootstrap estimate of Standard Error for the Hodges-Lehmann two-sample shift estimate;  $m = 6$ ,  $n = 9$ ; true distributions  $F$  and  $G$  both Uniform  $[0, 1]$ . The table shows summary statistics for  $\hat{\sigma}_B$ , over 100 trials of this situation.

where  $F$  and  $G$  are possibly different distributions on the real line. Suppose also that the statistic of interest is the Hodges-Lehmann shift estimate

$$\hat{\theta} = \text{median}\{F_j - U_i; \quad i = 1, 2, \dots, m, \quad j = 1, 2, \dots, n\}. \quad (5.2)$$

Having observed  $U_1 = u_1, U_2 = u_2, \dots, V_n = v_n$ , we desire an estimate for  $\sigma(F, G)$ , the standard error of  $\hat{\theta}$ .

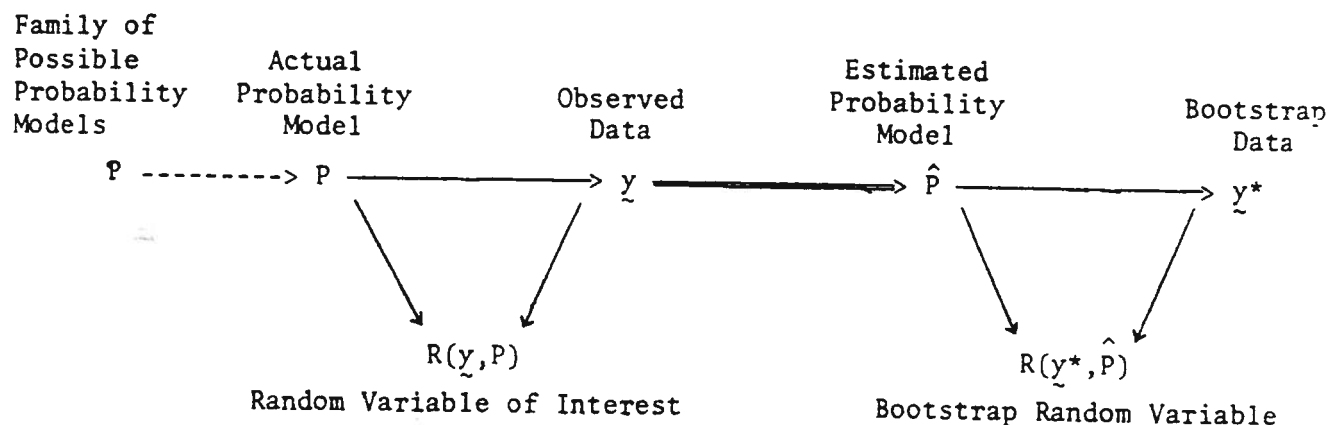
The bootstrap estimate of  $\sigma(F, G)$  is  $\hat{\sigma} = \sigma(\hat{F}, \hat{G})$ , where  $\hat{F}$  is the empirical distribution of  $u_1, u_2, \dots, u_m$ , and  $\hat{G}$  is the empirical distribution of  $v_1, v_2, \dots, v_n$ . It is easy to modify the Monte Carlo algorithm of Section 2 to numerically evaluate  $\hat{\sigma}$ . Let  $y = (u_1, u_2, \dots, v_n)$  be the observed data vector. A bootstrap sample  $y^* = (u_1^*, u_2^*, \dots, u_m^*, v_1^*, v_2^*, \dots, v_n^*)$  consists of a random sample  $U_1^*, \dots, U_m^*$  from  $\hat{F}$  and an independent random sample  $V_1^*, \dots, V_n^*$  from  $\hat{G}$ . With only this modification, steps (i) through (ii) of the Monte Carlo algorithm produce  $\hat{\sigma}_B$ , (2.4), approaching  $\hat{\sigma}$  as  $B \rightarrow \infty$ .

Table 4 reports on a simulation experiment investigating how well the bootstrap works on this problem. 100 trials of situation (5.1) were run, with  $m = 6$ ,  $n = 9$ ,  $F$  and  $G$  both Uniform  $[0, 1]$ . For each trial, both  $B = 100$  and  $B = 200$  bootstrap replications were generated. The bootstrap estimate  $\hat{\sigma}_B$  was nearly unbiased for the true standard error  $\sigma(F, G) = .167$  for either  $B = 100$  or  $B = 200$ , with a quite small standard deviation from trial to trial. The improvement in going from  $B = 100$  to  $B = 200$  is too small to show up in this experiment.

In practice, statisticians must often consider quite complicated data structures: time series models, multi-factor layouts, sequential sampling, censored and missing data, etc. Figure 8 illustrates how the bootstrap estimation process proceeds in a general situation. The actual probability mechanism  $P$  which generates the observed data  $y$  belongs to some fam-

ily  $\mathcal{P}$  of possible probability mechanism. In the Hodges-Lehmann example,  $P = (F, G)$ , a pair of distributions on the real line,  $\mathcal{P}$  equals the family of all such pairs, and  $y = (u_1, u_2, \dots, u_m, v_1, v_2, \dots, v_n)$  is generated by random sampling  $m$  times from  $F$  and  $n$  times from  $G$ .

We have a random variable of interest  $R(y, P)$ , which depends on both  $y$  and the unknown model  $P$ , and we wish to estimate some aspect of the distribution of  $R$ . In the Hodges-Lehmann example,  $R(y, P) = \hat{\theta}(y) - E_P\{\hat{\theta}\}$ , and we estimated  $\sigma(P) = E_P R(y, P)^2$ , the standard error of  $\hat{\theta}$ . As before, the notation  $E_P$  indicates expectation when  $y$  is generated according to mechanism  $P$ .



**Figure 8.** A schematic illustration of the bootstrap process for a general probability model  $P$ . The expectation of  $R(y, P)$  is estimated by the bootstrap expectation of  $R(y^*, \hat{P})$ . The double arrow indicates the crucial step in applying the bootstrap.

We assume that we have some way of estimating the entire probability model  $P$  from the data  $y$ , producing the estimate called  $\hat{P}$  in Figure 8. (In the two-sample problem,  $\hat{P} = (\hat{F}, \hat{G})$ , the pair of empirical distributions.) *This is the crucial step for the bootstrap.* It can be

carried out either parametrically or nonparametrically, by maximum likelihood or by some other estimation technique.

Once we have  $\hat{P}$ , we can use Monte Carlo methods to generate bootstrap data sets  $y^*$ , according to the same rules by which  $y$  is generated from  $P$ . The bootstrap random variable  $R(y^*, \hat{P})$  is observable, since we know  $\hat{P}$  as well as  $y^*$ , so the distribution of  $R(y^*, \hat{P})$  can be found by Monte Carlo sampling. The bootstrap estimate of  $E_P R(y, P)$  is then  $E_{\hat{P}} R(y^*, \hat{P})$ , and likewise for estimating any other aspect of  $R(y, P)$ 's distribution.

A regression model is a familiar example of a complicated data structure. We observe  $y = (y_1, y_2, \dots, y_n)$ , where

$$y_i = g(\beta, t_i) + \epsilon_i \quad i = 1, 2, \dots, n. \quad (5.3)$$

Here  $\beta$  is a vector of unknown parameters we wish to estimate; for each  $i$ ,  $t_i$  is an observed vector of covariates; and  $g$  is a known function of  $\beta$  and  $t_i$ , for instance  $e^{\beta' t_i}$ . The  $\epsilon_i$  are an i.i.d. sample from some unknown distribution  $F$  on the real line,

$$\epsilon_1, \epsilon_2, \dots, \epsilon_n \sim F, \quad (5.4)$$

where  $F$  is usually assumed to be centered at 0 in some sense, perhaps  $E\{\epsilon\} = 0$  or  $\text{Prob}\{\epsilon < 0\} = .5$ . The probability model is  $P = (\beta, F)$ ; (5.3) and (5.4) describe the step  $P \rightarrow y$  in Figure 5B. The covariates  $t_1, t_2, \dots, t_n$ , like the sample size  $n$  in the simple problem (1.1), are considered fixed at their observed values.

For every choice of  $\beta$  we have a vector  $g(\beta) = (g(\beta, t_1), g(\beta, t_2), \dots, g(\beta, t_n))$  of predicted values for  $y$ . Having observed  $y$ , we estimate  $\beta$  by minimizing some measure of distance between  $g(\beta)$  and  $y$ ,

$$\hat{\beta} : \min_{\beta} D(y, g(\beta)). \quad (5.5)$$

The most common choice of  $F$  is  $D(y, y) = \sum_{i=1}^n \{y_i - g(\beta, t_i)\}^2$ .

How accurate is  $\hat{\beta}$  as an estimate of  $\beta$ ? Let  $R(y, P)$  equal the vector  $\hat{\beta} - \beta$ . A familiar measure of accuracy is the mean square error matrix

$$\mathfrak{F}(P) = E_P(\hat{\beta} - \beta)(\hat{\beta} - \beta)' = E_P R(y, P) R(y, P)'. \quad (5.6)$$



The bootstrap estimate of accuracy  $\hat{\Phi} = \Phi(\hat{P})$  is obtained by following through Figure 8.

There is an obvious choice for  $\hat{P} = (\hat{\beta}, \hat{F})$  in this case. The estimate  $\hat{\beta}$  is obtained from (5.5). Then  $\hat{F}$  is the empirical distribution of the residuals.

$$\hat{F} : \text{mass } \frac{1}{n} \text{ on } \hat{\epsilon}_i \equiv y_i - g(\hat{\beta}, t_i), \quad i = 1, \dots, n. \quad (5.7)$$

A bootstrap sample  $y^*$  is obtained by following rules (5.3), (5.4),

$$y_i^* = g(\hat{\beta}, t_i) + \epsilon_i^*, \quad i = 1, 2, \dots, n, \quad (5.8)$$

where  $\epsilon_1^*, \epsilon_2^*, \dots, \epsilon_n^*$  is an i.i.d. sample from  $\hat{F}$ . [Notice that the  $\epsilon_n^*$  are independent bootstrap variates, even though the  $\hat{\epsilon}_i$  are not independent variates in the usual sense.]

Each bootstrap sample  $y^*(b)$  gives a bootstrap value  $\hat{\beta}^*(b)$ ,

$$\hat{\beta}^*(b) : \min_{\beta} D(y^*(b), g(\beta)), \quad (5.9)$$

as in (5.5). The estimate

$$\hat{\Phi}_B = \frac{\sum_{b=1}^B \{\hat{\beta}^*(b) - \hat{\beta}^*(\cdot)\} \{\hat{\beta}^*(b) - \hat{\beta}^*(\cdot)\}'}{B} \quad (5.10)$$

approaches the bootstrap estimate  $\hat{\Phi}$  as  $B \rightarrow \infty$ . (We could just as well divide by  $B - 1$  in (5.10).)

In the case of ordinary least squares regression, where  $g(\beta, t_i) = \beta' t_i$  and  $D(y, g) = \sum_{i=1}^n (y_i - g_i)^2$ , Section 7 of Efron (1979) shows that the bootstrap estimate,  $B = \infty$ , can be calculated without Monte Carlo sampling, and is

$$\hat{\Phi} = \hat{\sigma}^2 \left( \sum_{i=1}^n t_i t_i' \right)^{-1} \left[ \hat{\sigma}^2 \equiv \sum_{i=1}^n \hat{\epsilon}_i^2 / n \right]. \quad (5.11)$$

This is the usual Gauss-Markov answer, except for the divisor  $n$  in the definition of  $\hat{\sigma}^2$ .

There is another, simpler way to bootstrap a regression problem. We can consider each covariate-response pair  $x_i = (t_i, y_i)$  to be a single data point obtained by simple random sampling from a distribution  $F$ . If the covariate vector  $t_i$  is  $p$ -dimensional,  $F$  is a distribution on  $p + 1$  dimensions. Then we apply the bootstrap as described originally in Section 2 to the data set  $x_1, x_2, \dots, x_n \stackrel{\text{iid}}{\sim} F$ .

The two bootstrap methods for the regression problem are asymptotically equivalent, but can perform quite differently in small-sample situations. The class of possible probability models  $P$  is different for the two methods. The simple method, described last, takes less advantage of the special structure of the regression problem. It does *not* give answer (5.11) in the case of ordinary least squares. On the other hand the simple method gives a trustworthy estimate of  $\hat{\beta}$ 's variability *even if the regression model is not correct*. The bootstrap, as outlined in Figure 5B, is very general, but because of this generality there will often be more than one bootstrap solution for a given problem.

As the final example of this Section, we discuss *censored data*. The ages of 97 men at a California retirement center, Channing House, were observed either at death (an uncensored observation) or at the time the study ended (a censored observation). The data set  $y = \{(x_1, d_1), (x_2, d_2), \dots, (x_{97}, d_{97})\}$ , where  $x_i$  was the age of the  $i$ th man observed, and

$$d_i = \begin{cases} 1 & \text{if } x_i \text{ uncensored} \\ 0 & \text{if } x_i \text{ censored.} \end{cases} \quad (5.12)$$

Thus (777, 1) represents a Channing House man observed to die at age 777 months, while (843, 0) represents a man 843 months old when the study ended. His observation could be written "843+", and in fact  $d_i$  is just an indicator for the absence or presence of a "+".

A typical data point  $(X_i, D_i)$  can be thought of as generated in the following way: a real lifetime  $X_i^o$  is selected randomly according to a survival curve

$$S^o(t) \equiv \text{Prob}\{X_i^o > t\}, \quad (0 \leq t < \infty) \quad (5.13)$$

and a censoring time  $W_i$  is independently selected according to another survival curve

$$R(t) \equiv \text{Prob}\{W_i > t\}, \quad (0 \leq t < \infty). \quad (5.14)$$

The statistician gets to observe

$$X_i = \min\{X_i^o, W_i\} \quad (5.15)$$

and

$$D_i = \begin{cases} 1 & \text{if } X_i = X_i^o \\ 0 & \text{if } X_i = W_i. \end{cases} \quad (5.16)$$

Note:  $1 - S^\circ(t)$  and  $1 - R(t)$  are the cumulative distribution functions for  $X_i^\circ$  and  $W_i$  respectively; with censored data it is more convenient to consider survival curves than c.d.f.'s.

Under assumptions (5.12)–(5.15) there is a simple formula for the nonparametric MLE of  $S^\circ(t)$ , called the *Kaplan-Meier estimator*, Kaplan and Meier (1958). For convenience suppose  $x_1 < x_2 < x_3 \cdots < x_n$ ,  $n = 97$ . Then the Kaplan-Meier estimate is

$$\hat{S}^\circ(t) = \prod_{j=1}^{k_t} \left( \frac{n - i}{n - i + 1} \right)^{d_i}, \quad (5.17)$$

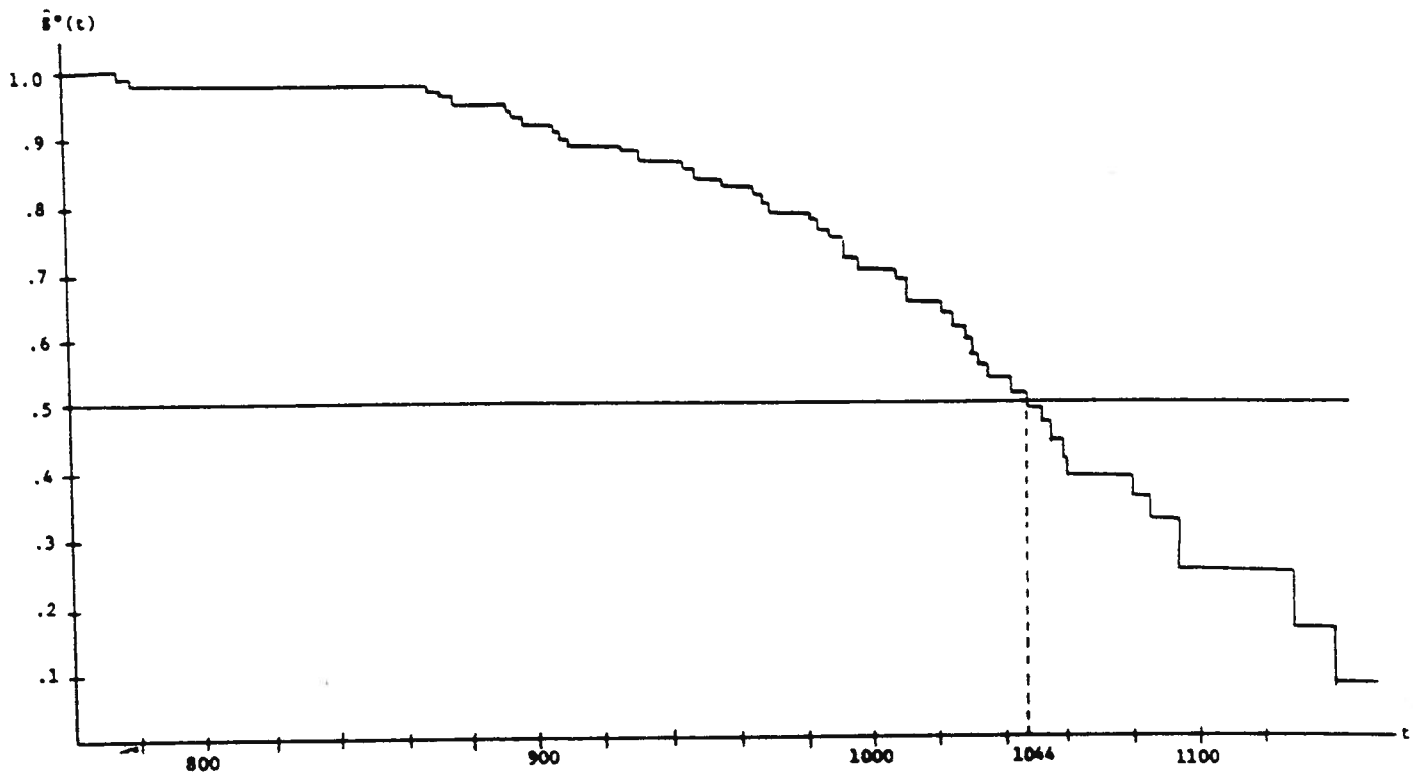
where  $k_t$  is the value of  $k$  such that  $t \in [x_k, x_{k+1})$ . In the case of no censoring,  $\hat{S}^\circ(t)$  is equivalent to the observed empirical distribution of  $x_1, x_2, \dots, x_n$ , but otherwise (5.16) corrects the empirical distribution to account for censoring. Likewise

$$\hat{R}(t) = \prod_{j=1}^{k_t} \left( \frac{n - i}{n - i + 1} \right)^{1-d_i} \quad (5.18)$$

is the Kaplan-Meier estimate of the censoring curve  $R(t)$ .

Figure 9 shows  $\hat{S}^\circ(t)$  for the Channing House men. It crosses the 50% survival level at  $\hat{\theta} = 1044$  months. Call this value the observed median lifetime. We can use the bootstrap to assign a standard error to the observed median.

The probability mechanism is  $P = (S^\circ, R)$ ;  $P$  produces  $(X_i^\circ, D_i)$  according to (5.12)–(5.15), and  $y = \{(x_1, d_1), \dots, (x_n, d_n)\}$  by  $n = 97$  independent repetitions of this process. An obvious choice of the estimate  $\hat{P}$  in Figure 8 is  $(\hat{S}^\circ, \hat{R})$ , (5.14), (5.15). The rest of the bootstrap process is automatic:  $\hat{S}^\circ$  and  $\hat{R}$  replace  $S^\circ$  and  $R$  in (5.12), (5.13);  $n$  pairs  $(X_i^*, D_i^*)$  are independently generated according to rules (5.12)–(5.15), giving the bootstrap data set  $y^* = \{(x_1^*, d_1^*), \dots, (x_n^*, d_n^*)\}$ ; and finally the bootstrap Kaplan-Meier curve  $\hat{S}^{**}$  is constructed according to formula (5.16), and the bootstrap observed median  $\hat{\theta}^*$  gave estimated standard error  $\hat{\sigma} = 14.0$  months for  $\hat{\theta}$ . An estimated bias of 4.1 months was calculated as at (4.4). Efron (1981c) gives a fuller description.



**Figure 9.** Kaplan-Meier estimated survival curve for the Channing House men;  $t$  = age in months. The median survival age is estimated to be 1,044 months (87 years).

Once again there is a simpler way to apply the bootstrap. Consider each pair  $y_i = (x_i, d_i)$  as an observed point obtained by simple random sampling from a bivariate distribution  $F$ , and apply the bootstrap as described in Section 2 to the data set  $y_1, y_2, \dots, y_n \stackrel{\text{iid}}{\sim} F$ . This method makes no use of the special structure (5.12)–(5.15). Surprisingly, it gives *exactly the same answers* as the more complicated bootstrap method described earlier, Efron (1981a).

## 6. Examples with more complicated data structures.

### Example 1: Autoregressive Time Series Model

This example illustrates an application of the bootstrap to a famous time series.

The data are the Wolfer annual sunspot numbers for the years 1770-1889 (taken from Anderson 1976). Let the count for the  $i$ th year be  $z_i$ . After centering the data, (replacing  $z_i$  by  $z_i - \bar{z}$ ) we fit a first order autoregressive model

$$z_i = \phi z_{i-1} + \epsilon_i \quad (6.1)$$

where  $\epsilon_i \sim \text{i.i.d. } N(0, \sigma^2)$ . The estimate  $\hat{\phi}$  turned out to be .815 with an estimated standard error, one over the square root of the Fisher information, of .053.

A bootstrap estimate of the standard error of  $\hat{\phi}$  can be obtained as follows. Define the residuals  $\hat{\epsilon}_i = z_i - \hat{\phi} z_{i-1}$  for  $i = 2, 3, \dots, 120$ . A bootstrap sample  $z_1^*, z_2^* \dots z_{120}^*$  is created by sampling  $\hat{\epsilon}_2^*, \hat{\epsilon}_3^* \dots \hat{\epsilon}_{120}^*$  with replacement from the residuals, then letting  $z_1^* = z_1$ , and  $z_i^* = \hat{\phi} z_{i-1}^* + \hat{\epsilon}_i^*$ ,  $i = 2, \dots, 120$ . Finally, after centering the time series  $z_1^*, z_2^*, \dots, z_{120}^*$ ,  $\hat{\phi}^*$  is the estimate of the autoregressive parameter for this new time series. (We could, if we wished, sample the  $\hat{\epsilon}_i^*$  from a fitted normal distribution.)

A histogram of 1000 such bootstrap values  $\hat{\phi}_1^*, \hat{\phi}_2^*, \dots, \hat{\phi}_{1000}^*$  is shown in Figure 10.

The bootstrap estimate of standard error was .055, agreeing nicely with the usual formula. Note however that the distribution is skewed to the left, so a confidence interval for  $\phi$  might be asymmetric about  $\hat{\phi}$ , as discussed in Sections 8 and 9.

In bootstrapping the residuals, we have assumed that the first order auto-regressive model is correct. (Recall the discussion of regression models in Section 5). In fact, the first order autoregressive model is far from adequate for this data. A fit of second-order autoregressive model

$$z_i = \alpha z_{i-1} + \theta z_{i-2} + \epsilon_i \quad (6.2)$$

gave estimates  $\hat{\alpha} = 1.37$ ,  $\hat{\theta} = -.677$ , both with an estimated standard error of .067, based on Fisher information calculations. We applied the bootstrap to this model, producing the

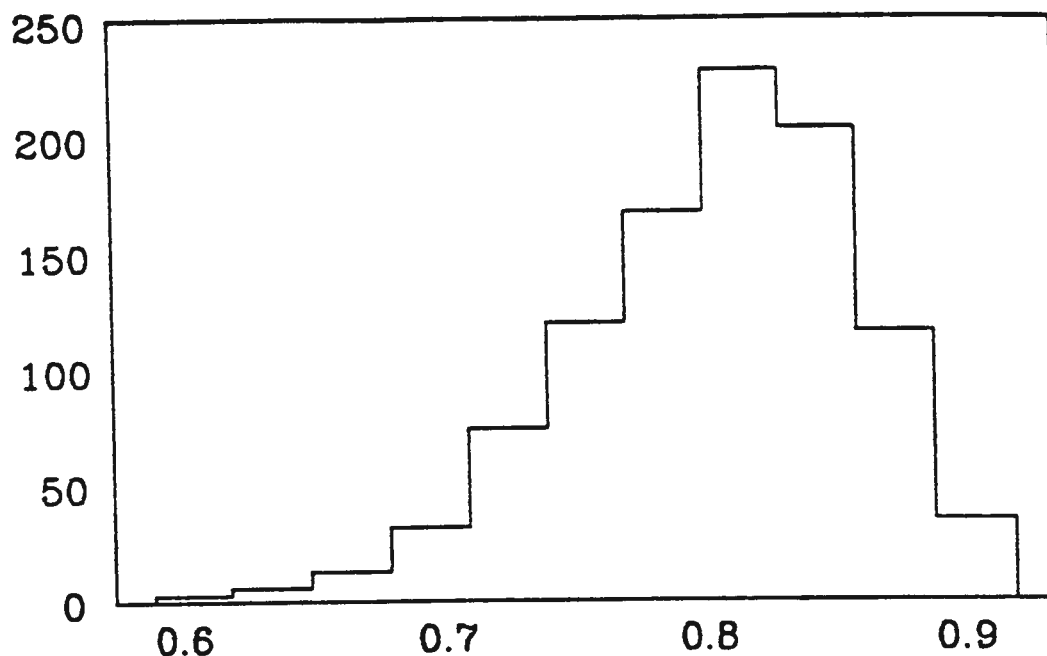


Figure 10. Bootstrap histogram of  $\hat{\phi}_1^*, \dots, \hat{\phi}_{1000}^*$  for the Wolfer sunspot data, model (6.1)

histograms for  $\alpha_1^*, \dots, \alpha_{1000}^*$  and  $\theta_1^*, \dots, \theta_{1000}^*$  shown in Figures 11 and 12 respectively.

The bootstrap standard errors were .070 and .068 respectively, both close to the usual value. Note that the additional term has reduced the skewness of the first coefficient.

### Example 2: Estimating a response transformation in regression

Box and Cox (1964) introduced a parametric family for estimating a transformation of the response in a regression. Given regression data  $\{(x_1, y_1), \dots, (x_n, y_n)\}$ , their model takes the form

$$z_i(\lambda) = x_i \cdot \beta + \epsilon_i \quad (6.3)$$

where  $z_i(\lambda) = (y_i^\lambda - 1)/\lambda$  for  $\lambda \neq 0$  and  $\log y_i$  for  $\lambda = 0$ , and  $\epsilon_i \sim \text{i.i.d } N(0, \sigma^2)$ . Estimates of

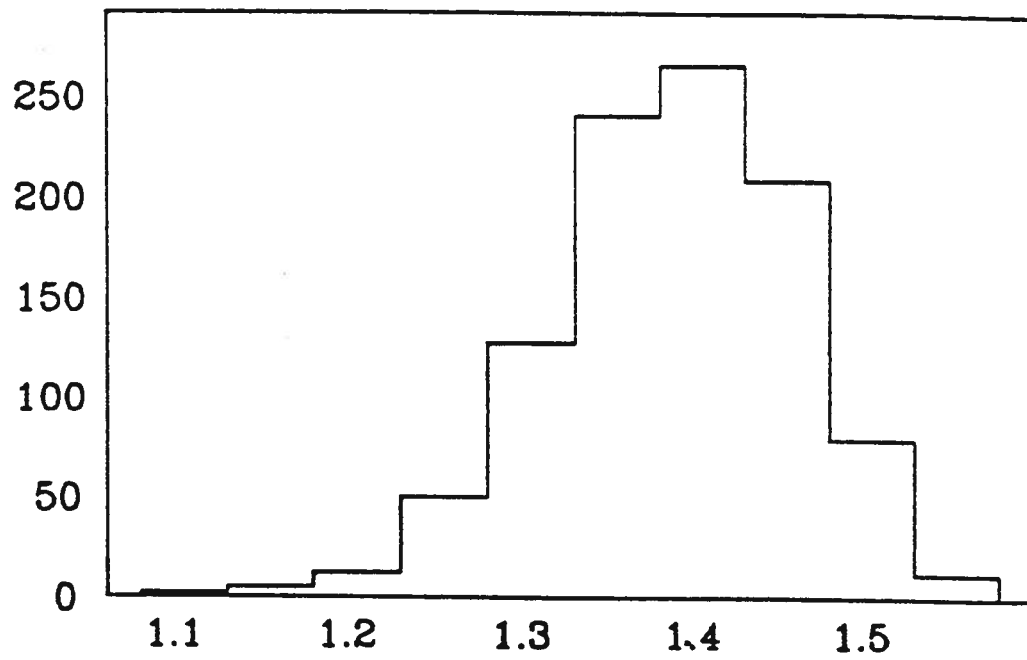


Figure 11. Bootstrap histogram of  $\hat{\alpha}^*, \dots, \hat{\alpha}_{1000}^*$  for the Wolfer sunspot data, model (6.2)

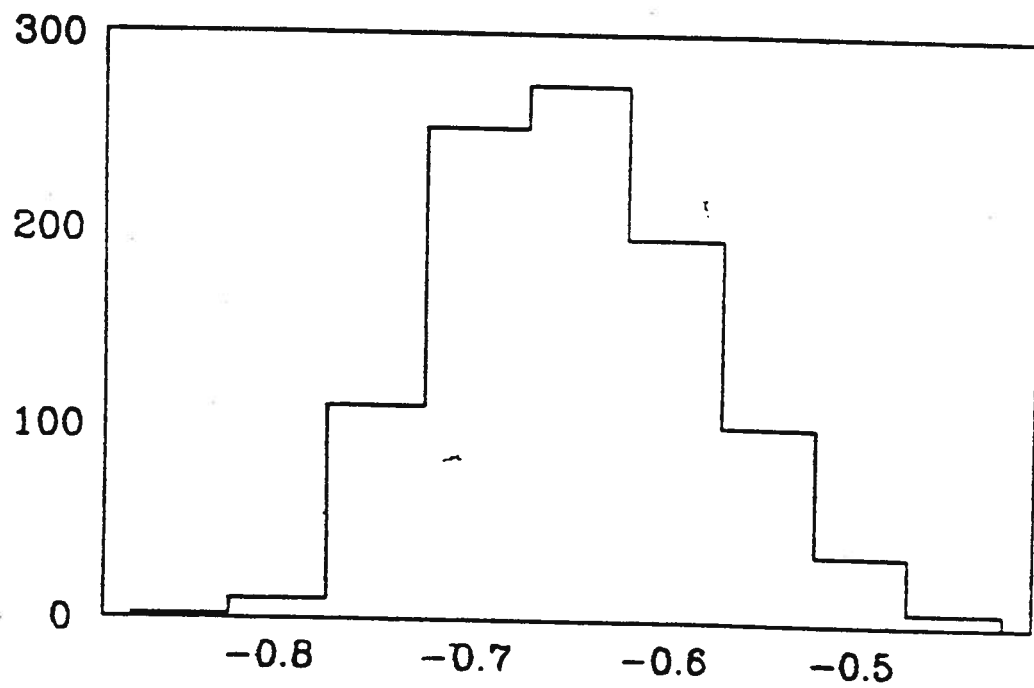


Figure 12. Bootstrap histogram of  $\hat{\theta}^*, \dots, \hat{\theta}_{1000}^*$  for the Wolfer sunspot data, model (6.2)

$\lambda$  and  $\beta$  are found by minimizing  $\sum_1^n (z_i - z_i \cdot \beta)^2$ .

Breiman and Friedman (1984) proposed a non-parametric solution for this problem. Their so-called ACE ("Alternating Conditional Expectation") model generalizes (6.3) to

$$s(y_i) = z_i \cdot \beta + \epsilon_i \quad (6.4)$$

where  $s(\cdot)$  is an unspecified smooth function. (In its most general form, ACE allows for transformations of the covariates as well). The function  $s(\cdot)$  and parameter  $\beta$  are estimated in an alternating fashion, utilizing a non-parametric smoother to estimate  $s(\cdot)$ .

In the following example, taken from Friedman and Tibshirani (1984), we compare the Box and Cox procedure to ACE and use the bootstrap to assess the variability of ACE.

The data, from Box and Cox (1964), consist of a 3x3x3 experiment on the strength of yarns, the response  $Y$  being number of cycles to failure, and the factors length of test specimen ( $X_1$ ) (250, 300 or 350 mm), amplitude of loading cycle ( $X_2$ ) (8, 9, or 10 mm), and load ( $X_3$ ) (40, 45 or 50 gm). As in Box and Cox, we treat the factors as quantitative and allow only a linear term for each. Box and Cox found that a logarithmic transformation was appropriate, with their procedure producing a value of -.06 for  $\hat{\lambda}$  with an estimated 95 percent confidence interval of (-.18,.06).

Figure 13 shows the transformation selected by the ACE algorithm. For comparison, the log function is plotted (normalized) on the same figure.

The similarity is truly remarkable! In order to assess the variability of the ACE curve, we can apply the bootstrap. Since the  $X$  matrix in this problem is fixed by design, we resampled from the residuals instead of from the  $(z_i, y_i)$  pairs. The bootstrap procedure was the following:



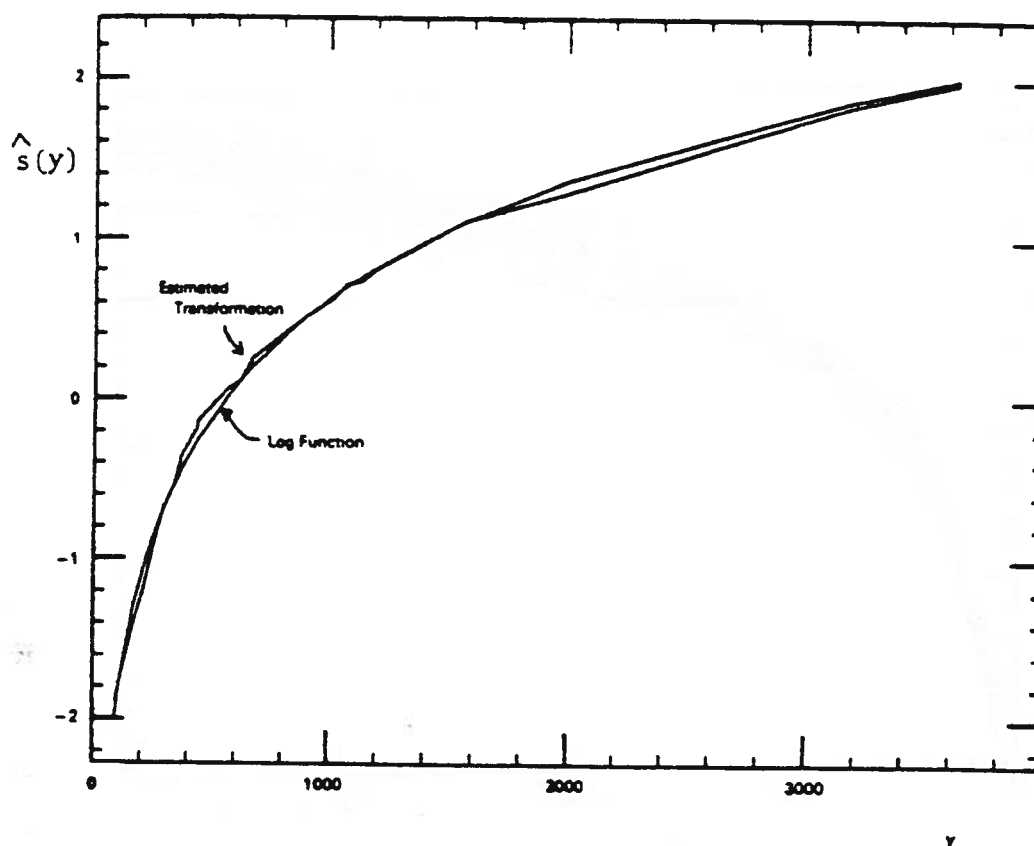


Figure 13. Estimated transformation from ACE and the log function, for Box and Cox example

Calculate residuals  $\hat{\epsilon}_i = \hat{s}(y_i) - \mathbf{z}_i \cdot \hat{\beta}$ ,  $i = 1, 2, \dots, n$

Repeat B times

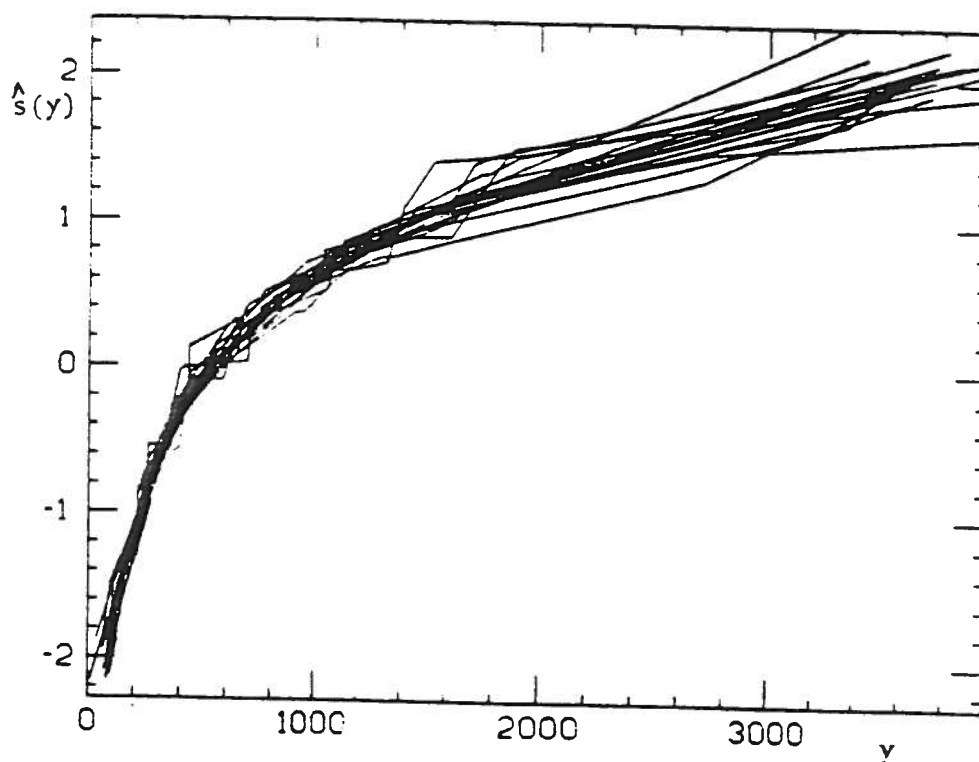
Choose a sample  $\tilde{\epsilon}_1^*, \dots, \tilde{\epsilon}_n^*$  with replacement from  $\hat{\epsilon}_1, \dots, \hat{\epsilon}_n$

Calculate  $y_i^* = \hat{s}^{-1}(\mathbf{z}_i \cdot \hat{\beta} + \tilde{\epsilon}_i^*)$ ,  $i = 1, 2, \dots, n$

Compute  $\hat{s}^*(\cdot) = \text{result of ACE algorithm applied to } (\mathbf{z}_1, y_1^*), \dots, (\mathbf{z}_n, y_n^*)$

End

The number of bootstrap replications B was 20. Note that the residuals are computed on the  $s(\cdot)$  scale, not the  $y$  scale, because it is on the  $s(\cdot)$  scale that the true residuals are assumed to be approximately i.i.d.. The 20 estimated transformations,  $\hat{s}_1^*(\cdot), \dots, \hat{s}_{20}^*(\cdot)$  are shown



**Figure 14.** Bootstrap replications of ACE transformations for Box and Cox example

in Figure 14.

The tight clustering of the smooths indicates that the original estimate  $\hat{s}(\cdot)$  has low variability, especially for smaller values of  $Y$ . This agrees qualitatively with the short confidence interval for  $\lambda$  in the Box and Cox analysis.

## 7. Bootstrap Confidence Intervals.

This section presents three closely related methods of using the bootstrap to set confidence intervals. The discussion is in terms of simple parametric models, where the logical basis of the bootstrap methods is easiest to see. Section 8 extends the methods to multiparameter and nonparametric models.

We have discussed obtaining  $\hat{\sigma}$ , the estimated standard error of an estimator  $\hat{\theta}$ . In practice,  $\hat{\theta}$  and  $\hat{\sigma}$  are usually used together to form the approximate confidence interval  $\theta \in \hat{\theta} \pm \hat{\sigma}z^{(\alpha)}$ , (1.7) is claimed to have approximate coverage probability  $1 - 2\alpha$ . For the law school example of Section 2, the values  $\hat{\theta} = .776$ ,  $\hat{\sigma} = .115$ ,  $z^{(.05)} = -1.645$ , give  $\theta \in [.587, .965]$  as an approximate 90% central interval for the true correlation coefficient.

We will call (1.7) the *standard interval for  $\theta$* . When working within parametric families like the bivariate normal,  $\hat{\sigma}$  in (1.7) is usually obtained by differentiating the log likelihood function, see Section 5a of Rao (1973), though in the context of this paper we might prefer to use the parametric bootstrap estimate of  $\sigma$ , e.g.  $\hat{\sigma}_{\text{NORM}}$  in Section 2.

The standard intervals are an immensely useful statistical tool. They have the great virtue of being automatic: a computer program can be written which produces (1.7) directly from the data  $y$  and the form of the density function for  $y$ , with no further input required from the statistician. Nevertheless the standard intervals can be quite inaccurate, as Table 5 shows. The standard interval (1.7), using  $\hat{\sigma}_{\text{NORM}}$ , (2.5), is strikingly different than the exact normal-theory interval based on the assumption of a bivariate normal sampling distribution  $F$ .

In this case it is well-known that it is better to make the transformation  $\hat{\phi} = \tanh^{-1}(\hat{\theta})$ ,  $\phi = \tanh^{-1}(\theta)$ , apply (1.7) on the  $\phi$  scale, and then transform back to the  $\theta$  scale. The resulting interval, line 3 of Table 7A, is moved closer to the exact interval. However, there is nothing automatic about the  $\tanh^{-1}$  transformation. For a different statistic than the correlation coefficient or a different distributional family than the bivariate normal, we might very well need other tricks to make (1.7) perform satisfactorily.

The bootstrap can be used to produce approximate confidence intervals in an automatic way. The following discussion is abridged from Efron (1984a and b) and Efron (1982, Chapter

|  |             |            |
|--|-------------|------------|
| 1. Exact (Normal Theory):                      | [.496,.898] | R/L = .44  |
| 2. Standard (1.7):                             | [.587,.965] | R/L = 1.00 |
| 3. Transformed Standard:                       | [.508,.907] | R/L = .49  |
| 4. Parametric Bootstrap (BC):                  | [.488,.900] | R/L = .43  |
| 5. Nonparametric Bootstrap (BC <sub>o</sub> ): | [.43,.92]   | R/L = .42  |

**Table 5.** Exact and approximate central 90% confidence intervals for  $\theta$ , the true correlation coefficient, from the law school data of Figure 1. R/L = ratio of right side of interval, measured from  $\hat{\theta} = .776$ , to left side. The exact interval is strikingly asymmetric about  $\hat{\theta}$ . Section 8 discusses the nonparametric method of line 5.

10). Line 4 of Table 5 shows that the parametric bootstrap interval for the correlation coefficient  $\theta$  is nearly identical to the exact interval. "Parametric" in this case means that the bootstrap algorithm begins from the bivariate normal MLE  $\hat{F}_{\text{NORM}}$ , as for the normal theory curve of Figure 2. This good performance is no accident. The bootstrap method used on line 4 in effect transforms  $\hat{\theta}$  to the best (most normal) scale. All of this is done automatically by the bootstrap algorithm, without requiring special intervention from the statistician. The price paid is a large amount of computing, perhaps  $B = 1000$  bootstrap replications, as discussed in Section 10.

Define  $\hat{G}(s)$  to be the parametric bootstrap c.d.f. of  $\hat{\theta}^*$ ,

$$\hat{G}(s) = \text{Prob.}\{\hat{\theta}^* < s\}, \quad (7.1)$$

where Prob. indicates probability computed according to the bootstrap distribution of  $\hat{\theta}^*$ . In Figure 2,  $\hat{G}(s)$  is obtained by integrating the normal theory curve. We will present three different kinds of bootstrap confidence intervals, in order of increasing generality. All three methods use percentiles of  $\hat{G}$  to define the confidence interval. They differ in which percentiles are used.

The simplest method is to take  $\theta \in [\hat{G}^{-1}(\alpha), \hat{G}^{-1}(1 - \alpha)]$  as an approximate  $1 - 2\alpha$  central interval for  $\theta$ . This is called the *percentile method* in Section 10.4 of Efron (1982). The percentile method interval is just the interval between the  $100 \cdot \alpha$  and  $100 \cdot (1 - \alpha)$  percentiles of the bootstrap distribution of  $\hat{\theta}^*$ .

We will use the notation of  $\theta[\alpha]$  for the  $\alpha$ -level endpoint of an approximate confidence interval for  $\theta$ , so  $\theta \in [\theta[\alpha], \theta[1 - \alpha]]$  is the central  $1 - 2\alpha$  interval. Subscripts will be used to indicate the various different methods. The percentile interval has the endpoints

$$\theta_p[\alpha] \equiv \hat{G}^{-1}(\alpha). \quad (7.2)$$

This compares with the standard interval,

$$\theta_S[\alpha] = \hat{\theta} + \hat{\sigma} z^{(\alpha)}. \quad (7.3)$$

Suppose the bootstrap c.d.f.  $\hat{G}$  is perfectly normal, say

$$\hat{G}(s) = \Phi\left(\frac{s - \hat{\theta}}{\hat{\sigma}}\right), \quad (7.4)$$

where  $\Phi(s) = \frac{1}{\sqrt{2\pi}} e^{-\frac{s^2}{2}} \int_{-\infty}^s (1\pi)^{-1/2} e^{-t^2/2} dt$ , the standard normal c.d.f. In other words, suppose that  $\hat{\theta}^*$  has bootstrap distribution  $N(\hat{\theta}, \hat{\sigma}^2)$ . In this case the standard method and the percentile method agree,  $\theta_S[\alpha] = \theta_p[\alpha]$ . In situations like that of Figure 2, where  $\hat{G}$  is markedly nonnormal, the standard interval is quite different from (7.2). Which is better?

To answer this question, consider the simplest possible situation, where for all  $\theta$

$$\hat{\theta} \sim N(\theta, \sigma^2). \quad (7.5)$$

That is, we have a single unknown parameter  $\theta$  with no nuisance parameters, and a single summary statistic  $\hat{\theta}$  normally distributed about  $\theta$  with constant standard error  $\sigma$ . In this case the parametric bootstrap c.d.f. is given by (7.4), so  $\theta_S[\alpha] = \theta_p[\alpha]$ . (The bootstrap estimate  $\hat{\sigma}$  equals  $\sigma$ .)

Suppose though that instead of (7.5) we have, for all  $\theta$ ,

$$\hat{\phi} \sim N(\phi, r^2), \quad (7.6)$$

for some monotone transformation  $\hat{\phi}g(\hat{\theta})$ ,  $\phi = g(\theta)$ , where  $r$  is a constant. In the correlation coefficient example the function  $g$  was  $\tanh^{-1}$ . The standard limits (7.2) can now be grossly inaccurate. However it is easy to verify that the percentile limits (7.2) are still correct. "Correct" here means that (7.2) is the mapping of the obvious interval for  $\phi$ ,  $\hat{\phi} \pm rz^{(\alpha)}$ , back to the

$\theta$  scale,  $\theta_P[\alpha] = g^{-1}(\hat{\phi} + rz^{(\alpha)})$ . It is also correct in the sense of having exactly the claimed average probability  $1 - 2\alpha$ .

Another way to state things is that the percentile intervals are transformation invariant,

$$\phi_P[\alpha] = g(\theta_P[\alpha]) \quad (7.7)$$

for any monotone transformation  $g$ . This implies that if the percentile intervals are correct on some transformed scale  $\phi = g(\theta)$ , then they must also be correct on the original scale  $\theta$ . The statistician doesn't need to know the normalizing transformation  $g$ , only that it exists. Definition (7.2) automatically takes care of the bookkeeping involved in the use of normalizing transformations for confidence intervals.

Fisher's theory of maximum likelihood estimation says that we are always in situation (7.5) to a first order of asymptotic approximation. However we are also in situation (7.6) for any choice of  $g$ , to the same order of approximation. Efron (1984a and b) uses higher order asymptotic theory to differentiate between the standard and bootstrap intervals. It is the higher order asymptotic terms which often make exact intervals strongly asymmetric about the MLE  $\hat{\theta}$ , as in Table 5. The bootstrap intervals are effective at capturing this asymmetry.

The percentile method automatically incorporates normalizing transformations, as in going from (7.5) to (7.6). It turns out that there are two other important ways that assumption (7.5) can be misleading, the first of which relates to possible bias in  $\hat{\theta}$ . For example consider  $f_{\theta}(\hat{\theta})$ , the family of densities for the observed correlation coefficient  $\hat{\theta}$  when sampling  $n = 15$  times from a bivariate normal distribution with true correlation  $\theta$ . In fact it is easy to see that no monotone mapping  $\hat{\phi} = g(\hat{\theta})$ ,  $\phi = g(\theta)$  transforms this family to  $\hat{\phi} \sim N(\phi, \tau^2)$ , as in (7.6). If there were such a  $g$ , then  $\text{Prob}_{\theta}\{\hat{\theta} < \theta\} = \text{Prob}_{\phi}\{\hat{\phi} < \phi\} = .50$ , but for  $\theta = .776$  integrating the density function  $f_{.776}(\hat{\theta})$  gives  $\text{Prob}_{\theta=.776}\{\hat{\theta} < \theta\} = .431$ .

The *bias-corrected percentile method* (BC method) makes an adjustment for this type of bias. Let

$$z_0 \equiv \Phi^{-1}\{\hat{G}(\hat{\theta})\}, \quad (7.8)$$

where  $\Phi^{-1}$  is the inverse function of the standard normal c.d.f. The BC method has  $\alpha$ -level

endpoint

$$\theta_{BC}[\alpha] \equiv \hat{G}^{-1}(\Phi\{2z_0 + z^{(\alpha)}\}). \quad (7.9)$$

Note: if  $\hat{G}(\hat{\theta}) = .50$ , that is if half of the bootstrap distribution of  $\hat{\theta}^*$  is less than the observed value  $\hat{\theta}$ , then  $z_0 = 0$  and  $\theta_{BC}[\alpha] = \theta_p[\alpha]$ . Otherwise definition (7.9) makes a bias correction.

Section 10.7 of Efron (1982) shows that the *BC* interval for  $\theta$  is exactly correct if

$$\hat{\phi} \sim N(\phi - z_0 r, r^2) \quad (7.10)$$

for some monotone transformation  $\hat{\phi} = g(\hat{\theta})$ ,  $\phi = g(\theta)$  and some constant  $z_0$ . It doesn't look like (7.10) is much more general than (7.6), but in fact the bias correction is often important.

In the example of Table 5, the percentile method (7.2) gives central 90% interval [.536, .911] compared to the *BC* interval [.488, .900]. By definition the endpoints .496 and .898 of the exact interval satisfy

$$\text{Prob}_{\theta=.496}\{\hat{\theta} > .776\} = .05 = \text{Prob}_{\theta=.898}\{\hat{\theta} < .776\}. \quad (7.11)$$

The corresponding quantities for the *BC* endpoints are

$$\text{Prob}_{\theta=.488}\{\hat{\theta} > .776\} = .0465, \quad \text{Prob}_{\theta=.900}\{\hat{\theta} < .776\} = .0475, \quad (7.12)$$

compared to

$$\text{Prob}_{\theta=.536}\{\hat{\theta} > .776\} = .0725, \quad \text{Prob}_{\theta=.911}\{\hat{\theta} < .776\} = .0293. \quad (7.13)$$

for the percentile endpoints. The bias correction is quite important in equalizing the error probabilities at the two endpoints. If  $z_0$  can be approximated accurately (as mentioned in Section 9), then it is preferable to use the *BC* intervals.

Table 6 shows a simple example where the *BC* method is less successful. The data consists of the single observation  $\hat{\theta} \sim \theta(\chi_{19}^2/19)$ , the notation indicating an unknown scale parameter  $\theta$ . In this case the *BC* interval based on  $\hat{\theta}$  is a definite improvement over the standard interval (1.7), but goes only about half as far as it should toward achieving the asymmetry of the exact interval.

It turns out that the parametric family  $\hat{\theta} \sim \theta(\chi_{19}^2/19)$  cannot be transformed into (7.10), not even approximately. The results of Efron (1982a) show that there does exist a monotone

|                         |  |              |
|-------------------------|--|--------------|
| 1. Exact                | $[\hat{\theta} \cdot .631, \hat{\theta} \cdot 1.88]$ | $R/L = 2.38$ |
| 2. Standard (1.7)       | $[\hat{\theta} \cdot .466, \hat{\theta} \cdot 1.53]$ | $R/L = 1.00$ |
| 3. $BC$ (7.9)           | $[\hat{\theta} \cdot .580, \hat{\theta} \cdot 1.69]$ | $R/L = 1.64$ |
| 4. $BC_*$ (7.15)        | $[\hat{\theta} \cdot .630, \hat{\theta} \cdot 1.88]$ | $R/L = 2.37$ |
| 5. Nonparametric $BC_*$ | $[\hat{\theta} \cdot .640, \hat{\theta} \cdot 1.68]$ | $R/L = 1.88$ |

**Table 6.** Central 90% confidence intervals for  $\theta$ , having observed  $\hat{\theta} \sim \theta(\chi_{19}^2/19)$ . The exact interval is sharply skewed to the right of  $\hat{\theta}$ . The  $BC$  method is only a partial improvement over the standard interval. The  $BC_*$  interval,  $a = .108$ , agrees almost perfectly with the exact interval.

transformation  $g$  such that  $\hat{\phi} = g(\hat{\theta})$ ,  $\phi = g(\theta)$  satisfy to a high degree of approximation

$$\hat{\phi} \sim N(\phi - z_0 r_\phi, r_\phi^2) \quad (r_\phi = 1 + a\phi). \quad (7.14)$$

The constants in (7.14) are  $z_0 = .1082$ ,  $a = .1077$ .

The  $BC_*$  method, Efron (1984b), is a method of assigning bootstrap confidence intervals which are exactly right for problems which can be mapped into form (7.14). This method has  $\alpha$ -level endpoint

$$\theta_{BC_*}[\alpha] \equiv \hat{G}^{-1} \left( \Phi \left\{ z_0 + \frac{z_0 + z^{(\alpha)}}{1 - a(z_0 + z^{(\alpha)})} \right\} \right). \quad (7.15)$$

If  $a = 0$  then  $\theta_{BC_*}[\alpha] = \theta_{BC}[\alpha]$ , but otherwise the  $BC_*$  intervals can be a substantial improvement over the  $BC$  method, as shown in Table 7B.

The constant  $z_0$  in (7.15) is given by  $z_0 = \Phi^{-1}\{\hat{G}(\hat{\theta})\}$ , (7.8), and so can be computed directly from the bootstrap distribution. How do we know  $a$ ? It turns out that in one-parameter families  $f_\theta(\hat{\theta})$ , a good approximation is

$$a \doteq \frac{\text{SKEW}_{\theta=\hat{\theta}}(\dot{\ell}_\theta(t))}{6}, \quad (7.16)$$

where  $\text{SKEW}_{\theta=\hat{\theta}}(\dot{\ell}_\theta(t))$  is the skewness at parameter value  $\theta = \hat{\theta}$  of the score statistic  $\dot{\ell}_\theta(t) = \frac{\partial}{\partial \theta} \log f_\theta(t)$ . For  $\hat{\theta} \sim \theta(\chi_{19}^2/19)$  this gives  $a \doteq .1081$ , compared to the actual value  $a = .1077$  derived in Efron (1984b). For the normal theory correlation of Table 5  $a \doteq 0$  which explains why the  $BC$  method, which takes  $a = 0$ , works so well there.



The advantage of formula (7.18) is that we needn't know the transformation  $g$  leading to (7.14) in order to approximate  $\alpha$ . In fact  $\theta_{BC_\alpha}[\alpha]$ , like  $\theta_{BC_\alpha}[\alpha]$  and  $\theta_P[\alpha]$ , is transformation invariant, as in (7.7). Like the bootstrap methods, the  $BC_\alpha$  intervals are computed directly from the form of the density function  $f_\theta(\cdot)$ , for  $\theta$  near  $\hat{\theta}$ .

Formula (7.16) applies to the case where  $\theta$  is the only parameter. Section 8 briefly discusses the more challenging problem of setting confidence intervals for a parameter  $\theta$  in a multiparameter family, and also in nonparametric situations where the number of nuisance parameters is effectively infinite.

To summarize this section, the progression from the standard intervals to the  $BC_\alpha$  method is based on a series of increasingly less restrictive assumptions, (7.5), (7.6), (7.10), and finally (7.14). Each step requires the statistician to do a greater amount of computation, first the bootstrap distribution  $\hat{G}$ , then the bias-correction constant  $z_0$ , and finally the constant  $\alpha$ . However all of these computations are algorithms in character, and can be carried out in an automatic fashion.

Chapter 10 of Efron (1982) discusses several other ways of using the bootstrap to construct approximate confidence intervals, which will not be presented here. One of these methods, the "bootstrap  $t$ ", was used in the blood serum example of Section 4.

## 8. Nonparametric and Multiparameter Confidence Intervals.

Section 7 focused on the simple case  $\hat{\theta} \sim f_\theta$ , where we have only a real-valued parameter  $\theta$  and a real-valued summary statistic  $\hat{\theta}$  from which we are trying to construct a confidence interval for  $\theta$ . Various favorable properties of the bootstrap confidence intervals were demonstrated in the simple case, but of course the simple case is where we least need a general method like the bootstrap.

Now we will discuss the more common situation where there are nuisance parameters besides the parameter of interest  $\theta$ ; or even more generally the nonparametric case, where the number of nuisance parameters is effectively infinite. The discussion is limited to a few brief examples. Efron (1984a and b) develops the theoretical basis of bootstrap confidence intervals for complicated situations, and gives many more examples.

|                                   | for $\theta$        | for $\phi$       |
|-----------------------------------|---------------------|------------------|
| 1. Exact (Fieller):               | [.29,.76]           | [1.32,3.50]      |
| 2. Parametric Boot ( <i>BC</i> ): | [.29,.76]           | [1.32,3.50]      |
| 3. Standard (1.7):                | [.27,.73]           | [1.08,2.92]      |
| MLE                               | $\hat{\theta} = .5$ | $\hat{\phi} = 2$ |

**Table 7.** Central 90% confidence intervals for  $\theta = \eta_2/\eta_1$  and for  $\phi = 1/\theta$ , having observed  $(y_1, y_2) = (8, 4)$  from a bivariate normal distribution  $\mathbf{y} \sim N_2(\boldsymbol{\eta}, \mathbf{I})$ . The *BC* intervals, line 2, are based on the parametric bootstrap distribution of  $\hat{\theta} = y_2/y_1$ .

### Example 1: Ratio Estimation

The data consists of  $\mathbf{y} = (y_1, y_2)$ , assumed to come from a bivariate normal distribution with unknown mean vector  $\boldsymbol{\eta}$  and covariance matrix the identity,

$$\mathbf{y} \sim N_2(\boldsymbol{\eta}, \mathbf{I}). \quad (8.1)$$

The parameter of interest, for which we desire a confidence interval, is the ratio

$$\theta = \eta_2/\eta_1. \quad (8.2)$$

Fieller (1954) provided well-known exact intervals for  $\theta$  having observed  $\mathbf{y} = (8, 4)$ . Also shown is the Fieller interval for  $\phi = 1/\theta = \eta_1/\eta_2$ , which equals  $[\cdot76^{-1}, \cdot29^{-1}]$ , the obvious transformation of the interval for  $\theta$ . The standard interval (1.7) is satisfactory for  $\theta$ , but not for  $\phi$ . Notice that the standard interval does not transform correctly from  $\theta$  to  $\phi$ .

Line 2 shows the *BC* intervals based on applying definitions (7.8), (7.9) to the parametric bootstrap distribution of  $\hat{\theta} = y_2/y_1$  (or  $\hat{\phi} = y_1/y_2$ ). This is the distribution of  $\hat{\theta}^* = y_2^*/y_1^*$  when sampling  $\mathbf{y}^* = (y_1^*, y_2^*)$  from  $\hat{F}_{\text{NORM}} \sim N_2((y_1, y_2), \mathbf{I})$ . The bootstrap intervals transform correctly, and in this case they agree with the exact interval to three decimal places.

### Example 2: Product of Normal Means

For most multiparameter situations, there do not exist exact confidence intervals for a

|                          | for $\theta$            | for $\phi$        |
|--------------------------|-------------------------|-------------------|
| 1. Almost Exact:         | [1.77, 17.03]           | [3.1, 290.0]      |
| 2. Parametric Boot (BC): | [1.77, 17.12]           | [3.1, 293.1]      |
| 3. Standard (1.7):       | [0.64, 15.36]           | [-53.7, 181.7]    |
| MLE                      | $\hat{\theta} \doteq 8$ | $\hat{\phi} = 64$ |

Table 8. Central 90% confidence intervals for  $\theta = \eta_1 \eta_2$  and  $\phi = \theta^2$  having observed  $\mathbf{y} = (2, 4)$ , where  $\mathbf{y} \sim N_2(\boldsymbol{\eta}, I)$ . The almost exact intervals are based on the high order approximation theory of Efron (1984a). The BC intervals of line 2 are based on the parametric bootstrap distribution of  $\hat{\theta} = y_1 y_2$ .

single parameter of interest. Suppose for instance that (8.2) is changed to

$$\theta = \eta_1 \eta_2 \quad (8.3)$$

still assuming (8.1). Table 8 shows approximate intervals for  $\theta$ , and also for  $\phi = \theta^2$ , having observed  $\mathbf{y} = (2, 4)$ . The "almost exact" intervals are based on an analogue of Fieller's argument, Efron (1984a), which with suitable care can be carried through to a high degree of accuracy. Once again, the parametric BC intervals are a close match to line 1. The fact that the standard intervals do not transform correctly is particularly obvious here.

The good performance of the parametric BC intervals is not accidental. The theory developed in Efron (1984a) shows that the BC intervals, based on bootstrapping the MLE  $\hat{\theta}$ , agree to high order with the almost exact intervals in the following class of problems: the data  $\mathbf{y}$  comes from a multiparameter family of densities  $f_{\boldsymbol{\eta}}(\mathbf{y})$ , both  $\mathbf{y}$  and  $\boldsymbol{\eta}$   $k$ -dimensional vectors; the real-valued parameter of interest  $\theta$  is a smooth function of  $\boldsymbol{\eta}$ ,  $\theta = t(\boldsymbol{\eta})$ ; and the family  $f_{\boldsymbol{\eta}}(\mathbf{y})$  can be transformed to multivariate normality, say

$$\mathbf{g}(\mathbf{y}) \sim N_k(h(\boldsymbol{\eta}), I), \quad (8.4)$$

by some one-to-one transformations  $\mathbf{g}$  and  $h$ .

Just as in Section 7, it is not necessary for the statistician to know the normalizing transformations  $\mathbf{g}$  and  $h$ , only that they exist. The BC intervals are obtained directly from the original densities  $f_{\boldsymbol{\eta}}$ : we find  $\hat{\boldsymbol{\eta}} = \hat{\boldsymbol{\eta}}(\mathbf{y})$ , the MLE of  $\boldsymbol{\eta}$ ; sample  $\mathbf{y}^* \sim f_{\hat{\boldsymbol{\eta}}}$ ; compute  $\hat{\theta}^*$ , the

bootstrap MLE of  $\theta$ ; calculate  $\hat{G}$ , the bootstrap c.d.f. of  $\hat{\theta}^*$ , usually by Monte Carlo sampling, and finally apply definitions (7.8), (7.9). This process gives the same interval for  $\theta$  whether or not the transformation to form (8.4) has been made.

Not all problems can be transformed as in (8.4) to a normal distribution with constant covariance. The case considered in Table 6 is a one-dimensional counterexample. As a result the  $BC$  intervals do not always work as well as in Tables 7 and 8, though they usually improve on the standard method. However in order to take advantage of the  $BC_a$  method, which is based on more general assumptions, we need to be able to calculate the constant  $a$ .

Efron (1984b) gives expressions for " $a$ " generalizing (7.16) to multiparameter families, and also to nonparametric situations. If (8.4) holds, then " $a$ " will have value zero, and the  $BC_a$  method reduces to the  $BC$  case. Otherwise the two intervals differ.

Here we will discuss only the nonparametric situation: the observed data  $y = (x_1, x_2, \dots, x_n)$  consists of i.i.d. observations  $X_1, X_2, \dots, X_n \sim F$ , where  $F$  can be any distribution on the sample space  $X$ ; we want a confidence interval for  $\theta = t(F)$ , some real-valued functional of  $F$ ; and the bootstrap interval are based on bootstrapping  $\hat{\theta} = t(\hat{F})$ , which is the nonparametric MLE of  $\theta$ . In this case a good approximation to the constant  $a$  is given in terms of the empirical influence function  $U_i^*$ , defined in Section 10 at (10.11),

$$a \doteq \frac{1}{6} \frac{\sum_{i=1}^n (U_i^*)^3}{\{\sum_{i=1}^n (U_i^*)^2\}^{3/2}}. \quad (8.5)$$

This is a convenient formula, since it is easy to numerically evaluate the  $U_i^*$  by simply substituting a small value of  $\theta$  into (10.11).

### Example 3: The Law School Data

For  $\hat{\theta}$  the correlation coefficient, the values of  $U_i^*$  corresponding to the 15 data points shown in Figure 1 are -1.507, .168, .273, .004, .525, -.049, -.100, .477, .310, .004, -.526, -.091, .323, .125, -.048. (Notice how influential law school 1 is.) Formula (8.5) gives  $a \doteq -.0817$ .  $B = 100,000$  bootstrap replications, about 100 times more than was actually necessary, see Section 10, gave  $z_0 = -.0927$ , and the central 90% interval  $\theta \in [.43, .92]$  shown in Table 7. The nonparametric  $BC_a$  interval is quite reasonable in this example, particularly considering

that there is no guarantee that the true law school distribution  $F$  is anywhere near a bivariate normal.

#### Example 4: Mouse Leukemia Data (the first example in Section 3)

The standard central 90% interval for  $\beta$  in formula (3.1) is  $[\cdot 835, 2.18]$ . The bias-correction constant  $z_0 \doteq .0275$ , giving  $BC$  interval  $[1.00, 2.39]$ . This is shifted far right of the standard interval, reflecting the long right tail of the bootstrap histogram seen in Figure 3. We can calculate " $a$ " from (8.5), considering each of the  $n = 42$  data points to be a triple  $(y_i, x_i, \delta_i)$ :  $a \doteq -.152$ . Because  $a$  is negative, the  $BC_a$  interval is shifted back to the left, equaling  $[\cdot 788, 2.10]$ . This contrasts with the law school example, where  $a$ ,  $z_0$ , and the skewness of the bootstrap distribution added to each other rather than cancelling out, resulting in a  $BC_a$  interval much different than the standard normal.

Efron (1984b) provides some theoretical support for the nonparametric  $BC_a$  method. However the problem of setting approximate nonparametric confidence intervals is still far from well understood, and all methods should be interpreted with some caution. We end this section with a cautionary example.

#### Example 5: The Variance

Suppose  $X$  is the real line, and  $\theta = \text{Var}_F X$ , the variance. Line 5 of Table 2 shows the result of applying the nonparametric  $BC_a$  method to data sets  $x_1, x_2, \dots, x_{20}$  which were actually i.i.d. samples from a  $N(0, 1)$  distribution. The number .640 for example is the average of  $\theta_{BC_a}[\cdot 05]/\hat{\theta}$  over 40 such data sets,  $B = 4000$  bootstrap replications per data set. The upper limit  $1.68 \cdot \hat{\theta}$  is noticeably small, as pointed out by Schenker (1983). The reason is simple: the nonparametric bootstrap distribution which is a scaled  $\chi^2_1$  random variable. The results of Beran (1984), Bickel and Friedman (1981), and Singh (1981) show that the nonparametric bootstrap distribution is highly accurate asymptotically, but of course that isn't a guarantee of good small-sample behavior. Bootstrapping from a smoothed version of  $\hat{F}$ , as in lines 3, 4, and 5 of Table 2 alleviates the problem in this particular example.

## 9. Bootstrap Sample Sizes.

How many bootstrap replications must we take? Consider the standard error estimate  $\hat{\sigma}_B$  based on  $B$  bootstrap replications, (2.4). As  $B \rightarrow \infty$ ,  $\hat{\sigma}_B$  approaches  $\hat{\sigma}$ , the bootstrap estimate of standard error as originally defined in (2.3). Because  $\hat{F}$  does not estimate  $F$  perfectly,  $\hat{\sigma} = \sigma(\hat{F})$  will have a non-zero coefficient of variation ( $CV$ ) for estimating the true standard error  $\sigma = \sigma(F)$ ;  $\hat{\sigma}_B$  will have a larger  $CV$  because of the randomness added by the Monte Carlo bootstrap sampling.

It is easy to derive the following approximation,

$$CV(\hat{\sigma}_B) \doteq \left\{ CV(\hat{\sigma})^2 + \frac{E\{\hat{\delta}\} + 2}{4B} \right\}^{1/2}, \quad (9.1)$$

where  $\hat{\delta}$  is the kurtosis of the bootstrap distribution of  $\hat{\sigma}^*$ , given the data  $y$ , and  $E\{\hat{\delta}\}$  its expected value averaged over  $y$ . For typical situations,  $CV(\hat{\sigma})$  lies between .10 and .30. For example if  $\hat{\theta} = \bar{x}$ ,  $n = 20$ ,  $x_i \stackrel{\text{iid}}{\sim} N(0, 1)$ , then  $CV(\hat{\sigma}) \doteq .16$ .

Table 9 shows  $CV(\hat{\sigma}_B)$  for various values of  $B$  and  $CV(\hat{\sigma})$ , assuming  $E\{\hat{\delta}\} = 0$  in (9.1). For values of  $CV(\hat{\sigma}) > .10$ , there is little improvement past  $B = 100$ . In fact  $B$  as small as 25 gives reasonable results. Even smaller values of  $B$  can be quite informative, as we saw in the Stanford Heart Transplant Data, Figure of Section 3.

|                    | $B \rightarrow$ |     |     |     |     |          |
|--------------------|-----------------|-----|-----|-----|-----|----------|
|                    |                 | 25  | 50  | 100 | 200 | $\infty$ |
| $CV(\hat{\sigma})$ | .25             | .29 | .27 | .26 | .25 | .25      |
| $\downarrow$       | .20             | .24 | .22 | .21 | .21 | .20      |
|                    | .15             | .21 | .18 | .17 | .16 | .15      |
|                    | .05             | .15 | .11 | .09 | .07 | .05      |
|                    | 0               | .14 | .10 | .07 | .05 | 0        |

Table 9. Coefficient of variation of  $\hat{\sigma}_B$ , the bootstrap estimate of standard error based on  $B$  Monte Carlo replications, as a function of  $B$  and  $CV(\hat{\sigma})$ , the limiting  $CV$  as  $B \rightarrow \infty$ . Based on (9.1), assuming  $E\{\hat{\delta}\} = 0$ .

The situation is quite different for setting bootstrap confidence intervals. The calculations of Efron (1984b), Section 8, show that  $B = 1000$  is a rough minimum for the number of Monte Carlo bootstraps necessary to compute the  $BC$  or  $BC_0$  intervals. Somewhat smaller values, say  $B = 250$ , can give a useful percentile interval, the difference being that then the constant  $z_0$  need not be computed. Confidence intervals are a fundamentally more ambitious measure of statistical accuracy than standard errors, so it is not surprising that they require more computational effort.

## 10. The Jackknife and the Delta Method.

This section returns to the simple case of assigning a standard error to  $\hat{\theta}(y)$ , where  $y = (x_1, \dots, x_n)$  is obtained by random sampling from a single unknown distribution,  $x_1, \dots, x_n \stackrel{\text{iid}}{\sim} F$ . We will give another description of the bootstrap estimate  $\hat{\sigma}$ , which illustrates the bootstrap's relationship to older techniques of assigning standard errors, like the jackknife and the delta method.

For a given bootstrap sample  $y^* = (x_1^*, \dots, x_n^*)$ , as described in step (i) of the algorithm in Section 2, let  $p_i^*$  indicate the proportion of the bootstrap sample equal to  $x_i$ ,

$$p_i^* = \frac{\#\{x_j^* = x_i\}}{n} \quad i = 1, 2, \dots, n, \quad (10.1)$$

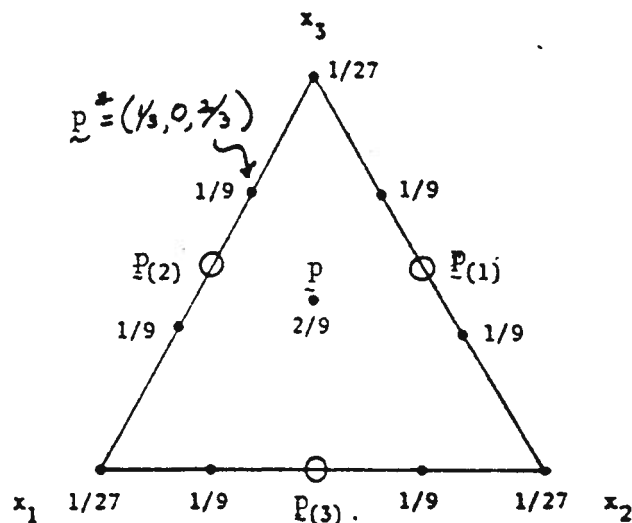
$p^* = (p_1^*, p_2^*, \dots, p_n^*)$ . The vector  $p^*$  has a rescaled multinomial distribution

$$p^* \sim \text{Mult}_n(n, p^*)/n \quad (p^* = (1/n, 1/n, \dots, 1/n)), \quad (10.2)$$

where the notation indicates the proportions observed from  $n$  random draws on  $n$  categories, each with probability  $1/n$ .

For  $n = 3$  there are 10 possible bootstrap vectors  $p^*$ . These are indicated in Figure 15 along with their multinomial probabilities, from (10.2). For example,  $p^* = (1/3, 0, 2/3)$ , corresponding to  $x^* = (x_1, x_3, x_3)$  or any permutation of these values, has bootstrap probability  $1/9$ .

To make our discussion easier, suppose that the statistic of interest  $\hat{\theta}$  is of functional form:  $\hat{\theta} = \theta(\hat{F})$ , where  $\theta(F)$  is a functional assigning a real number to any distribution  $F$



**Figure 15.** The bootstrap and jackknife sampling points in the case  $n = 3$ . The bootstrap points ( $\cdot$ ) are shown with their probabilities.

on the sample space  $\mathcal{X}$ . The mean, the correlation coefficient, and the trimmed mean are all of functional form. Statistics of functional form have the same value as a function of  $\hat{F}$ , no matter what the sample size  $n$  may be, which is convenient for discussing the jackknife and delta method.

For any vector  $p = (p_1, p_2, \dots, p_n)$  having non-negative weights summing to 1, define the weighted empirical distribution

$$\hat{F}(p) : \text{probability } p_i \text{ on } x_i; \quad i = 1, \dots, n. \quad (10.3)$$

For  $p = p^* = 1/n$ , the weighted empirical distribution equals  $\hat{F}$ , (1.4).

Corresponding to  $p$  is a resampled value of  $\hat{\theta}$ ,

$$\hat{\theta}(p) \equiv \theta(\hat{F}(p)). \quad (10.4)$$

The shortened notation  $\hat{\theta}(p)$  assumes that the data  $(x_1, x_2, \dots, x_n)$  is considered fixed. Notice



that  $\hat{\theta}(\mathbf{p}^0) = \theta(\hat{F})$  is the observed value of the statistic of interest. The bootstrap estimate  $\hat{\sigma}$ , (2.3), can then be written

$$\hat{\sigma} = [\text{Var. } \hat{\theta}(\mathbf{p}^*)]^{1/2}, \quad (10.5)$$

where  $\text{Var.}$  indicates variance with respect to distribution (10.2). In terms of Figure 15,  $\hat{\sigma}$  is the standard deviation of the ten possible bootstrap values  $\hat{\theta}(\mathbf{p}^*)$ , weighted as shown.

It looks like we could always calculate  $\hat{\sigma}$  simply by doing a finite sum. Unfortunately the number of bootstrap points is  $\binom{2n-1}{n}$ , 77,558,710 for  $n = 15$  so straightforward calculation of  $\hat{\sigma}$  is usually impractical. That is why we have emphasized Monte Carlo approximations to  $\hat{\sigma}$ . Therneau (1983) considers the question of methods more efficient than pure Monte Carlo, but at present there is no generally better method available.

However there is another approach to approximating (10.5): we can replace the usually complicated function  $\hat{\theta}(\mathbf{p})$  by an approximation linear in  $\mathbf{p}$ , and then use the well-known formula for the multinomial variance of a linear function. The *jackknife approximation*  $\hat{\theta}_J(\mathbf{p})$  is the linear function of  $\mathbf{p}$  which matches  $\hat{\theta}(\mathbf{p})$ , (10.4), at the  $n$  points corresponding to the deletion of a single  $x_i$  from the observed data set  $x_1, x_2, \dots, x_n$ ,

$$\mathbf{p}_{(i)} = \frac{1}{n-1}(1, 1, \dots, 1, 0, 1, \dots, 1) \quad (10.6)$$

$i = 1, 2, \dots, n$ . Figure 7A indicates the jackknife points for  $n = 3$ ; because  $\hat{\theta}$  is the functional form, (10.4), it doesn't matter that the jackknife points correspond to sample size  $n-1$  rather than  $n$ .

The linear function  $\hat{\theta}_J(\mathbf{p})$  is calculated to be

$$\hat{\theta}_J(\mathbf{p}) = \hat{\theta}_{(i)} + (\mathbf{p} - \mathbf{p}^0) \cdot \mathbf{U} \quad (10.7)$$

where in terms of  $\hat{\theta}_{(i)} \equiv \hat{\theta}(\mathbf{p}_{(i)})$ ,  $\hat{\theta}_{(\cdot)} = \sum_{i=1}^n \hat{\theta}_{(i)} / n$ , and  $\mathbf{U}$  is the vector with  $i$ th coordinate

$$U_i = (n-1)(\hat{\theta}_{(\cdot)} - \hat{\theta}_{(i)}). \quad (10.8)$$

The jackknife estimate of standard error, Tukey (1958), Miller (1974), is

$$\hat{\sigma}_J \equiv \left[ \frac{n-1}{n} \sum_{i=1}^n \{ \hat{\theta}_{(i)} - \hat{\theta}_{(\cdot)} \}^2 \right]^{1/2} = \left[ \frac{\sum_{i=1}^n U_i^2}{n(n-1)} \right]^{1/2}. \quad (10.9)$$

A standard multinomial calculation gives the following theorem (Efron 1982),

**Theorem.** The jackknife estimate of standard error equals  $[n/(n-1)]^{1/2}$  times the bootstrap estimate of standard error for  $\hat{\theta}_J$ ,

$$\hat{\sigma}_J = \left[ \frac{n}{n-1} \text{Var. } \hat{\theta}_J(p^*) \right]^{1/2}. \quad (10.10)$$

In other words, the jackknife estimate is itself almost a bootstrap estimate applied to a linear approximation of  $\hat{\theta}$ . The factor  $[n/(n-1)]^{1/2}$  in (10.10) makes  $\hat{\sigma}_J^2$  unbiased for  $\sigma^2$  in the case where  $\hat{\theta} = \bar{x}$ , the sample mean. We could multiply the bootstrap estimate  $\hat{\sigma}$  by this same factor, and achieve the same unbiasedness, but there doesn't seem to be any consistent advantage to doing so. The jackknife requires  $n$ , rather than  $B = 50$  to 200 resamples, at the expense of adding a linear approximation to the standard error estimate. Tables 1 and 2 indicate that there is some estimating efficiency lost in making this approximation. For statistic like the sample median which are difficult to approximate linearly, the jackknife is useless, see Section 3.4 of Efron (1982).

There is a more obvious linear approximation to  $\hat{\theta}(p)$  than  $\hat{\theta}_J(p)$ . Why not use the first-order Taylor series expansion for  $\hat{\theta}(p)$  about the point  $p = p^*$ ? This is the idea of Jaeckel's *infinitesimal jackknife* (1972). The Taylor series approximation turns out to be

$$\hat{\theta}_T(p) = \hat{\theta}(p^*) + (p - p^*)' U^*$$

where

$$U_i^* = \lim_{\epsilon \rightarrow 0} \frac{\hat{\theta}((1-\epsilon)p^* + \epsilon \delta_i) - \hat{\theta}(p^*)}{\epsilon}, \quad (10.11)$$

$\delta_i$  being the  $i$ th coordinate vector. This suggests the infinitesimal jackknife estimate of standard error

$$\hat{\sigma}_{IJ} \equiv [\text{Var. } \hat{\theta}_T(p^*)]^{1/2} = [\Sigma U_i^{*2}/n^2]^{1/2} \quad (10.12)$$

with  $\text{Var.}$  still indicating variance under (10.2). The ordinary jackknife can be thought of as taking  $\epsilon = -1/(n-1)$  in the definition of  $U_i^*$ , while the infinitesimal jackknife lets  $\epsilon \rightarrow 0$ , thereby earning the name.

The  $U_i^*$  are values of what Mallows (1974) calls the empirical influence function. Their

definition is a nonparametric estimate of the true influence function

$$IF(z) = \lim_{\epsilon \rightarrow 0} \frac{\theta((1-\epsilon)F + \epsilon\delta_z) - \theta(F)}{\epsilon},$$

$\delta_z$  being the degenerate distribution putting mass 1 on  $z$ . The right side of (10.12) is then the obvious estimate of the influence function approximation to the standard error of  $\hat{\theta}$ , (Hampel 1974),  $\sigma(F) \doteq [\int IF^2(z)dF(z)/n]^{1/2}$ . The empirical influence function method and the infinitesimal jackknife give identical estimates of standard error.

How have statisticians gotten along for so many years without methods like the jackknife and bootstrap? The answer is the delta method, which is still the most commonly used device for approximating standard errors. The method applies to statistics of the form  $t(\bar{Q}_1, \bar{Q}_2, \dots, \bar{Q}_A)$ , where  $t(\cdot, \dots, \cdot)$  is a known function and each  $\bar{Q}_a$  is an observed average,  $\bar{Q}_a = \sum_{i=1}^n Q_a(X_i)/n$ . For example the correlation  $\hat{\theta}$  is a function of  $A = 5$  such averages; the average of the first coordinate values, the second coordinates, the first coordinates squared, the second coordinates squares, and the cross-products.

In its nonparametric formulation, the delta method works by (a) expanding  $t$  in a linear Taylor series using the usual expressions for variances and covariances of averages; and (b) substituting  $\gamma(\hat{F})$  for any unknown quantity  $\gamma(F)$  occurring in (c). For example, the nonparametric delta method estimates the standard error of the correlation  $\hat{\theta}$  by

$$\left\{ \frac{\hat{\theta}^2}{4n} \left[ \frac{\hat{\mu}_{40}}{\hat{\mu}_{40}^2} + \frac{\hat{\mu}_{40}}{\hat{\mu}_{02}^2} + \frac{2\hat{\mu}_{22}}{\hat{\mu}_{20}\hat{\mu}_{02}} + \frac{4\hat{\mu}_{22}}{\hat{\mu}_{11}^2} - \frac{4\hat{\mu}_{31}}{\hat{\mu}_{11}\hat{\mu}_{02}} - \frac{4\hat{\mu}_{13}}{\hat{\mu}_{11}\hat{\mu}_{02}} \right] \right\}^{1/2}$$

where, in terms of  $z_i = (y_i, z_i)$ ,  $\hat{\mu}_{gh} \equiv \Sigma(y_i - \bar{y})^g(z_i - \bar{z})^h/n$  (Cramer 1946, p. 359).

**Theorem.** For statistics of the form  $\hat{\theta} = t(\bar{Q}_1, \dots, \bar{Q}_A)$ , the nonparametric delta method and the infinitesimal jackknife give the same estimate of standard error (Efron 1981b).

The infinitesimal jackknife, the delta method, and the empirical influence function approach are three names for the same method. Notice that the results reported in line 7 of Table 2 show a severe downward bias. Efron and Stein (1981) show that the ordinary jackknife is always biased upwards, in a sense made precise in that paper. In the authors' opinion the ordinary jackknife is the method of choice if one does not want to do the bootstrap computations.

## Appendix

### BOOTSTRAP PROGRAM

The following FORTRAN program bootstraps the statistic defined by the user-specified function THETA. Comments in italics are not part of the FORTRAN code. Note that the random number subroutines IRAND and RAND will be installation dependent.

```

REAL Y(100), YSTAR(100), THSTAR(1000)
EXTERNAL THETA
N=100 sample size
NBOOT=1000 number of bootstraps
DO 10 I=1,N
    READ(5,*) Y(I) read in data
10  CONTINUE
    TEMP=THETA(N,Y)
    WRITE(6,100) TEMP write out value of theta for original sample
100  FORMAT(' THETA= ', f13.5)
    READ(5,*) ISEED read in seed for random number generator
    CALL IRAND(ISEED) initialize random number generator
    DO 20 I=1,NBOOT
        DO 30 J=1,N
            U=RAND() get a random number between 0 and 1
            II=INT(U*N) + 1 convert it to a random integer between 1 and N
            YSTAR(J)=Y(II) assign the jth element of bootstrap sample
30      CONTINUE
            THSTAR(I)=THETA(N,YSTAR) compute bootstrap value
20      CONTINUE
            THBAR=0
            DO 40 I=1,NBOOT
                THBAR=THBAR+THSTAR(I)/NBOOT compute bootstrap mean
40      CONTINUE
            THVAR=0
            DO 50 I=1,NBOOT
                THVAR=THVAR+(THSTAR(I)-THBAR)**2 compute bootstrap variance
50      CONTINUE
            SDBOOT=SQRT(THVAR/(NBOOT-1)) compute bootstrap estimate of standard error

```

```
WRITE(6,102) SDBOOT
102  FORMAT(' BOOTSTRAP ESTIMATE OF STANDARD ERROR= ', f13.6)
    WRITE(6,*)
    WRITE(6,103)
103  FORMAT(' BOOTSTRAP VALUES OF THETA: ')
    DO 60 I=1,NBOOT
    WRITE(6,*) THSTAR(I) write out bootstrap values for further analysis
60   CONTINUE
    STOP
    END
```

```
REAL FUNCTION THETA(N,Y)
REAL Y(N)
```

*compute statistic of interest for the sample y(1), y(2)...y(n)*

```
RETURN
END
```

## References

- Anderson, O. D. (1976). *Time Series Analysis and Forecasting*.
- Beran, R. (1984). Bootstrap methods in statistics. *Jber. d. Dt. Math. Verein* 86, 14-30.
- Bickel, P. J. and Freedman, D. A. (1981). Some asymptotic theory for the bootstrap. *Ann. Statist.* 9, 1196-1217.
- Box, G. E. P., and Cox, D. R. (1964). An analysis of transformations. *JRSS B* 26, 211-252.
- Breiman, L. and Friedman, J. H. (1984). Estimating optimal correlations for multiple regression and correlation. To appear *J. Amer. Statist. Assoc.*, March 1985.
- Cox, D. R. (1972). Regression models and life tables. *JRSS B* 34, 187-202.
- Cramér, H. (1966). *Mathematical Methods of Statistics*. Princeton University Press, New Jersey.
- Efron, B. (1979a). Bootstrap methods: another look at the jackknife. *Ann. Statist.* 7, 1-26.
- Efron, B. (1979b). Computers and the theory of statistics: thinking the unthinkable. *SIAM Review* 21, 460-480.
- Efron, B. (1981a). Censored data and the bootstrap. *J. Amer. Statist. Assoc.* 76, 312-319.
- Efron, B. (1981b). Maximum likelihood and decision theory. *Ann. Statist.* 9, 340-356.
- Efron, B. (1981c). Nonparametric estimates of standard error: the jackknife, the bootstrap, and other resampling methods. *Biometrika* 68, 589-599.
- Efron, B. (1982). The jackknife, the bootstrap, and other resampling plans. *SIAM CBMS-NSF Monograph* 38.
- Efron, B. (1982a). Transformation theory: how normal is a one parameter family of distributions? *Ann. Statist.* 10, 323-339.
- Efron, B. (1983). Estimating the error rate of a prediction rule: improvements in cross-validation. *J. Amer. Statist. Assoc.* 78, 316-331.
- Efron, B. (1984a). Bootstrap confidence intervals for a class of parametric problems. To appear in *Biometrika*.
- Efron, B. (1984b). Better bootstrap confidence intervals. Department of Statistics, Stanford University Technical Report No. 226.
- Efron, B. and Gong, G. (1983). A leisurely look at the bootstrap, the jackknife, and cross validation. *The American Statistician* 37, 36-48.
- Efron, B. and Stein, C. (1981). The jackknife estimate of variance. *Ann. Statist.* 9, 586-596.
- Fieller, E. C. (1954). Some problems in interval estimation. *JRSSB* 16, 175-183.
- Friedman, J. H. and Stuetzle, W. (1981). Projection pursuit regression. *J. Amer. Statist. Assoc.* 76, 817-823.
- Friedman, J. H. and Tibshirani, R. J. (1984). The monotone smoothing of scatterplots. *Technometrics* 26, 3, 243-250.

- Hampel, F. R. (1974). The influence curve and its role in robust estimation. *J. Amer. Statist. Assoc.* **69**, 383-383.
- Hastie, T. J. and Tibshirani, R. J. (1984). Discussion of Peter Huber's "Projection Pursuit". To appear in *Ann. Statist.*, August 1985.
- Hinkley, D. V. (1978). Improving the jackknife with special reference to correlation estimation. *Biometrika* **65**, 13-22.
- Jaekel, L. (1972). The infinitesimal jackknife. Memorandum MM 72-1215-11, Bell Laboratories, Murray Hill, New Jersey.
- Johnson, N. and Kotz, S. (1970). *Continuous Univariate Distributions*. Houghton Mifflin, Boston.
- Kaplan, E. L. and Meier, P. (1958). Nonparametric estimation from incomplete samples. *J. Amer. Statist. Assoc.* **53**, 457-481.
- Kiefer, J. and Wolfowitz, J. (1956). Consistency of the maximum likelihood estimator in the presence of infinitely many incidental parameters. *Ann. Math. Statist.* **27**, 887-906.
- Mallows, C. (1974). On some topics in robustness. Memorandum, Bell Laboratories, Murray Hill, New Jersey.
- Miller, R. G. (1974). The jackknife - a review. *Biometrika* **61**, 1-17.
- Miller, R. G. and Halpern, J. (1982). Regression with censored data. *Biometrika* **69**, 521-531.
- Rao, C. R. (1973). *Linear Statistical Inference and Its Applications*. Wiley, New York.
- Singh, K. (1981). On the asymptotic accuracy of Efron's bootstrap. *Ann. Statist.* **9**, 1187-1195.
- Therneau, T. (1983). Variance Reduction Techniques for the Bootstrap. Ph.D. Thesis, Stanford University, Department of Statistics.
- Tibshirani, R. J. and Hastie, T. J. (1984). Local likelihood estimation. Department of Statistics, Stanford University, Technical Report No. 97.
- Tukey, J. (1958). Bias and confidence in not quite large samples, abstract. *Ann. Math. Statist.* **29**, p. 614.