

# What Is Bootstrapping?

### 1.1. BACKGROUND

The bootstrap is a form of a larger class of methods that resample from the original data set and thus are called resampling procedures. Some resampling procedures similar to the bootstrap go back a long way [e.g., the jackknife goes back to Quenouille (1949), and permutation methods go back to Fisher and Pitman in the 1930s]. Use of computers to do simulation also goes back to the early days of computing in the late 1940s.

However, it was Efron (1979a) who unified ideas and connected the simple nonparametric bootstrap, for independent and identically distributed (IID) observations, which “resamples the data with replacement,” with earlier accepted statistical tools for estimating standard errors such as the jackknife and the delta method. This first method is now commonly called the nonparametric IID bootstrap. It was only after the later papers by Efron and Gong (1983), Efron and Tibshirani (1986), and Diaconis and Efron (1983) and the monograph Efron (1982a) that the statistical and scientific community began to take notice of many of these ideas, appreciate the extensions of the methods and their wide applicability, and recognize their importance.

After the publication of the Efron (1982a) monograph, research activity on the bootstrap grew exponentially. Early on, there were many theoretical developments on the asymptotic consistency of bootstrap estimates. In some of these works, cases where the bootstrap estimate failed to be a consistent estimator for the parameter were uncovered.

Real-world applications began to appear. In the early 1990s the emphasis shifted to finding applications and variants that would work well in practice. In the 1980s along with the theoretical developments, there were many simulation studies that compared the bootstrap and its variants with other competing estimators for a variety of different problems. It also became clear that

although the bootstrap had significant practical value, it also had some limitations.

A special conference of the Institute of Mathematical Statistics was held in Ann Arbor Michigan in May 1990, where many of the prominent bootstrap researchers presented papers exploring the applications and limitations of the bootstrap. The proceedings of this conference were compiled in the book *Exploring the Limits of Bootstrap*, edited by LePage and Billard and published by Wiley in 1992.

A second similar conference, also held in 1990 in Tier, Germany, covered many developments in bootstrapping. The European conference covered Monte Carlo methods, bootstrap confidence bands and prediction intervals, hypothesis tests, time series methods, linear models, special topics, and applications. Limitations of the methods were not addressed at this conference. Its proceedings were published in 1992 by Springer-Verlag. The editors for the proceedings were Jöckel, Rothe, and Sendler.

Although Efron introduced his version of the bootstrap in a 1977 Stanford University Technical Report [later published in a well-known paper in the *Annals of Statistics* (Efron, 1979a)], the procedure was slow to catch on. Many of the applications only began to be covered in textbooks in the 1990s.

Initially, there was a great deal of skepticism and distrust regarding bootstrap methodology. As mentioned in Davison and Hinkley (1997, p. 3): “In the simplest nonparametric problems, we do literally sample from the data, and a common initial reaction is that this is a fraud. In fact it is not.” The article in *Scientific American* (Diaconis and Efron, 1983) was an attempt to popularize the bootstrap in the scientific community by explaining it in layman’s terms and exhibiting a variety of important applications. Unfortunately, by making the explanation simple, technical details were glossed over and the article tended to increase the skepticism rather than abate it.

Other efforts to popularize the bootstrap that were partially successful with the statistical community were Efron (1982a), Efron and Gong (1981), Efron and Gong (1983), Efron (1979b), and Efron and Tibshirani (1986). Unfortunately it was only the *Scientific American* article that got significant exposure to a wide audience of scientists and researchers.

While working at the Aerospace Corporation in the period from 1980 to 1988, I observed that because of the *Scientific American* article, many of the scientist and engineers that I worked with had misconceptions about the methodology. Some supported it because they saw it as a way to use simulation in place of additional sampling (a misunderstanding of what kind of information the Monte Carlo approximation to the bootstrap actually gives). Others rejected it because they interpreted the *Scientific American* article as saying that the technique allowed inferences to be made from data without assumptions by replacing the need for additional “real” data with “simulated” data, and they viewed this as phony science (this is a misunderstanding that comes about because of the oversimplified exposition in the article).

Both views were expressed by my engineering colleagues at the Aerospace Corporation, and I found myself having to try to dispel both of these notions. In so doing, I got to thinking about how the bootstrap could help me in my own research and I saw there was a need for a book like this one. I also felt that in order for articles or books to popularize bootstrap techniques among the scientist, engineers, and other potential practitioners, some of the mathematical and statistical justification had to be presented and any text that skimmed over this would be doomed for failure.

The monograph by Mooney and Duvall (1993) presents only a little of the theory and in my view fails to provide the researcher with even an intuitive feel for why the methodology works. The text by Efron and Tibshirani (1993) was the first attempt at presenting the general methodology and applications to a broad audience of social scientists and researchers. Although it seemed to me to do a very good job of reaching that broad audience, Efron mentioned that he felt that parts of the text were still a little too technical to be clear to everyone in his intended audience.

There is a fine line to draw between being too technical to be understood by those without a strong mathematical background and being too simple to provide a true picture of the methodology devoid of misconceptions. To explain the methodology to those who do not have the mathematical background for a deep understanding of the bootstrap theory, we must avoid technical details on stochastic convergence and other advanced probability tools. But we cannot simplify it to the extent of ignoring the theory because that leads to misconceptions such as the two main ones previously mentioned.

In the late 1970s when I was a graduate student at Stanford University, I saw the theory develop first-hand. Although I understood the technique, I failed to appreciate its value. I was not alone, since many of my fellow graduate students also failed to recognize its great potential. Some statistics professors were skeptical about its usefulness as an addition to the current parametric, semiparametric, and nonparametric techniques.

Why didn't we give the bootstrap more consideration? At that time the bootstrap seemed so simple and straightforward. We did not see it as a part of a revolution in statistical thinking and approaches to data analysis. But today it is clear that this is exactly what it was!

A second reason why some graduate students at Stanford, and possibly other universities, did not elect the bootstrap as a topic for their dissertation research (including Naihua Duan, who was one of Efron's students at that time) is that the key asymptotic properties of the bootstrap appeared to be very difficult to prove. The mathematical approaches and results only began to be known when the papers by Bickel and Freedman (1981) and Singh (1981) appeared, and this was two to three years after many of us had graduated.

Gail Gong was one of Efron's students and the first Stanford graduate student to do a dissertation on the bootstrap. From that point on, many

students at Stanford and other universities followed as the flood gates opened to bootstrap research. Rob Tibshirani was another graduate student of Efron who did his dissertation research on the bootstrap and followed it up with the statistical science article (Efron and Tibshirani, 1986), a book with Trevor Hastie on general additive models, and the text with Efron on the bootstrap (Efron and Tibshirani, 1993). Other Stanford dissertations on bootstrap were Therneau (1983) and Hesterberg (1988). Both dealt with variance reduction techniques for reducing the number of bootstrap iterations necessary to get the Monte Carlo approximation to the bootstrap estimate to achieve a desired level of accuracy with respect to the bootstrap estimate (which is the limit as the number of bootstrap iterations approaches infinity).

My interest in bootstrap research began in earnest in 1983 after I read Efron's paper (Efron, 1983) on the bias adjustment in error rate estimation for classification problems. This applied directly to some of the work I was doing on target discrimination at the Aerospace Corporation and also later at Nichols Research Corporation. This led to a series of simulation studies that I published with Carlton Nealy and Krishna Murthy.

In the late 1980s I met Phil Good, who is an expert on permutation methods and was looking for a way to solve a particular problem that he was having trouble setting up in the framework of a permutation test. I suggested a straightforward bootstrap approach, and this led to comparisons of various procedures to solve the problem. It also opened up a dialogue between us about the virtues of permutation methods, bootstrap methods and other resampling methods, and the basic conditions for their applicability. We recognized that bootstrap and permutation tests were both part of the various resampling procedures that were becoming so useful but were not taught in the introductory statistics courses. That led him to write a series of books on permutation tests and resampling methods and led me to write the first edition of this text and later to incorporate the bootstrap in an introductory course in biostatistics and the text that Professor Robert Friis and I subsequently put together for the course (Chernick and Friis, 2002).

In addition to both being resampling methods, bootstrap and permutation methods could be characterized as computer-intensive, depending on the application. Both approaches avoid unverified parametric assumptions, by relying solely on the original sample. Both require minimal assumptions such as exchangeability of the observations under the null hypothesis. Exchangeability is a property of a random sample that is slightly weaker than the assumption that observations are independent and identically distributed. To be mathematically formal, for a sequence of  $n$  observations the sequence is exchangeable if the probability distribution of any  $k$  consecutive observations ( $k = 1, 2, 3, \dots, n$ ) does not change when the order of the observations is changed through a permutation.

The importance of the bootstrap is now generally recognized as has been noted in the article in the supplemental volume of the *Encyclopedia of Statistical Sciences* (1989 Bootstrapping—II by David Banks, pp. 17–22), the

inclusion of Efron's 1979 *Annals of Statistics* paper in *Breakthroughs in Statistics*, Volume II: *Methodology and Distribution*, S. Kotz and N. L. Johnson, editors (1992, pp. 565–595 with an introduction by R. Beran), and Hall's 1988 *Annals of Statistics* paper in *Breakthroughs in Statistics*, Volume III, S. Kotz and N. L. Johnson, editors (1997, pp. 489–518 with an introduction by E. Mammen). We can also find the bootstrap referenced prominently in the *Encyclopedia of Biostatistics*, with two entries in Volume I: (1) "Bootstrap Methods" by DeAngelis and Young (1998) and (2) "Bootstrapping in Survival Analysis" by Sauerbrei (1998).

The bibliography in the first edition contained 1650 references, and I have only expanded it as necessary. In the first edition I put an asterisk next to each of the 619 references that were referenced directly in the text and also numbered them in the alphabetical order that they were listed. In this edition I continue to use the asterisk to identify those books and articles referenced directly in the text but no longer number them.

The idea of sampling with replacement from the original data did not begin with Efron. Also even earlier than the first use of bootstrap sampling, there were a few related techniques that are now often referred to as resampling techniques. These other techniques predate Efron's bootstrap. Among them are the jackknife, cross-validation, random subsampling, and permutation procedures. Permutation tests have been addressed in standard books on nonparametric inference and in specialized books devoted exclusively to permutation tests including Good (1994, 2000), Edgington (1980, 1987, 1995), and Manly (1991, 1997).

The idea of resampling from the empirical distribution to form a Monte Carlo approximation to the bootstrap estimate may have been thought of and used prior to Efron. Simon (1969) has been referenced by some to indicate his use of the idea as a tool in teaching elementary statistics prior to Efron. Bruce and Simon have been instrumental in popularizing the bootstrap approach through their company Resampling Stats Inc. and their associated software. They also continue to use the Monte Carlo approximation to the bootstrap as a tool for introducing statistical concepts in a first elementary course in statistics [see Simon and Bruce (1991, 1995)]. Julian Simon died several years ago; but Peter Bruce continues to run the company and in addition to teaching resampling in online courses, he has set up a faculty to teach a variety of online statistics courses.

It is clear, however, that widespread use of the methods (particularly by professional statisticians) along with the many theoretical developments occurred only after Efron's 1979 work. That paper (Efron, 1979a) connected the simple bootstrap idea to established methods for estimating the standard error of an estimator, namely, the jackknife, cross-validation, and the delta method, thus providing the theoretical underpinnings that that were then further developed by Efron and other researchers.

There have been other procedures that have been called bootstrap that differ from Efron's concept. I mention two of them in Section 1.4. Whenever

I refer to the bootstrap in this text, I will be referring to Efron's version. Even Efron's bootstrap has many modifications. Among these are the double bootstrap, the smoothed bootstrap, the parametric bootstrap (discussed in Chapter 6), and the Bayesian bootstrap (which was introduced by Rubin in the missing data application described in Section 8.7). Some of the variants of the bootstrap are discussed in Section 2.1.2, including specialized methods specific to the classification problem [e.g., the 632 estimator introduced in Efron (1983) and the convex bootstrap introduced in Chernick, Murthy, and Neale (1985)].

In May 1998 a conference was held at Rutgers University, organized by Kesar Singh, a Rutgers statistics professor who is a prominent bootstrap researcher. The purpose of the conference was to provide a collection of papers on recent bootstrap developments by key bootstrap researchers and to celebrate the approximately 20 years of research since Efron's original work [first published as a Stanford Technical Report in 1977 and subsequently in the *Annals of Statistics* (Efron, 1979a)]. Abstracts of the papers presented were available from the Rutgers University Statistics Department web site.

Although no proceedings were published for the conference, I received copies of many of the papers by direct request to the authors. The presenters at the meeting included Michael Sherman, Brad Efron, Gutti Babu, C. R. Rao, Kesar Singh, Alastair Young, Dimitris Politis, J.-J. Ren, and Peter Hall. The papers that I received are included in the bibliography. They are Babu, Pathak, and Rao (1998), Sherman and Carlstein (1997), Efron and Tibshirani (1998), and Babu (1998).

This book is organized as follows. Chapter 1 introduces the key ideas and describes the wide range of applications. Chapter 2 deals with estimation and particularly the bias-adjusted estimators with emphasis on error rate estimation for discriminant functions. It shows through simulation studies how the bootstrap and variants such as the 632 estimator perform compared to the more traditional methods when the number of training samples is small. Also discussed are ratio estimates, estimates of medians, standard errors, and quantiles.

Chapter 3 covers confidence intervals and hypothesis tests. The 1–1 correspondence between confidence intervals and hypothesis tests is used to construct hypothesis tests based on bootstrap confidence intervals. We cover two so-called percentile methods and show how more accurate and correct bootstrap confidence intervals can be constructed. In particular, the hierarchy of percentile methods improved by bias correction BC and then BCa is given along with the rate of convergence for these methods and the weakening assumptions required for the validity of the method.

An application in a clinical trial to demonstrate the efficacy of the Tendril DX steroid lead in comparison to nonsteroid leads is also presented. Also covered is a very recent application to adaptive design clinical trials. In this application, proof of concept along with dose–response model identification methods and minimum effective dose estimates are included based on an

adaptive design. The author uses the MED as a parameter to generate “semi-parametric” bootstrap percentile methods.

Chapter 4 covers regression problems, both linear and nonlinear. An application of bootstrap estimates in nonlinear regression of the standard errors of parameters is given for a quasi-optical experiment. New in this edition is the coverage of bootstrap methods applied to outlier detection in least-squares regression.

Chapter 5 addresses time series models and related forecasting problems. This includes model based bootstrap and the various forms of block bootstrap. At the time of the first edition, the moving block bootstrap had been developed but was not very mature. Over the eight intervening years, there have been additional variations on the block bootstrap and more theory and applications. Recently, these developments have been well summarized in the text Lahiri (2003a). We have included some of those block bootstrap methods as well as the sieve bootstrap.

Chapter 6 provides a comparison with other resampling methods and recommends the preferred approach when there is clear evidence in the literature, either through theory or simulation, of its superiority. This was a unique feature of the book when the first edition was published. We have added to our list of resampling methods the  $m$  out of  $n$  bootstrap that we did not cover in the first edition. Although the  $m$  out of  $n$  bootstrap had been considered as a method to consider, it has only recently been proven to be important as a way to remedy inconsistency problems of the naïve bootstrap in many cases.

Chapter 7 deals with simulation methods, emphasizing the variety of available variance reduction techniques and showing the applications for which they can effectively be applied. This chapter is essentially the same as in the first edition.

Chapter 8 gives an account of a variety of miscellaneous topics. These include kriging (a form of smoothing in the analysis of spatial data) and other applications to spatial data, survey sampling, subset selection in both regression and discriminant analysis, analysis of censored data,  $p$ -value adjustment for multiplicity, estimation of process capability indices (measures of manufacturing process performance in quality assurance work), application of the Bayesian bootstrap in missing data problems, and the estimation of individual and population bioequivalence in pharmaceutical studies (often used to get acceptance of a generic drug when compared to a similar market-approved drug).

Chapter 9 describes examples in the literature where the ordinary bootstrap procedures fail. In many instances, modifications have been devised to overcome the problem, and these are discussed. In the first edition, remedies for the case of simple random sampling were discussed. In this edition, we also include remedies for extreme values including the result of Zelterman (1993) and the use of the  $m$  out of  $n$  bootstrap.

Bootstrap diagnostics are also discussed in Chapter 9. Efron’s jackknife-after-bootstrap is discussed because it is the first tool devised to help identify



whether or not a nonparametric bootstrap will work in a given application. The work from Efron (1992c) is described in Section 9.7.

Chapter 9 differs from the other chapters in that it goes into some of the technical probability details that the practitioner lacking this background may choose to skip. The practitioner may not need 1992c to understand exactly why these cases fail but should have a general awareness of the cases where the ordinary bootstrap fails and whether or not remedies have been found.

Each chapter (except Chapter 6) has a historical notes section. This section is intended as a guide to the literature related to the chapter and puts the results into their chronological order of development. I found that this was a nice feature in several earlier bootstrap books, including Hall (1992a), Efron and Tibshirani (1993), and Davison and Hinkley (1997). Although related references are cited throughout the text, the historical notes are intended to provide a perspective regarding when the techniques were originally proposed and how the key developments followed chronologically.

One notable change in the second edition is the increased description of techniques, particularly in Chapters 8 and 9.

## 1.2. INTRODUCTION

Two of the most important problems in applied statistics are the determination of an estimator for a particular parameter of interest and the evaluation of the accuracy of that estimator through estimates of the standard error of the estimator and the determination of confidence intervals for the parameter. Efron, when introducing his version of the “bootstrap” (Efron, 1979a), was particularly motivated by these two problems. Most important was the estimation of the standard error of the parameter estimator, particularly when the estimator was complex and standard approximations such as the delta methods were either not appropriate or too inaccurate.

Because of the bootstrap’s generality, it has been applied to a much wider class of problems than just the estimation of standard errors and confidence intervals. Applications include error rate estimation in discriminant analysis, subset selection in regression, logistic regression, and classification problems, cluster analysis, kriging (i.e., a form of spatial modeling), nonlinear regression, time series analysis, complex surveys,  $p$ -value adjustment in multiple testing problems, and survival and reliability analysis.

It has been applied in various disciplines including psychology, geology, econometrics, biology, engineering, chemistry, and accounting. It is our purpose to describe some of these applications in detail for the practitioner in order to exemplify its usefulness and illustrate its limitations. In some cases the bootstrap will offer a solution that may not be very good but may still be used for lack of an alternative approach. Since the publication of the first edition of this text, research has emphasized applications and has added to the long list of applications including particular applications in the pharma-



ceutical industry. In addition, modifications to the bootstrap have been devised that overcome some of the limitations that had been identified.

Before providing a formal definition of the bootstrap, here is an informal description of how it works. In its most general form, we have a sample of size  $n$  and we want to estimate a parameter or determine the standard error or a confidence interval for the parameter or even test a hypothesis about the parameter. If we do not make any parametric assumptions, we may find this difficult to do. The bootstrap provides a way to do this.

We look at the sample and consider the empirical distribution. The empirical distribution is the probability distribution that has probability  $1/n$  assigned to each sample value. The bootstrap idea is simply to replace the unknown population distribution with the known empirical distribution.

Properties of the estimator such as its standard error are then determined based on the empirical distribution. Sometimes these properties can be determined analytically, but more often they are approximated by Monte Carlo methods (i.e., we sample with replacement from the empirical distribution).

Now here is a more formal definition. Efron's bootstrap is defined as follows: Given a sample of  $n$  independent identically distributed random vectors  $X_1, X_2, \dots, X_n$  and a real-valued estimator  $(X_1, X_2, \dots, X_n)$  (denoted by  $\hat{\theta}$ ) of the parameter  $\theta$ , a procedure to assess the accuracy of  $\hat{\theta}$  is defined in terms of the empirical distribution function  $F_n$ . This empirical distribution function assigns probability mass  $1/n$  to each observed value of the random vectors  $X_i$  for  $i = 1, 2, \dots, n$ .

The empirical distribution function is the maximum likelihood estimator of the distribution for the observations when no parametric assumptions are made. The bootstrap distribution for  $\hat{\theta} - \theta$  is the distribution obtained by generating  $\hat{\theta}$ 's by sampling independently with replacement from the empirical distribution  $F_n$ . The bootstrap estimate of the standard error of  $\hat{\theta}$  is then the standard deviation of the bootstrap distribution for  $\hat{\theta} - \theta$ .

It should be noted here that almost any parameter of the bootstrap distribution can be used as a "bootstrap" estimate of the corresponding population parameter. We could consider the skewness, the kurtosis, the median, or the 95th percentile of the bootstrap distribution for  $\hat{\theta}$ .

Practical application of the technique usually requires the generation of bootstrap samples or resamples (i.e., samples obtained by independently sampling with replacement from the empirical distribution). From the bootstrap sampling, a Monte Carlo approximation of the bootstrap estimate is obtained. The procedure is straightforward.

1. Generate a sample with replacement from the empirical distribution (a bootstrap sample),
2. Compute \* the value of  $\hat{\theta}$  obtained by using the bootstrap sample in place of the original sample,
3. Repeat steps 1 and 2  $k$  times.

For standard error estimation,  $k$  is recommended to be at least 100. This recommendation can be attributed to the article Efron (1987). It has recently been challenged in a paper by Booth and Sarkar (1998). Further discussion on this recommendation can be found in Chapter 7.

By replicating steps 1 and 2  $k$  times, we obtain a Monte Carlo approximation to the distribution of  $\theta^*$ . The standard deviation of this Monte Carlo distribution of  $\theta^*$  is the Monte Carlo approximation to the bootstrap estimate of the standard error for  $\hat{\theta}$ . Often this estimate is simply referred to as the bootstrap estimate, and for  $k$  very large (e.g., 500) there is very little difference between the bootstrap estimator and this Monte Carlo approximation.

What we would like to know for inference is the distribution of  $\hat{\theta} - \theta$ . What we have is a Monte Carlo approximation to the distribution of  $\theta^* - \hat{\theta}$ . The key idea of the bootstrap is that for  $n$  sufficiently large, we expect the two distributions to be nearly the same.

In a few cases, we are able to compute the bootstrap estimator directly without the Monte Carlo approximation. For example, in the case of the estimator being the mean of the distribution of a real-valued random variable, Efron (1982a, p. 2) states that the bootstrap estimate of the standard error of is  $\hat{\sigma}_{\text{BOOT}} = [(n-1)/n]^{1/2} \hat{\sigma}$ , where  $\hat{\sigma}$  is defined as

$$\hat{\sigma} = \left[ \frac{1}{n(n-1)} \sum_{i=1}^n (x_i - \bar{x})^2 \right]^{1/2},$$

where  $x_i$  is the value of the  $i$ th observation and  $\bar{x}$  is the mean of the sample. As a second example, consider the case of testing the hypothesis of equality of distributions for censored matched pairs (i.e., observations whose values may be truncated). The bootstrap test applied to paired differences is equivalent to the sign test and the distribution under the null hypothesis is binomial with  $p = 1/2$ . So no bootstrap sampling is required to determine the critical region for the test.

The bootstrap is often referred to as a computer-intensive method. It gets this label because in most practical problems where it is deemed to be useful the estimation is complex and bootstrap samples are required. In the case of confidence interval estimation and hypothesis testing problems, this may mean at least 1000 bootstrap replications (i.e.,  $k = 1000$ ). In Section 7.1, we address the important practical issue of what value to use for  $k$ .

Methods for reducing the computer time by more efficient Monte Carlo sampling are discussed in Section 7.2. The examples above illustrate that there are cases for which the bootstrap is not computer-intensive at all!

Another point worth emphasizing here is that the bootstrap samples differ from the original sample because some of the observations will be repeated once, twice, or more in a bootstrap sample. There will also be some observations that will not appear at all in a particular bootstrap sample. Consequently, the values for  $\theta^*$  will vary from one bootstrap sample to the next.

The actual probability that a particular  $X_i$  will appear  $j$  times in a bootstrap sample for  $j = 0, 1, 2, \dots, n$ , can be determined using the multinomial distribution or alternatively by using classical occupancy theory. For the latter approach see (Chernick and Murthy, 1985). Efron (1983) calls these probabilities the repetition rates and discusses them in motivating the use of the .632 estimator (a particular bootstrap type estimator) for classification error rate estimation. A general account of the classical occupancy problem can be found in Johnson and Kotz (1977).

The basic idea behind the bootstrap is the variability of  $\theta^*$  (based on  $F_n$ ) around  $\hat{\theta}$  will be similar to (or mimic) the variability of  $\hat{\theta}$  (based on the true population distribution  $F$ ) around the true parameter value,  $\theta$ . There is good reason to believe that this will be true for large sample sizes, since as  $n$  gets larger and larger,  $F_n$  comes closer and closer to  $F$  and so sampling with replacements from  $F_n$  is almost like random sampling from  $F$ .

The strong law of large numbers for independent identically distributed random variables implies that with probability one,  $F_n$  converges to  $F$  pointwise [see Chung (1974, pp. 131–132) for details]. Strong laws pertaining to the bootstrap can be found in Athreya (1983). A stronger result, the Glivenko–Cantelli theorem [see Chung (1974, p. 133)], asserts that the empirical distribution converges uniformly with probability 1 to the distribution  $F$  when the observations are independent and identically distributed. Although not stated explicitly in the early bootstrap literature, this fundamental theoretical result lends credence to the bootstrap approach. The theorem was extended in Tucker (1959) to the case of a random sequence from a strictly stationary stochastic process.

In addition to the Glivenko–Cantelli theorem, the validity of the bootstrap requires that the estimator (a functional of the empirical distribution function) converge to the “true parameter value” (i.e., the functional for the “true” population distribution). A functional is simply a mapping that assigns a real value to a function. Most commonly used parameters of distribution functions can be expressed as functionals of the distribution, including the mean, the variance, the skewness, and the kurtosis.

Interestingly, sample estimates such as the sample mean can be expressed as the same functional applied to the empirical distribution. For more discussion of this see Chernick (1982), who deal with a form of a functional derivative called an influence function. The concept of an influence function was first introduced by Hampel (1974) as a method for comparing robust estimators.

Influence functions have had uses in robust statistical methods and in the detection of outlying observations in data sets. Formal treatment of statistical functionals can be found in Fernholtz (1983). There are also connections for the influence function with the jackknife and the bootstrap as shown by Efron (1982a).

Convergence of the bootstrap estimate to the appropriate limit (consistency) requires some sort of smoothness condition on the functional corresponding to the estimator. In particular, conditions given in Hall (1992a)

employ asymptotic normality for the functional and further allow for the existence of an Edgeworth expansion for its distribution function. So there is more needed. For independent and identically distributed observations we require (1) the convergence of  $F_n$  to  $F$  (this is satisfied by virtue of the Glivenko–Cantelli theorem), (2) an estimate that is the corresponding functional of  $F_n$  as the parameter is of  $F$  (satisfied for means, standard deviations, variances, medians, and other sample quantiles of the distribution), and (3) a smoothness condition on the functional. Some of the consistency proofs also make use of the well-known Berry–Esseen theorem [see Lahiri (2003a, pp. 21–22, Theorem 2.1) for the sample mean]. When the bootstrap fails (i.e., bootstrap estimates are inconsistent), it is often because the smoothness condition is not satisfied (e.g., extreme order statistics such as the minimum or maximum of the sample).

These Edgeworth expansions along with the Cornish–Fisher expansions not only can be used to assure the consistency of the bootstrap, but they also provide asymptotic rates of convergence. Examples where the bootstrap fails asymptotically, due to a lack of smoothness of the functional, are given in Chapter 9.

Also, the original bootstrap idea applies to independent identically distributed observations and is guaranteed to work only in large samples. Using the Monte Carlo approximation, bootstrapping can be applied to many practical problems such as parameter estimation in time series, regression, and analysis of variance problems, and even to problems involving small samples.

For some of these problems, we may be on shaky ground, particularly when small sample sizes are involved. Nevertheless, through the extensive research that took place in the 1980s and 1990s, it was discovered that the bootstrap sometimes works better than conventional approaches even in small samples (e.g., the case of error rate estimation for linear discriminant functions to be discussed in Section 2.1.2).

There is also a strong temptation to apply the bootstrap to a number of complex statistical problems where we cannot resort to classical theory to resort to. At least for some of these problems, we recommend that the practitioner try the bootstrap. Only for cases where there is theoretical evidence that the bootstrap leads us astray would we advise against its use.

The determination of variability in subset selection for regression, logistic regression, and its use in discriminant analysis problems provide examples of such complex problems. Another example is the determination of the variability of spatial contours based on the method of kriging. The bootstrap and alternatives in spatial problems are treated in Cressie (1991). Other books that cover spatial data problems are Mardia, Kent, and Bibby (1979) and Hall (1988c). Tibshirani (1992) provides some examples of the usefulness of the bootstrap in complex problems.

Diaconis and Efron (1983) demonstrate, with just five bootstrap sample contour maps, the value of the bootstrap approach in uncovering the vari-

ability in the contours. These problems that can be addressed by the bootstrap approach are discussed in more detail in Chapter 8.

### 1.3. WIDE RANGE OF APPLICATIONS

As mentioned at the end of the last section, there is a great deal of temptation to apply the bootstrap in a wide number of settings. In the regression case, for example, we may treat the vector including the dependent variable and the explanatory variable as independent random vectors, or alternatively we may compute residuals and bootstrap them. These are two distinct approaches to bootstrapping in regression problems which will be discussed in detail in Chapter 5.

In the case of estimating the error rate of a linear discriminant function, Efron showed in Efron (1982a, pp. 49–58) and Efron (1983) that the bootstrap could be used to (1) estimate the bias of the “apparent error rate” estimate (a naïve estimate of error rate that is also referred to as the resubstitution estimate) and (2) produce an improved error rate estimate by adjusting for the bias.

The most attractive feature of the bootstrap and the permutation tests described in Good (1994) is the freedom they provide from restrictive parametric assumptions and simplified models. There is no need to force Gaussian or other parametric distributional assumptions on the data.

In many problems, the data may be skewed or have a heavy-tailed distribution or may even be multimodal. The model does not need to be simplified to some “linear” approximation, and the estimator itself can be complicated.

We do not require an analytic expression for the estimator. The bootstrap Monte Carlo approximation can be applied as long as there is a computational method for deriving the estimator. That means that we can numerical integrate using iterative schemes to calculate the estimator. The bootstrap doesn’t care. The only price we pay for such complications is in the time and cost for the computer usage (which is becoming cheaper and faster).

Another feature that makes the bootstrap approach attractive is its simplicity. We can formulate bootstrap simulations for almost any conceivable problem. Once we program the computer to carry out the bootstrap replications, we let the computer do all the work. A danger to this approach is that a practitioner might bootstrap at will, without consulting a statistician (or considering the statistical implications) and without giving careful thought to the problem.

This book will aid the practitioner in the proper use of the bootstrap by acquainting him with its advantages and limitations, lending theoretical support where available and Monte Carlo results where the theory is not yet available. Theoretical counterexamples to the consistency of bootstrap estimates also provide guidelines to its limitations and warn the practitioner when not to

apply the bootstrap. Some simulation studies also provide such negative results.

However, over the past 9 years, modifications to the basic or naïve bootstrap that fails due to inconsistency have been constructed to be consistent. One notable approach to be covered in Chapter 9 is the  $m$ -out-of- $n$  bootstrap. Instead of sampling  $n$  times with replacement from the empirical distribution where  $n$  is the original sample size, the  $m$ -out-of- $n$  bootstrap samples  $m$  times with replacement from the empirical distribution where  $m$  is chosen to be less than  $n$ . In the asymptotic theory both  $m$  and  $n$  tend to infinity but  $m$  increases at a slower rate. The rate to choose depends on the application.

I believe, as do many others now, that many simulation studies indicate that the bootstrap can safely be applied to a large number of problems even where strong theoretical justification does not yet exist. For many problems where realistic assumptions make other statistical approaches impossible or at least intractable, the bootstrap at least provides a solution even if it is not a very good one. For some people in certain situations, even a poor solution is better than no solution.

Another problem that creates difficulties for the scientist and engineer is that of missing data. In designing an experiment or a survey, we may strive for balance in the design and choose specific samples sizes in order to make the planned inferences from the data. The correct inference can be made only if we observe the complete data set.

Unfortunately, in the real world, the cost of experimentation, faulty measurement, or lack of response from those selected for the survey may lead to incomplete and possibly unbalanced designs. Milliken and Johnson (1984) refer to such problem data as messy data.

In Milliken and Johnson (1984, 1989) they provide ways to analyze messy data. When data are missing or censored, bootstrapping provides another approach for dealing with the messy data (see Section 8.4 for more details on censored data, and see Section 8.7 for an application to missing data).

The bootstrap alerts the practitioner to variability in his data, of which he or she may not be aware. In regression, logistic regression, or discriminant analysis, stepwise subset selection is a commonly used method available in most statistical computer packages. The computer does not tell the user how arbitrary the final selection actually is. When a large number of variables or features are included and many are correlated or redundant, there can be a great deal of variability to the selection. The bootstrap samples enable the user to see how the chosen variables or features change from bootstrap sample to bootstrap sample and provide some insight as to which variables or features are really important and which ones are correlated and easily substituted for by others. This is particularly well illustrated by the logistic regression problem studied in Gong (1986). This problem is discussed in detail in Section 8.2.

In the case of kriging, spatial contours of features such as pollution concentration are generated based on data at monitoring stations. The method is a

form of interpolation between the stations based on certain statistical spatial modeling assumptions. However, the contour maps themselves do not provide the practitioner with an understanding of the variability of these estimates. Kriging plots for different bootstrap samples provide the practitioner with a graphical display of this variability and at least warn him of variability in the data and analytic results. Diaconis and Efron (1983) make this point convincingly, and I will demonstrate this application in Section 8.1. The practical value of this cannot be underestimated!

Babu and Feigelson (1996) discuss applications in astronomy. They devote a whole chapter (Chapter 5, pp. 93–103) to resampling methods, emphasizing the importance of the bootstrap.

In clinical trials, sample sizes are determined based on achieving a certain power for a statistical hypothesis of efficacy of the treatment. In Section 3.3, I show an example of a clinical trial for a pacemaker lead (Pacesetter's Tendril DX model). In this trial, the sample sizes for the treatment and control leads were chosen to provide an 80% chance of detecting a clinically significant improvement (decrease of 0.5 volts) in the average capture threshold at the three-month follow-up for the experimental Tendril DX lead (model 1388T) compared to the respective control lead (Tendril model 1188T) when applying a one-sided significance test at the 5% significance level. This was based on the standard normal distribution theory. In the study, nonparametric methods were also considered. Bootstrap confidence intervals based on Efron's percentile method were used to do the hypothesis test without needing parametric assumptions. The Wilcoxon rank sum test was another nonparametric procedure that was used to test for a statistically significant change in capture threshold.

A similar study for a passive fixation lead, the Passive Plus DX lead, was conducted to get FDA approval for the steroid eluting version of this type of lead. In addition to comparing the investigational (steroid eluting) lead with the non-steroid control lead, using both the bootstrap (percentile method) and Wilcoxon rank sum tests, I also tried the bootstrap percentile  $t$  confidence intervals for the test. This method theoretically can give a more accurate confidence interval. The results were very similar and conclusive at showing the efficacy of the steroid lead. The percentile  $t$  method of confidence interval estimation is described in Section 3.1.5.

However, the statistical conclusion for such a trial is based on a single test at the three-month follow-up after all 99 experimental and 33 control leads have been implanted, and the patients had threshold tests at the three-month follow-up.

In the practice of clinical trials, the investigators do not want to wait for all the patients to reach their three-month follow-up before doing the analysis. Consequently, it is quite common to do interim analyses at some point or points in the trial (it could be one in the middle of the trial or two at the one-third and two-thirds points in the trial). Also, separate analyses are sometimes done on subsets of the population. Furthermore, sometimes separate analyses



are done on subsets of the population. These examples are all situations where multiple testing is involved. Multiple testing requires specific techniques for controlling the type I error rate (in this context the so-called family-wise error rate is the error rate that is controlled. Equivalent to controlling the family-wise type I error rate the  $p$ -values for the individual tests can be adjusted. Probability bounds such as the Bonferroni can be used to give conservative estimates of the  $p$ -value or simultaneous inference methods can be used [see Miller (1981b) for a thorough treatment of this subject].

An alternative approach would be to estimate the  $p$ -value adjustment by bootstrapping. This idea has been exploited by Westfall and Young and is described in detail in Westfall and Young (1993). We will attempt to convey the key concepts. The application of bootstrap  $p$ -value adjustment to the Passive Plus DX clinical trial data is covered in Section 8.5. Consult Miller (1981b), Hsu (1996), and/or Westfall and Young (1993) for more details on multiple testing,  $p$ -value adjustment, and multiple comparisons.

In concluding this section, we wish to emphasize that the bootstrap is not a panacea. There are certainly practical problems where classical parametric methods are reasonable and provide either more efficient estimates or more powerful hypothesis tests. Even for some parametric problems, the parametric bootstrap, as discussed by Davison and Hinkley (1997, p. 3) and illustrated by them on pages 148 and 149, can be useful.

What the bootstrap does do is free the scientist from restrictive modeling and distributional assumptions by using the power of the computer to replace difficult analysis. In an age when computers are becoming more and more powerful, inexpensive, fast, and easy to use, the future looks bright for additional use of these so-called computer-intensive statistical methods, as we have seen over the past decade.

#### 1.4. HISTORICAL NOTES

It should be pointed out that bootstrap research began in the late 1970s, although many key related developments can be traced back to earlier times. Most of the important theoretical development; took place in the 1980s after Efron (1979a). The first proofs of the consistency of the bootstrap estimate of the sample mean came in 1981 with the papers of Singh (1981) and Bickel and Freedman (1981).

Regarding this seminal paper by Efron (1979a), Davison and Hinkley (1997) write “The publication in 1979 of Bradley Efron’s first article on bootstrap methods was a major event in Statistics, at once synthesizing some of the earlier resampling ideas and establishing a new framework for simulation-based statistical analysis. The idea of replacing complicated and often inaccurate approximations to biases, variances, and other measures of uncertainty by computer simulations caught the imagination of both theoretical researchers and users of statistical methods.”

As mentioned earlier in this chapter, a number of related techniques are often referred to as resampling techniques. These other resampling techniques predate Efron's bootstrap. Among these are the jackknife, cross-validation, random subsampling, and the permutation test procedures described in Good (1994), Edgington (1980, 1987, 1995), and Manly (1991, 1997).

Makinodan, Albright, Peter, Good, and Heidrick (1976) apply permutation tests to study the effect of age in mice on the mediation of immune response. Due to the fact that an entire factor was missing, the model and the permutation test provides a clever way to deal with imbalance in the data. A detailed description is given in Good (1994, pp. 58–59).

Efron himself points to some of the early work of R. A. Fisher (in the 1920s) on maximum likelihood estimation as the inspiration for many of the basic ideas. The jackknife was introduced by Quenouille (1949) and popularized by Tukey (1958), and Miller (1974) provides an excellent review of the jackknife methods. Extensive coverage of the jackknife can be found in the book by Gray and Schucany (1972).

Bickel and Freedman (1981) and Singh (1981) presented the first results demonstrating the consistency of the bootstrap under certain mathematical conditions. Bickel and Freedman (1981) also provide a counterexample for consistency of the nonparametric bootstrap, and this is also illustrated by Schervish (1995, p. 330, Example 5.80). Gine and Zinn (1989) provide necessary conditions for the consistency of the bootstrap for the mean.

Athreya (1987a,b), Knight (1989), and Angus (1993) all provide examples where the bootstrap failed to be consistent due to its inability to meet certain necessary mathematical conditions. Hall, Hardle, and Simar (1993) showed that estimators for bootstrap distributions can also be inconsistent.

The general subject of empirical processes is related to the bootstrap and can be used as a tool to demonstrate consistency (see Csorgo, 1983; Shorack and Wellner, 1986; van der Vaart and Wellner, 1996). Fernholtz (1983) provides the mathematical theory of statistical functionals and functional derivatives (such as influence functions) that relate to bootstrap theory.

Quantile estimation via bootstrapping appears in Helmers, Janssen, and Veraverbeke (1992) and in Falk and Kaufmann (1991). Csorgo and Mason (1989) bootstrap the empirical distribution and Tu (1992) uses jackknife pseudovalue to approximate the distribution of a general standardized functional statistic.

Subsampling methods began with Hartigan (1969, 1971, 1975) and McCarthy (1969). These papers are discussed briefly in the development of bootstrap confidence intervals in Chapter 3. A more recent account is given by Babu (1992).

Young and Daniels (1990) discuss the bias that is introduced in Efron's nonparametric bootstrap by the use of the empirical distribution as a substitute for the true unknown distribution.

Diaconis and Holmes (1994) show how to avoid the Monte Carlo approximation to the bootstrap by cleverly enumerating all possible bootstrap samples using what are called Gray codes.

The term bootstrap has been used in other similar contexts which predate Efron's work, but these methods are not the same and some confusion occurs. When I gave a presentation on the bootstrap at the Aerospace Corporation in 1983 a colleague, Dr. Ira Weiss, mentioned that he used the bootstrap in 1970 long before Efron coined the term. After looking at Ira's paper, I realized that it was a different procedure with a similar idea.

Apparently, control theorists came up with a procedure for applying Kalman filtering with an unknown noise covariance which they also named the bootstrap. Like Efron, they were probably thinking of the old adage "picking yourself up by your own bootstraps" (as was attributed to the fictional Baron von Munchausen as a trick for climbing out from the bottom of a lake) when they chose the term to apply to an estimation procedure that avoids a priori assumptions and uses only the data at hand. A survey and comparison of procedures for dealing with the problem of unknown noise covariance including this other bootstrap technique is given in Weiss (1970). The term bootstrap has also been used in totally different contexts by computer scientists.

An entry on bootstrapping in the *Encyclopedia of Statistical Science* (1981, Volume 1, p. 301) is provided by the editors and is very brief. In 1981 when that volume was published, the true value of bootstrapping was not fully appreciated. The editors subsequently remedied this with an article in the supplemental volume.

The point, however, is that the original entry cited only three references. The first, Efron's *SIAM Review* article (Efron, 1979b), was one of the first published works describing Efron's bootstrap. The second article from *Technometrics* by Fuchs (1978) does not appear to deal with the bootstrap at all! The third article by LaMotte (1978) and also in *Technometrics* does refer to a bootstrap but does not mention any of Efron's ideas and appears to be discussing a different bootstrap.

Because of these other bootstraps, we have tried to refer to the bootstrap as Efron's bootstrap; a few others have done the same, but it has not caught on. In the statistical literature, reference to the bootstrap will almost always mean Efron's bootstrap or some derivative of it. In the engineering literature an ambiguity may exist and we really need to look at the description of the procedure in detail to determine precisely what the author means.

The term bootstrap has also commonly appeared in the computer science literature, and I understand that mathematicians use the term to describe certain types of numerical solutions to partial differential equations. Still it is my experience that if I search for articles in mathematical or statistical indices using the keyword "bootstrap," I would find that the majority of the articles referred to Efron's bootstrap or a variant of it. I wrote the preceding statement back in 1999 when the first edition was published. Now in 2007, I formed the basis for the second bibliography of the text by searching the Current Index

to Statistics (CIS) for the years 1999 to 2007 with only the keyword “bootstrap” required to appear in the title or the list of key words. Of the large number of articles and books that I found from this search, all of the references were referring to Efron’s bootstrap or a method derived from the original idea of Efron. The term “boofstrap” is used these days as a noun or a verb.

However, I have no similar experience with the computer science literature or the engineering literature. But Efron’s bootstrap now has a presence in these two fields as well. In computer science there have been many meetings on the interface between computer science and statistics, and much of the common ground involves computer-intensive methods such as the bootstrap. Because of the rapid growth of bootstrap application in a variety of industries, the “statistical” bootstrap now appears in some of the physics and engineering journals including the IEEE journals. In fact the article I include in Chapter 4, an application of nonlinear regression to a quasi-optical experiment, I coauthored with three engineers and the article appeared in the *IEEE Transactions on Microwave Theory and Techniques*.

Efron (1983) compared several variations to the bootstrap estimate. He considered simulation of Gaussian distributions for the two-class problem (with equal covariances for the classes) and small sample sizes (e.g., a total of, say, 14–20 training samples split equally among the two populations). For linear discriminant functions, he showed that the bootstrap and in particular the .632 estimator are superior to the commonly used leave-one-out estimate (also called cross-validation by Efron). Subsequent simulation studies will be summarized in Section 2.1.2 along with guidelines for the use of some of the bootstrap estimates.

There have since been a number of interesting simulation studies that show the value of certain bootstrap variants when the training sample size is small (particularly the estimator referred to as the .632 estimate). In a series of simulations studies, Chernick, Murthy, and Nealy (1985, 1986, 1988a,b) confirmed the results in Efron (1983). They also showed that the .632 estimator was superior when the populations were not Gaussian but had finite first moments. In the case of Cauchy distributions and other heavy-tailed distributions from the Pearson VII family of distributions which do not have finite first moments, they showed that other bootstrap approaches were better than the .632 estimator.

Other related simulation studies include Chatterjee and Chatterjee (1983), McLachlan (1980), Snapinn and Knoke (1984, 1985a,b, 1988), Jain, Dubes, and Chen (1987) and Efron and Tibshirani (1997a). We summarize the results of these studies and provide guidelines to the use of the bootstrap procedures for linear and quadratic discriminant functions in Section 2.1.2. McLachlan (1992) also gives a good summary treatment to some of this literature. Additional theoretical results can be found in Davison and Hall (1992). Hand (1986) is another good survey article on error rate estimation. The 632+ estimator proposed by Efron and Tibshirani (1997a) was applied to an ecological

problem by Furlanello, Merler, Chemini, and Rizzoli (1998). Ueda and Nakano (1995) apply the bootstrap and cross-validation to error rate estimation for neural network-type classifiers. Hand (1981, p. 189; 1982, pp. 178–179) discusses the bootstrap approach to estimating the error rates in discriminant analysis.

In the late 1980s and the 1990s, a number of books appeared that covered some aspect of bootstrapping at least partially. Noreen's book (Noreen, 1989) deals with the bootstrap in very elementary ways for hypothesis testing only.

There are now several survey articles on bootstrapping in general, including Babu and Rao (1993), Young (1994), Stine (1992), Efron (1982b), Efron and LePage (1992), Efron and Tibshirani (1985, 1986, 1996a, 1997b), Hall (1994), Manly (1993), Gonzalez-Manteiga, Prada-Sanchez, and Romo (1993), Politis (1998), and Hinkley (1984, 1988). Overviews on the bootstrap or special aspects of bootstrapping include Beran (1984b), Leger, Politis, and Romano (1992), Pollack, Simon, Bruce, Borenstein, and Lieberman (1994), and Fiellin and Feinstein (1998) on the bootstrap in general; Babu and Bose (1989), DiCiccio and Efron (1996), and DiCiccio and Romano (1988, 1990) on confidence intervals; Efron (1988b) on regression; Falk (1992a) on quantile estimation; and DeAngelis and Young (1992) on smoothing. Lanyon (1987) reviews the jackknife and bootstrap for applications to ornithology. Efron (1988c) gives a general discussion of the value of bootstrap confidence intervals aimed at an audience of psychologists.

The latest edition of *Kendall's Advanced Theory of Statistics*, Volume I, deals with the bootstrap as a tool for estimating standard errors in Chapter 10 [see Stuart and Ord (1993, pp. 365–368)].

The use of the bootstrap to compute standard errors for estimates and to obtain confidence intervals for multilevel linear models is given in Goldstein (1995, pp. 60–63). Waclawiw and Liang (1994) give an example of parametric bootstrapping using generalized estimating equations. Other works involving the bootstrap and jackknife in estimating equation models include Lele (1991a,b).

Lehmann and Casella (1998) mention the bootstrap as a tool in reducing the bias of an estimator (p. 144) and in the attainment of higher order efficiency (p. 519). Lehmann (1999, Section 6.5, pp. 420–435) presents some details on the asymptotic properties of the bootstrap.

In the context of generalized least-squares estimation of regression parameters Carroll and Ruppert (1988, pp. 26–28) describe the use of the bootstrap to get confidence intervals. In a brief discussion, Nelson (1990) mentions the bootstrap as a potential tool in regression models with right censoring of data for application to accelerated lifetime testing. Srivastava and Singh (1989) deal with the application of bootstrap in multiplicative models. Bickel and Ren (1996) employ an  $m$ -out-of- $n$  bootstrap for goodness of fit tests with doubly censored data.

McLachlan and Basford (1988) discuss the bootstrap in a number of places as an approach for determining the number of distributions or modes in a

mixture model. Another excellent text on mixture models is Titterington, Smith, and Makov (1985). Efron and Tibshirani (1996b) take a novel approach to bootstrapping that can be applied to the determination of the number of modes in a density function and the number of variables in a model. In addition to determining the number of modes, Romano (1988c) uses the bootstrap to determine the location of a mode.

Linhart and Zucchini (1986, pp. 22–23) describe how the bootstrap can be used for model selection. Thompson (1989, pp. 42–43) mentions the use of bootstrap techniques for estimating parameters in growth models (i.e., a non-linear regression problem). McDonald (1982) shows how smoothed or ordinary bootstrap samples can be drawn to obtain regression estimates.

Rubin (1987, pp. 44–46) discusses his “Bayesian” bootstrap for problems of imputation. The original paper on the Bayesian bootstrap is Rubin (1981). Banks (1988) provides a modification to the Bayesian bootstrap. Other papers involving the Bayesian bootstrap are Lo (1987, 1988, 1993a) and Weng (1989). Geisser (1993) discusses the bootstrap with respect to predictive distributions (another Bayesian concept). Ghosh and Meeden (1997, pp. 140–149) discuss applications of the Bayesian bootstrap to finite population sampling. The Bayesian bootstrap is often applied to imputation problems. Rubin (1996) is a survey article detailing the history of multiple imputation. At the time of the article the method of multiple imputation had been studied for more than 18 years.

Rey (1983) devotes Chapter 5 of his monograph to the bootstrap. He is using it in the context of robust estimation. His discussion is particularly interesting because he mentions both the pros and the cons and is critical of some of the early claims made for the bootstrap [particularly in Diaconis and Efron (1983)].

Staudte and Sheather (1990) deal with the bootstrap as an approach to estimating standard errors of estimates. They are particularly interested in the standard errors of robust estimators. Although they do deal with hypothesis testing, they do not use the bootstrap for any hypothesis testing problems. Their book includes a computer disk that has Minitab macros for bootstrapping in it. Minitab computer code for these macros is presented in Appendix D of their book.

Barnett and Lewis (1995) discuss the bootstrap as it relates to checking modeling assumptions in the face of outliers. Agresti (1990) discusses the bootstrap as it can be applied to categorical data.

McLachlan and Krishnan (1997) discuss the bootstrap in the context of robust estimation of a covariance matrix. Beran and Srivastava (1985) provide bootstrap tests for functions of a covariance matrix. Other papers covering the theory of the bootstrap as it relates to robust estimators are Babu and Singh (1984b) and Arcones and Gine (1992). Lahiri (1992a) does bootstrapping of *M*-estimators (a type of robust location estimator).

The text by van der Vaart and Wellner (1996) is devoted to weak convergence and empirical processes. Empirical process theory can be applied to



obtain important results in bootstrapping, and van der Vaart and Wellner illustrate this in Section 3.6 of their book (14 pages devoted to the subject of bootstrapping, pp. 345–359).

Hall (1992a) considers functionals that admit Edgeworth expansions. Edgeworth expansions provide insight into the accuracy of bootstrap confidence intervals, the value of bootstrap hypothesis tests, and use of the bootstrap in parametric regression. It also provides guidance to the practitioner regarding the variants of the bootstrap and the Monte Carlo approximations. Some articles relating Edgeworth expansions to applications of the bootstrap include Abramovitch and Singh (1985), Bhattacharya and Qumsiyeh (1989), Babu and Singh (1989), and Bai and Rao (1991, 1992).

Chambers and Hastie (1991) discuss applications of statistical models through the use of the *S* language. They discuss the bootstrap in various places.

Gifi (1990) applies the bootstrap to multivariate problems. Other uses of the bootstrap in branches of multivariate analysis are documented Diaconis and Efron (1983), who apply the bootstrap to principal component analysis, and Greenacre (1984), who covers the use of bootstrapping in correspondence analysis.

One of the classic texts on multivariate analysis is Anderson (1959), which was the first to provide extensive coverage of the theory based on the multivariate normal model. In the second edition of the text, Anderson (1984), he introduces some bootstrap applications. Flury (1997) provides another recent account of multivariate analysis. Flury (1988) is a text devoted to the principal components technique and so is Jolliffe (1986). Seber (1984), Gnandesikan (1977, 1997), Hawkins (1982), and Mardia, Kent, and Bibby (1979) all deal with the subject of multivariate analysis and multivariate data.

Scott (1992, pp. 257–260) discusses the bootstrap as a tool in estimating standard errors and confidence intervals in the context of multivariate density estimation. Other articles where the bootstrap appears as a density estimation tool are Faraway and Jhun (1990), Falk (1992b), and Taylor and Thompson (1992).

Applications in survival analysis include Burr (1994), Hsieh (1992), LeBlanc and Crowley (1993) and Gross and Lai (1996a).

An application of the double bootstrap appears in McCullough and Vinod (1998). Application to the estimation of correlation coefficients can be found in Lunneborg (1985) and Young (1988a).

General discussion of bootstrapping related to nonparametric procedures include Romano (1988a), Romano (1989b), and Simonoff (1986), where goodness of fit of distributions in sparse multinomial data problems is addressed using the bootstrap. Tu, Burdick, and Mitchell (1992) apply bootstrap resampling to nonparametric rank estimation.

Hahn and Meeker (1991) briefly discuss bootstrap confidence intervals. Frangos and Schucany (1990) discuss the technical aspects of estimating the acceleration constant for Efron's  $BC_a$  confidence interval method. Bickel and



Krieger (1989) use the bootstrap to attain confidence bands for a distribution function, and Wang and Wahba (1995) get bootstrap confidence bands for smoothing splines and compare them to bands constructed using Bayesian methods.

Bailey (1992) provides a form of bootstrapping for order statistics and other random variables whose distributions can be represented as convolutions of other distributions. By substituting the empirical distributions for the distributions in the convolution, a “bootstrap” distribution for the random variable is derived.

Beran (1982) compares the bootstrap with various competitive methods in estimating sampling distributions. Bau (1984) does bootstrapping for statistics involving linear combinations. Parr (1983) is an early reference comparing the bootstrap, the jackknife, and the delta method in the context of bias and variance estimation. Hall (1988d) deals with the rate of convergence for bootstrap approximations.

Applications to directional data include Fisher and Hall (1989) and Ducharme, Jhun, Romano, and Troung (1985). Applications to finite population sampling include Chao and Lo (1985), Booth, Butler, and Hall (1994), Kuk (1987, 1989), and Sitter (1992b).

Applications have appeared in a variety of disciplines. These include Choi, Nam, and Park (1996), quality assurance (for process capability indices); Jones, Wortberg, Kreissig, Hammock, and Rocke (1996), engineering; Bajgier (1992), Seppala, Moskowitz, Plante, and Tang (1995) and Liu and Tang (1996), process control; Chao and Huwang (1987), reliability; Coakley (1996), image processing; Bar-Ness and Punt (1996), communications; and Zoubir and Iskander (1996) and Zoubir and Boashash (1998), signal processing. Ames and Muralidhar (1991) and Biddle, Bruton, and Siegel (1990) provide applications in auditing. Robeson (1995) applies the bootstrap in meteorology, Tambour and Zethraeus (1998) in economics, and Tran (1996) in sports medicine. Roy (1994) and Schafer (1992) provide applications in chemistry, Rothery (1985) and Lanyon (1987) in ornithology. Das Peddada and Chang (1992) give an application in physics. Mooney (1996) covers bootstrap applications in political science. Adams, Gurevitch, and Rosenberg (1997) and Shipley (1996) apply the bootstrap to problems in ecology; Andrieu, Caraux, and Gascuel (1997) in evolution; and Aastveit (1990), Felsenstein (1985), Sanderson (1989, 1995), Sitnikova, Rzhetsky, and Nei (1995), Leal and Ott (1993), Tivang, Nienhuis, and Smith (1994), Schork (1992), Zharkikh and Li (1992, 1995) in genetics. Lunneborg (1987) gives us applications in the behavioral sciences. Abel and Berger (1986) and Brey (1990) give applications in biology. Aegerter, Muller, Nakache, and Boue (1994), Baker and Chu (1990), Barlow and Sun (1989), Mapleson (1986), Tsodikov, Hasenclever, and Loeffler (1998), and Wahrendorf and Brown (1980) apply the bootstrap to a variety of medical problems.

The first monograph on the bootstrap was Efron (1982a). In the 1990s there were a number of books introduced that are dedicated to bootstrapping and/or

related resampling methods. These include Beran and Ducharme (1991), Chernick (1999), Davison and Hinkley (1997), Efron and Tibshirani (1993), Hall (1992a), Helmers (1991b), Hjorth (1994), Janas (1993), Mammen (1992b), Manly (1997), Mooney and Duval (1993), Shao and Tu (1995), and Westfall and Young (1993). Schervish (1995) devotes a section and Sprent (1998) has a whole chapter on the bootstrap. In addition to the bootstrap chapter, the bootstrap is discussed throughout Sprent (1998) because it is one of a few data-driven statistical methods that are the theme of the text. Chernick and Friis (2002) introduce bootstrapping in a biostatistics text for health science students, Hesterberg, Moore, Monaghan, Clipson and Epstein (2003) is a chapter for an introductory statistics text that covers bootstrap and permutation methods and it has been incorporated as Chapter 18 of Moore, McCabe, Duckworth and Sclove (2003) as well as Chapter 14 of the on-line 5th Edition of Moore and McCabe (2005).

Efron has demonstrated the value of the bootstrap in a number of applied and theoretical contexts. In Efron (1988a), he provides three examples of the value of inference through computer-intensive methods. In Efron (1992b) he shows how the bootstrap has impacted theoretical statistics by raising six basic theoretical questions.

Davison and Hinkley (1997) provide a computer diskette with a library of useful SPLUS functions that can be used to implement bootstrapping in a variety of problems. These routines can be used with the commercial Version 3.3 of SPLUS, and they are described in Chapter 11 of the book. Barbe and Bertail (1995) deal with weighted bootstraps.

Two conferences were held in 1990, one in Michigan and the other in Trier, Germany. These conferences specialized in research developments in bootstrap and related techniques. Proceedings from these conferences were published in LePage and Billard (1992) for the Michigan conference, and in Jockel, Rothe, and Sandler (1992) for the conference in Trier.

In 2003, a portion of an issue of the journal *Statistical Science* was devoted to the bootstrap on its Silver Anniversary. It included articles by Efron, Casella, and others. The text by Lahiri (2003a) covers the dependent cases in detail emphasizing block bootstrap methods for time series and spatial data. It also covers model-based methods and provides some coverage of the independent case. Next to this text, Lahiri (2003a) provides the most recent coverage on bootstrap methods. It provides detailed descriptions of the methodology along with rigorous proofs of important theorems. It also uses simulations for comparison of various methods.

## 1.5. SUMMARY

In this chapter, I have given a basic explanation of Efron's nonparametric bootstrap. I have followed this up with explanations as to why the procedure can be expected to work in a wide variety of applications and also have given a historical perspective to the development of the bootstrap and its early

acceptance or lack thereof. I have also pointed out some of the sections in subsequent chapters and additional references that will provide more details than the brief discussions given in this chapter.

I have tried to make the discussion casual and friendly with each concept described as simply as possible and each definition stated as clearly as I can make them. However, it was necessary for me to mention some advanced concepts including statistical functionals, influence functions, Edgeworth and Cornish–Fisher expansions, and stationary stochastic processes. All these topics are well covered in the statistical literature on the bootstrap.

Since these concepts involve advanced probability and mathematics for a detailed description, I deliberately avoided such mathematical development to try to keep the text at a level for practitioners who do not have a strong mathematical background. Readers with an advanced mathematical background who might be curious about these concepts can refer to the references given throughout the chapter. In addition, Serfling (1980) is a good advanced text that provides much asymptotic statistical theory.

For the practitioner with less mathematical background, these details are not important. It is important to be aware that such theory exists to justify the use of the bootstrap in various contexts, but a deeper understanding is not necessary and for some it is not desirable.

This approach is really no different from the common practice, in elementary statistics texts, to mention the central limit theorem as justification for the use of the normal distribution to approximate the sampling distribution of sums or averages of random variables without providing any proof of the theorem such as Glivenko–Cantelli or Berry–Esseen or of related concepts such as convergence in distribution, triangular arrays, and Lindeberg–Feller conditions.