

Statistics and Computing

Series Editors:

J. Chambers

W. Eddy

W. Härdle

S. Sheather

L. Tierney

Statistics and Computing

Gentle: Numerical Linear Algebra for Applications in Statistics.

Gentle: Random Number Generation and Monte Carlo Methods.

Härdle/Klinke/Turlach: XploRe: An Interactive Statistical Computing Environment.

Krause/Olson: The Basics of S and S-PLUS.

Ó Ruanaidh/Fitzgerald: Numerical Bayesian Methods Applied to Signal Processing.

Pannatier: VARIOWIN: Software for Spatial Data Analysis in 2D.

Venables/Ripley: Modern Applied Statistics with S-PLUS, 2nd edition.

Lange: Numerical Analysis for Statisticians.

James E. Gentle

Numerical Linear Algebra for Applications in Statistics

With 21 Illustrations



Springer

James E. Gentle
Institute for Computational
Sciences and Informatics
George Mason University
Fairfax, VA 22030-4444
USA

Series Editors:

J. Chambers
Bell Labs, Lucent
Technologies
600 Mountain Ave.
Murray Hill, NJ 07974
USA

W. Eddy
Department of Statistics
Carnegie-Mellon University
Pittsburgh, PA 15213
USA

W. Härdle
Institut für Statistik und
Ökonometrie
Humboldt-Universität zu Berlin
Spandauer Str. 1
D-10178 Berlin
Germany

S. Sheather
Australian Graduate School
of Medicine
PO Box 1
Kensington
New South Wales 2033
Australia

L. Tierney
School of Statistics
University of Minnesota
Vincent Hall
Minneapolis, MN 55455
USA

Library of Congress Cataloging-in-Publication Data

Gentle, James E., 1943–

Numerical linear algebra for applications in statistics / James E.
Gentle.

p. cm. – (Statistics and computing)

Includes bibliographical references and indexes.

ISBN 978-1-4612-6842-0 ISBN 978-1-4612-0623-1 (eBook)

DOI 10.1007/978-1-4612-0623-1

1. Algebras, Linear. 2. Linear models (Statistics) I. Title.

II. Series.

QA184.G45 1998

512'.5--DC21

98-3959

Printed on acid-free paper.

©1998 Springer Science+Business Media New York
Originally published by Springer-Verlag New York, Inc. in 1998
Softcover reprint of the hardcover 1st edition 1998

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher Springer Science+Business Media, LLC, except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use of general descriptive names, trade names, trademarks, etc., in this publication, even if the former are not especially identified, is not to be taken as a sign that such names, as understood by the Trade Marks and Merchandise Marks Act, may accordingly be used freely by anyone.

Production managed by Lesley Poliner; manufacturing supervised by Joe Quatela.
Camera-ready copy prepared from the author's LaTeX files.

9 8 7 6 5 4 3 2 1

ISBN 978-1-4612-6842-0

To María

Preface

Linear algebra is one of the most important mathematical and computational tools in the sciences. While linear models and linear transformations are important in their own right, the basic role that linear algebra plays in nonlinear models, in optimization, and in other areas of statistics also makes an understanding of linear methods one of the most fundamental requirements for research in statistics or in application of statistical methods.

This book presents the aspects of numerical linear algebra that are important to statisticians. An understanding of numerical linear algebra requires both some basic understanding of computer arithmetic and some knowledge of the basic properties of vectors and matrices, so parts of the first two chapters are devoted to these fundamentals.

There are a number of books that cover linear algebra and matrix theory, with an emphasis on applications in statistics. In this book I state most of the propositions of linear algebra that are useful in data analysis, but in most cases I do not give the details of the proofs. Most of the proofs are rather straightforward and are available in the general references cited.

There are also a number of books on numerical linear algebra for numerical analysts. This book covers the computational aspects of vectors and matrices with an emphasis on statistical applications.

Books on statistical linear models also often address the computational issues; in this book the computations are central.

Throughout this book, the difference between an expression and a computing method is emphasized. For example, often we may write the solution to the linear system $Ax = b$ as $A^{-1}b$. Although this is the solution (so long as A is square and full rank), solving the linear system does not involve computing A^{-1} . We may write $A^{-1}b$, but we know we can compute the solution without inverting the matrix.

Chapter 1 provides some basic information on how data are stored and manipulated in a computer. Some of this material is rather tedious, but it is important to have a general understanding of computer arithmetic before considering computations for linear algebra. The impatient reader may skip or just skim Chapter 1, but the reader should be aware that the way the computer stores numbers and performs computations has far-reaching consequences. Computer arithmetic differs from ordinary arithmetic in many ways; for exam-

ple, computer arithmetic lacks associativity of addition and multiplication, and series often converge even when they are not supposed to. (On the computer, a straightforward evaluation of $\sum_{x=1}^{\infty} x$ converges!) Much of Chapter 1 is presented in the spirit of Forsythe (1970), “Pitfalls in computation, or why a math book isn’t enough.”

I emphasize the difference in the abstract number systems called the reals, \mathbb{R} , and the integers, \mathbb{Z} , from the computer number systems \mathbb{F} , the floating-point numbers, and \mathbb{I} , the fixed-point numbers. (Appendix A provides definitions for the notation.)

Chapter 1 also covers some of the fundamentals of algorithms, such as iterations, recursion, and convergence. It also discusses software development. Software issues are revisited in Chapter 5.

In Chapter 2, before considering numerical linear algebra, I begin with some basic properties of linear algebra. Except for some occasional asides, this material in the lengthy Section 2.1 is not in the area of *numerical* linear algebra. Knowledge of this material, however, is assumed in many places in the rest of the book. This section also includes topics, such as condition numbers, that would not be found in the usual books on “matrix algebra for statisticians”. Section 2.1 can be considered as a mini-reference source on vectors and matrices for statisticians.

In Section 2.2, building on the material from Chapter 1, I discuss how vectors and matrices are represented and manipulated in a computer.

Chapters 3 and 4 cover the basic computations for decomposing matrices, solving linear systems, and extracting eigenvalues and eigenvectors. These are the fundamental operations of numerical linear algebra. The need to solve linear systems arises much more often in statistics than does the need for eigenanalysis, and consequently Chapter 3 is longer and more detailed than Chapter 4.

Chapter 5 provides a brief introduction to software available for computations with linear systems. Some specific systems mentioned include the IMSL Libraries for Fortran and C, Matlab, and S-Plus. All of these systems are easy to use, and the best way to learn them is to begin using them for simple problems.

Throughout the text, the methods particularly useful in statistical computations are emphasized; and in Chapter 6, a few applications in statistics are discussed specifically.

Appendix A collects the notation used in this book. It is generally “standard” notation, but one thing the reader must become accustomed to is the lack of notational distinction between a vector and a scalar.

All vectors are “column” vectors, although I may write them as horizontal lists of their elements. (Whether vectors are “row” vectors or “column” vectors is generally only relevant for how we write expressions involving vector/matrix multiplication or partitions of matrices.)

I write algorithms in various ways, sometimes in a form that looks similar to Fortran or C, and sometimes as a list of numbered steps. I believe all of the descriptions used are straightforward and unambiguous.

One of the most significant developments in recent years, along with the general growth of computing power, has been the growth of data. It is now common to search through massive datasets and compute summary statistics from various items that may indicate relationships that were not previously recognized. The individual items or the relationships among them may not have been of primary interest when the data were originally collected. This process of prowling through the data is sometimes called *data mining* or *knowledge discovery in databases* (KDD). The objective is to discover characteristics of the data that may not be expected based on the existing theory.

Data must be stored; it must be transported; it must be sorted, searched, and otherwise rearranged; and computations must be performed on it. The size of the dataset largely determines whether these actions are feasible. For processing such massive datasets, the order of computations is a key measure of feasibility. Massive datasets make seemingly trivial improvements in algorithms important. The speedup of Strassen's method of an $O(n^3)$ algorithm to an $O(n^{2.81})$ algorithm, for example, becomes relevant for very large datasets. (Strassen's method is described on page 83.)

We now are beginning to encounter datasets of size 10^{10} and larger. We can quickly determine that a process whose computations are $O(n^2)$ cannot be reasonably contemplated for such massive datasets. If computations can be performed at a rate of 10^{12} per second (teraflop), it would take roughly three years to complete the computations. (A rough order of magnitude for quick "year" computations is $\pi \times 10^7$ seconds equals approximately one year.)

Advances in computer hardware continue to expand what is computationally feasible. It is interesting to note, however, that the order of time required for computations are determined by the problem to be solved and the algorithm to be used, not by the capabilities of the hardware. Advances in algorithm design have reduced the order of computations for many standard problems, while advances in hardware have not changed the order of the computations. Hardware advances have changed only the constant in the order of time.

This book has been assembled from lecture notes I have used in various courses in the computational and statistical sciences over the past few years. I believe the topics addressed in this book constitute the most important material for an introductory course in statistical computing, and should be covered in every such course. There are a number of additional topics that could be covered in a course in scientific computing, such as random number generation, optimization, and quadrature and solution of differential equations. Most of these topics require an understanding of computer arithmetic and of numerical linear algebra as covered in this book, so this book could serve as a basic reference for courses on other topics in statistical computing.

This book could also be used as a supplementary reference text for a course in linear regression that emphasizes the computational aspects.

The prerequisites for this text are minimal. Obviously some background in mathematics is necessary. Some background in statistics or data analysis and some level of scientific computer literacy are also required.

I do not use any particular software system in the book; but I do assume the ability to program in either Fortran or C, and the availability of either S-Plus, Matlab, or Maple. For some exercises the required software can be obtained from either `statlib` or `netlib` (see the bibliography).

Some exercises are Monte Carlo studies. I do not discuss Monte Carlo methods in this text, so the reader lacking background in that area may need to consult another reference in order to work those exercises.

I believe examples are very important. When I have taught this material, my lectures have consisted in large part of working through examples. Some of those examples have become exercises in the present text. The exercises should be considered an integral part of the book.

Acknowledgments

Over the years, I have benefited from associations with top-notch statisticians, numerical analysts, and computer scientists. There are far too many to acknowledge individually, but four stand out. My first real mentor — who was a giant in each of these areas — was Hoh Hartley. My education in statistical computing continued with Bill Kennedy, as I began to find my place in the broad field of statistics. During my years of producing software used by people all over the world, my collaborations with Tom Aird helped me better to understand some of the central issues of mathematical software. Finally, during the past several years, my understanding of computational statistics has been honed through my association with Ed Wegman. I thank these four people especially.

I thank my wife María, to whom this book is dedicated, for everything.

I used $\text{T}_{\text{E}}\text{X}$ via $\text{L}^{\text{A}}\text{T}_{\text{E}}\text{X}$ to write the book. I did all of the typing, programming, etc., myself (mostly early in the morning or late at night), so all mistakes are mine.

Material relating to courses I teach in the computational sciences is available over the World Wide Web at the URL,

<http://www.science.gmu.edu/>

Notes on this book, including errata, are available at

<http://www.science.gmu.edu/~jgentle/linbk/>

Notes on a larger book in computational statistics are available at

<http://www.science.gmu.edu/~jgentle/cmpstbk/>

Contents

Preface	vii
1 Computer Storage and Manipulation of Data	1
1.1 Digital Representation of Numeric Data	3
1.2 Computer Operations on Numeric Data	18
1.3 Numerical Algorithms and Analysis	26
Exercises	41
2 Basic Vector/Matrix Computations	47
2.1 Notation, Definitions, and Basic Properties	48
2.1.1 Operations on Vectors; Vector Spaces	48
2.1.2 Vectors and Matrices	52
2.1.3 Operations on Vectors and Matrices	55
2.1.4 Partitioned Matrices	58
2.1.5 Matrix Rank	59
2.1.6 Identity Matrices	60
2.1.7 Inverses	61
2.1.8 Linear Systems	62
2.1.9 Generalized Inverses	63
2.1.10 Other Special Vectors and Matrices	64
2.1.11 Eigenanalysis	67
2.1.12 Similarity Transformations	69
2.1.13 Norms	70
2.1.14 Matrix Norms	72
2.1.15 Orthogonal Transformations	74
2.1.16 Orthogonalization Transformations	74
2.1.17 Condition of Matrices	75
2.1.18 Matrix Derivatives	79
2.2 Computer Representations and Basic Operations	81
2.2.1 Computer Representation of Vectors and Matrices	81
2.2.2 Multiplication of Vectors and Matrices	82
Exercises	84

3	Solution of Linear Systems	87
3.1	Gaussian Elimination	87
3.2	Matrix Factorizations	92
3.2.1	LU and LDU Factorizations	92
3.2.2	Cholesky Factorization	93
3.2.3	QR Factorization	95
3.2.4	Householder Transformations (Reflections)	97
3.2.5	Givens Transformations (Rotations)	99
3.2.6	Gram-Schmidt Transformations	102
3.2.7	Singular Value Factorization	102
3.2.8	Choice of Direct Methods	103
3.3	Iterative Methods	103
3.3.1	The Gauss-Seidel Method with Successive Overrelaxation	103
3.3.2	Solution of Linear Systems as an Optimization Problem; Conjugate Gradient Methods	104
3.4	Numerical Accuracy	107
3.5	Iterative Refinement	109
3.6	Updating a Solution	109
3.7	Overdetermined Systems; Least Squares	111
3.7.1	Full Rank Coefficient Matrix	112
3.7.2	Coefficient Matrix Not of Full Rank	113
3.7.3	Updating a Solution to an Overdetermined System	114
3.8	Other Computations for Linear Systems	115
3.8.1	Rank Determination	115
3.8.2	Computing the Determinant	115
3.8.3	Computing the Condition Number	115
	Exercises	117
4	Computation of Eigenvectors and Eigenvalues and the Singular Value Decomposition	123
4.1	Power Method	124
4.2	Jacobi Method	126
4.3	QR Method for Eigenanalysis	129
4.4	Singular Value Decomposition	131
	Exercises	134
5	Software for Numerical Linear Algebra	137
5.1	Fortran and C	138
5.1.1	BLAS	140
5.1.2	Fortran and C Libraries	142
5.1.3	Fortran 90 and 95	146
5.2	Interactive Systems for Array Manipulation	148
5.2.1	Matlab	148
5.2.2	S, S-Plus	151

5.3	High-Performance Software	153
5.4	Test Data	155
	Exercises	157
6	Applications in Statistics	161
6.1	Fitting Linear Models with Data	162
6.2	Linear Models and Least Squares	163
6.2.1	The Normal Equations and the Sweep Operator	165
6.2.2	Linear Least Squares Subject to Linear Equality Constraints	166
6.2.3	Weighted Least Squares	166
6.2.4	Updating Linear Regression Statistics	167
6.2.5	Tests of Hypotheses	169
6.2.6	D-Optimal Designs	170
6.3	Ill-Conditioning in Statistical Applications	172
6.4	Testing the Rank of a Matrix	173
6.5	Stochastic Processes	175
	Exercises	176
	Appendices	183
	A Notation and Definitions	183
	B Solutions and Hints for Selected Exercises	191
	Bibliography	197
	Literature in Computational Statistics	198
	World Wide Web, News Groups, List Servers, and Bulletin Boards	199
	References	202
	Author Index	213
	Subject Index	217