# Computational Statistics

First Edition

**Geof H. Givens and Jennifer A. Hoeting**

# *Preface*

This book covers most topics needed to develop a broad and thorough working knowledge of modern statistical computing and computational statistics. We seek to develop a practical understanding of how and why existing methods work, enabling readers to use modern statistical methods effectively. Since many new methods are built from components of existing techniques, our ultimate goal is to provide scientists with the tools they need to contribute new ideas to the field.

Achieving these goals requires familiarity with diverse topics in statistical computing, computational statistics, computer science, and numerical analysis. Our choice of topics reflects our view of what is central to this evolving field, and what will be interesting and useful for our readers. We pragmatically assigned priority to topics that can be of the most benefit to students and researchers most quickly.

Some topics we omitted represent important areas of past and present research in the field, but their priority here is lowered by the availability of high-quality software. For example, the generation of pseudo-random numbers is a classic topic, but one that we prefer to address by giving students reliable software. Some topics, such as numerical linear algebra, are on the borderline. Such topics are critical for many applications, yet good routines are generally available. In our judgment, the frequency with which one must shelve the routines and dig into the details of numerical linear algebra falls (barely) below the threshold we set for inclusion in this book. Among the classic topics we have chosen to cover are optimization and numerical integration. We include these because (i) they are cornerstones of frequentist and Bayesian inference; (ii) routine application of available software often fails for hard problems; and (iii) the

methods themselves are often secondary components of other statistical computing algorithms.

Our use of the adjective *modern* is potentially troublesome: there is no way that this book can cover all the latest, greatest techniques. We have not even tried. Some topics, such as heuristic search and Markov chain Monte Carlo, simply move too quickly. We have instead tried to offer a reasonably up-to-date survey of a broad portion of the field, while leaving room for diversions and esoterica. Some topics (e.g., principal curves and tabu search) are included simply because they are interesting and provide very different perspectives on familiar problems. Perhaps a future researcher may draw ideas from such topics to design a creative and effective new algorithm.

Our target audience includes graduate students in statistics and related fields, working statisticians, and quantitative empirical scientists in other fields. We hope such readers may use the book when applying standard methods and developing new methods.

The level of mathematics expected of the reader does not extend much beyond Taylor series and linear algebra. Breadth of mathematical training is more helpful than depth. Essential review is provided in Chapter 1. More advanced readers will find greater mathematical detail in the wide variety of high-quality books available on specific topics, many of which are referenced in the text. Other readers caring less about analytical details may prefer to focus on our descriptions of algorithms and examples.

The expected level of statistics is equivalent to that obtained by a graduate student in his or her first year of study of the theory of statistics and probability. An understanding of maximum likelihood methods, Bayesian methods, elementary asymptotic theory, Markov chains, and linear models is most important. Many of these topics are reviewed in Chapter 1.

With respect to computer programming, we find that good students can learn as they go. However, a working knowledge of a suitable language allows implementation of the ideas covered in this book to progress much more quickly. We have chosen to forgo any language-specific examples, algorithms, or coding. For those wishing to learn a language while they study this book, we recommend you choose a high-level, interactive package that permits the flexible design of graphical displays and includes supporting statistics and probability functions. At the time of writing, we recommend S-Plus, R, and MATLAB.[1] These are the sort of languages often used by researchers during the development of new statistical computing techniques, and are suitable for implementing all the methods we describe, except in some cases for problems of vast scope or complexity. Of course, lower-level languages such as C++ can also be used, and are favored for professional grade implementation of algorithms after researchers have refined the methodology.

Even adept computer programmers may have little understanding of how mathematics is carried out in the binary world of a computer. Mysterious problems with

---

[1] Websites for these software packages are www.insightful.com, www.r-project.org, and www.mathworks.com, respectively. R is free software reimplementing portions of S-Plus; the others are commercial.

full-rank matrices that appear noninvertible, integrals and likelihoods that vanish, numerical approximations that appear more precise than they really are, and other oddities are not unusual. While not dismissing the importance of computer arithmetic and numerically stable computation, we prefer to focus on the big picture of how algorithms work and to sweep under the rug some of the nitty-gritty numerical computation details.

The book is organized into three major parts: optimization (Chapters 2, 3, and 4), integration (Chapters 5, 6, 7, and 8), and smoothing (Chapters 10, 11, and 12). Chapter 9 adds another essential topic, the bootstrap. The chapters are written to stand independently, so a course can be built by selecting the topics one wishes to teach. For a one-semester course, our selection typically weights most heavily topics from Chapters 2, 5, 6, 7, 9, 10, and 11. With a leisurely pace or more thorough coverage, a shorter list of topics could still easily fill a semester course. There is sufficient material here to provide a thorough one-year course of study, notwithstanding any supplemental topics one might wish to teach.

A variety of homework problems are included at the end of each chapter. Some are straightforward, while others require the student to develop a thorough understanding of the model/method being used, to carefully (and perhaps cleverly) code a suitable technique, and to devote considerable attention to the interpretation of results.

The datasets discussed in the text and problems are available from the book website, *www.stat.colostate.edu/computationalstatistics*. The errata will also be found there. Responsibility for all errors lies with us.

Geof H. Givens and Jennifer A. Hoeting

*Fort Collins, Colorado*

# *Contents*