Wiley Series in Probability and Statistics

Foundations of Linear and Generalized Linear Models

Alan Agresti

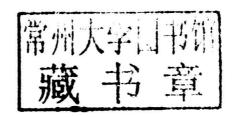
WILEY

Foundations of Linear and Generalized Linear Models

ALAN AGRESTI

Distinguished Professor Emeritus University of Florida Gainesville, FL

Visiting Professor Harvard University Cambridge, MA





Copyright © 2015 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey. Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470, or on the web at www.copyright.com. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at http://www.wiley.com/go/permission.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic formats. For more information about Wiley products, visit our web site at www.wiley.com.

Library of Congress Cataloging-in-Publication Data

Agresti, Alan, author.

Foundations of linear and generalized linear models / Alan Agresti. pages cm. – (Wiley series in probability and statistics)

pages cm. – (whey series in probability and st

Includes bibliographical references and index.

ISBN 978-1-118-73003-4 (hardback)

 Mathematical analysis-Foundations. 2. Linear models (Statistics) I. Title. OA299.8.A37 2015

003'.74-dc23

2014036543

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

Foundations of Linear and Generalized Linear Models

WILEY SERIES IN PROBABILITY AND STATISTICS

Established by WALTER A. SHEWHART and SAMUEL S. WILKS

Editors: David J. Balding, Noel A. C. Cressie, Garrett M. Fitzmaurice, Geof H. Givens, Harvey Goldstein, Geert Molenberghs, David W. Scott, Adrian F. M. Smith, Ruey S. Tsay, Sanford Weisberg Editors Emeriti: J. Stuart Hunter, Iain M. Johnstone, Joseph B. Kadane, Jozef L. Teugels

A complete list of the titles in this series appears at the end of this volume.

此为试读,需要完整PDF请访问: www.ertongbook.com

To my statistician friends in Europe



Preface

PURPOSE OF THIS BOOK

Why yet another book on linear models? Over the years, a multitude of books have already been written about this well-traveled topic, many of which provide more comprehensive presentations of linear modeling than this one attempts. My book is intended to present an overview of the key ideas and foundational results of linear and generalized linear models. I believe this overview approach will be useful for students who lack the time in their program for a more detailed study of the topic. This situation is increasingly common in Statistics and Biostatistics departments. As courses are added on recent influential developments (such as "big data," statistical learning, Monte Carlo methods, and application areas such as genetics and finance), programs struggle to keep room in their curriculum for courses that have traditionally been at the core of the field. Many departments no longer devote an entire year or more to courses about linear modeling.

Books such as those by Dobson and Barnett (2008), Fox (2008), and Madsen and Thyregod (2011) present fine overviews of both linear and generalized linear models. By contrast, my book has more emphasis on the theoretical foundations—showing how linear model fitting projects the data onto a model vector subspace and how orthogonal decompositions of the data yield information about effects, deriving likelihood equations and likelihood-based inference, and providing extensive references for historical developments and new methodology. In doing so, my book has less emphasis than some other books on practical issues of data analysis, such as model selection and checking. However, each chapter contains at least one section that applies the models presented in that chapter to a dataset, using R software. The book is not intended to be a primer on R software or on the myriad details relevant to statistical practice, however, so these examples are relatively simple ones that merely convey the basic concepts and spirit of model building.

The presentation of linear models for continuous responses in Chapters 1–3 has a geometrical rather than an algebraic emphasis. More comprehensive books on linear models that use a geometrical approach are the ones by Christensen (2011) and by

xii Preface

Seber and Lee (2003). The presentation of generalized linear models in Chapters 4–9 includes several sections that focus on discrete data. Some of this significantly abbreviates material from my book, *Categorical Data Analysis* (3rd ed., John Wiley & Sons, 2013). Broader overviews of generalized linear modeling include the classic book by McCullagh and Nelder (1989) and the more recent book by Aitkin et al. (2009). An excellent book on statistical modeling in an even more general sense is by Davison (2003).

USE AS A TEXTBOOK

This book can serve as a textbook for a one-semester or two-quarter course on linear and generalized linear models. It is intended for graduate students in the first or second year of Statistics and Biostatistics programs. It also can serve programs with a heavy focus on statistical modeling, such as econometrics and operations research. The book also should be useful to students in the social, biological, and environmental sciences who choose Statistics as their minor area of concentration.

As a prerequisite, the reader should be familiar with basic theory of statistics, such as presented by Casella and Berger (2001). Although not mandatory, it will be helpful if readers have at least some background in applied statistical modeling, including linear regression and ANOVA. I also assume some linear algebra background. In this book, I recall and briefly review fundamental statistical theory and matrix algebra results where they are used. This contrasts with the approach in many books on linear models of having several chapters on matrix algebra and distribution theory before presenting the main results on linear models. Readers wanting to improve their knowledge of matrix algebra can find on the Web (e.g., with a Google search of "review of matrix algebra") overviews that provide more than enough background for reading this book. Also helpful as background for Chapters 1-3 on linear models are online lectures, such as the MIT linear algebra lectures by G. Strang at http://ocw.mit.edu/courses/mathematics on topics such as vector spaces, column space and null space, independence and a basis, inverses, orthogonality, projections and least squares, eigenvalues and eigenvectors, and symmetric and idempotent matrices. By not including separate chapters on matrix algebra and distribution theory, I hope instructors will be able to cover most of the book in a single semester or in a pair of quarters.

Each chapter contains exercises for students to practice and extend the theory and methods and also to help assimilate the material by analyzing data. Complete data files for the text examples and exercises are available at the text website, http://www.stat.ufl.edu/~aa/glm/data/. Appendix A contains supplementary data analysis exercises that are not tied to any particular chapter. Appendix B contains solution outlines and hints for some of the exercises.

I emphasize that this book is not intended to be a complete overview of linear and generalized linear modeling. Some important classes of models are beyond its scope; examples are transition (e.g., Markov) models and survival (time-to-event) models. I intend merely for the book to be an overview of the *foundations* of this subject—that is, core material that should be part of the background of any statistical scientist. I

PREFACE XIII

invite readers to use it as a stepping stone to reading more specialized books that focus on recent advances and extensions of the models presented here.

ACKNOWLEDGMENTS

This book evolved from a one-semester course that I was invited to develop and teach as a visiting professor for the Statistics Department at Harvard University in the fall terms of 2011–2014. That course covers most of the material in Chapters 1–9. My grateful thanks to Xiao-Li Meng (then chair of the department) for inviting me to teach this course, and likewise thanks to Dave Harrington for extending this invitation through 2014. (The book's front cover, showing the Zakim bridge in Boston, reflects the Boston-area origins of this book.) Special thanks to Dave Hoaglin, who besides being a noted statistician and highly published book author, has wonderful editing skills. Dave gave me detailed and helpful comments and suggestions for my working versions of all the chapters, both for the statistical issues and the expository presentation. He also found many errors that otherwise would have found their way into print!

Thanks also to David Hitchcock, who kindly read the entire manuscript and made numerous helpful suggestions, as did Maria Kateri and Thomas Kneib for a few chapters. Hani Doss kindly shared his fine course notes on linear models (Doss 2010) when I was organizing my own thoughts about how to present the foundations of linear models in only two chapters. Thanks to Regina Dittrich for checking the R code and pointing out errors. I owe thanks also to several friends and colleagues who provided comments or datasets or other help, including Pat Altham, Alessandra Brazzale, Jane Brockmann, Phil Brown, Brian Caffo, Leena Choi, Guido Consonni, Brent Coull, Anthony Davison, Kimberly Dibble, Anna Gottard, Ralitza Gueorguieva, Alessandra Guglielmi, Jarrod Hadfield, Rebecca Hale, Don Hedeker, Georg Heinze, Jon Hennessy, Harry Khamis, Eunhee Kim, Joseph Lang, Ramon Littell, I-Ming Liu, Brian Marx, Clint Moore, Bhramar Mukherjee, Dan Nettleton, Keramat Nourijelyani, Donald Pierce, Penelope Pooler, Euijung Ryu, Michael Schemper, Cristiano Varin, Larry Winner, and Lo-Hua Yuan. James Booth, Gianfranco Lovison, and Brett Presnell have generously shared materials over the years dealing with generalized linear models. Alex Blocker, Jon Bischof, Jon Hennessy, and Guillaume Basse were outstanding and very helpful teaching assistants for my Harvard Statistics 244 course, and Jon Hennessy contributed solutions to many exercises from which I extracted material at the end of this book. Thanks to students in that course for their comments about the manuscript. Finally, thanks to my wife Jacki Levine for encouraging me to spend the terms visiting Harvard and for support of all kinds, including helpful advice in the early planning stages of this book.

ALAN AGRESTI

Contents

Preface

1	Intro	oduction to Linear and Generalized Linear Models	1
	1.1	Components of a Generalized Linear Model, 2	
	1.2	Quantitative/Qualitative Explanatory Variables and Interpreting Effects,	6
	1.3	Model Matrices and Model Vector Spaces, 10	
	1.4	Identifiability and Estimability, 13	
	1.5	Example: Using Software to Fit a GLM, 15	
		oter Notes, 20	
		cises, 21	
2	Line	ear Models: Least Squares Theory	26
	2.1	Least Squares Model Fitting, 27	
	2.2	Projections of Data Onto Model Spaces, 33	
	2.3	Linear Model Examples: Projections and SS Decompositions, 41	
	2.4	Summarizing Variability in a Linear Model, 49	
	2.5	Residuals, Leverage, and Influence, 56	
	2.6	Example: Summarizing the Fit of a Linear Model, 62	
	2.7	Optimality of Least Squares and Generalized Least Squares, 67	
	Chap	pter Notes, 71	
	Exer	rcises, 71	
3	Normal Linear Models: Statistical Inference		80
	3.1	Distribution Theory for Normal Variates, 81	
	3.2	Significance Tests for Normal Linear Models, 86	
	3.3	Confidence Intervals and Prediction Intervals for Normal Linear Models, 95	

xi

viii CONTENTS

3.4 Example: Normal Linear Model Inference, 99

		Multiple Comparisons: Bonferroni, Tukey, and FDR Methods, 107				
	^	Chapter Notes, 111				
	Exercises, 112					
1	Gen	Generalized Linear Models: Model Fitting and Inference				
	4.1	Exponential Dispersion Family Distributions for a GLM, 120				
	4.2	Likelihood and Asymptotic Distributions for GLMs, 123				
	4.3	Likelihood-Ratio/Wald/Score Methods of Inference for GLM Parameters, 128				
	4.4	Deviance of a GLM, Model Comparison, and Model Checking, 132				
	4.5	Fitting Generalized Linear Models, 138				
	4.6	Selecting Explanatory Variables for a GLM, 143				
	4.7	Example: Building a GLM, 149				
	App	endix: GLM Analogs of Orthogonality Results for Linear Models, 156				
	Chap	pter Notes, 158				
	Exer	rcises, 159				
5	Models for Binary Data		165			
	5.1	Link Functions for Binary Data, 165				
	5.2	Logistic Regression: Properties and Interpretations, 168				
	5.3	Inference About Parameters of Logistic Regression Models, 172				
	5.4	Logistic Regression Model Fitting, 176				
	5.5	Deviance and Goodness of Fit for Binary GLMs, 179				
	5.6	Probit and Complementary Log-Log Models, 183				
	5.7	Examples: Binary Data Modeling, 186				
	Chaj	pter Notes, 193				
	Exercises, 194					
5	Mul	tinomial Response Models	202			
	6.1	Nominal Responses: Baseline-Category Logit Models, 203				
	6.2	Ordinal Responses: Cumulative Logit and Probit Models, 209				
	6.3	Examples: Nominal and Ordinal Responses, 216				
	Chapter Notes, 223					
	Exer	reises, 223				
7	Mod	Models for Count Data				
	7.1	Poisson GLMs for Counts and Rates, 229				
	7.2	Poisson/Multinomial Models for Contingency Tables, 235				

ix CONTENTS

8

9

	7.3	Negative Binomial GLMS, 247				
	7.4	Models for Zero-Inflated Data, 250				
	7.5	Example: Modeling Count Data, 254				
	Chapt	Chapter Notes, 259				
	Exerc	ises, 260				
8	Quas	i-Likelihood Methods	268			
	8.1	Variance Inflation for Overdispersed Poisson and Binomial GLMs, 26	9			
	8.2	Beta-Binomial Models and Quasi-Likelihood Alternatives, 272				
	8.3	Quasi-Likelihood and Model Misspecification, 278				
	Chapt	ter Notes, 282				
	Exerc	ises, 282				
9	Mode	eling Correlated Responses	286			
	9.1	Marginal Models and Models with Random Effects, 287				
	9.2	Normal Linear Mixed Models, 294				
		Fitting and Prediction for Normal Linear Mixed Models, 302				
		Binomial and Poisson GLMMs, 307				
		GLMM Fitting, Inference, and Prediction, 311				
		Marginal Modeling and Generalized Estimating Equations, 314				
		Example: Modeling Correlated Survey Responses, 319				
	Chapter Notes, 322					
	Exerc	tises, 324				
10	Bay	esian Linear and Generalized Linear Modeling	333			
	10.1					
	10.2					
	10.3					
	10.4					
		pter Notes, 357				
	Exe	rcises, 359				
11	Exte	ensions of Generalized Linear Models	364			
	11.1		365			
	11.2					
	11.3	Smoothing, Generalized Additive Models, and Other GLM Extensions, 378				
		pter Notes, 386				
	Exe	rcises, 388				

Appendix A	Supplemental Data Analysis Exercises	391
Appendix B	Solution Outlines for Selected Exercises	396
References		410
Author Index	427	
Example Ind	433	
Subject Index		
Website		

Data sets for the book are at www.stat.ufl.edu/~aa/glm/data

X

CONTENTS

Introduction to Linear and Generalized Linear Models

This is a book about *linear models* and *generalized linear models*. As the names suggest, the linear model is a special case of the generalized linear model. In this first chapter, we define generalized linear models, and in doing so we also introduce the linear model.

Chapter 2 and 3 focus on the linear model. Chapter 2 introduces the *least squares* method for fitting the model, and Chapter 3 presents statistical inference under the assumption of a *normal* distribution for the response variable. Chapter 4 presents analogous model-fitting and inferential results for the generalized linear model. This generalization enables us to model non-normal responses, such as categorical data and count data.

The remainder of the book presents the most important generalized linear models. Chapter 5 focuses on models that assume a binomial distribution for the response variable. These apply to binary data, such as "success" and "failure" for possible outcomes in a medical trial or "favor" and "oppose" for possible responses in a sample survey. Chapter 6 extends the models to multicategory responses, assuming a multinomial distribution. Chapter 7 introduces models that assume a Poisson or negative binomial distribution for the response variable. These apply to count data, such as observations in a health survey on the number of respondent visits in the past year to a doctor. Chapter 8 presents ways of weakening distributional assumptions in generalized linear models, introducing quasi-likelihood methods that merely focus on the mean and variance of the response distribution. Chapters 1-8 assume independent observations. Chapter 9 generalizes the models further to permit correlated observations, such as in handling multivariate responses. Chapters 1-9 use the traditional frequentist approach to statistical inference, assuming probability distributions for the response variables but treating model parameters as fixed, unknown values. Chapter 10 presents the Bayesian approach for linear models and generalized linear models, which treats the model parameters as random variables having their own distributions. The final chapter introduces extensions of the models that handle more complex situations, such as *high-dimensional* settings in which models have enormous numbers of parameters.

1.1 COMPONENTS OF A GENERALIZED LINEAR MODEL

The ordinary linear regression model uses linearity to describe the relationship between the mean of the response variable and a set of explanatory variables, with inference assuming that the response distribution is normal. *Generalized linear models* (GLMs) extend standard linear regression models to encompass non-normal response distributions and possibly nonlinear functions of the mean. They have three components.

- Random component: This specifies the response variable y and its probability distribution. The observations $y = (y_1, \dots, y_n)^T$ on that distribution are treated as independent.
- Linear predictor: For a parameter vector $\boldsymbol{\beta} = (\beta_1, \beta_2, \dots, \beta_p)^T$ and a $n \times p$ model matrix X that contains values of p explanatory variables for the n observations, the linear predictor is $X\boldsymbol{\beta}$.
- *Link function*: This is a function *g* applied to each component of *E*(*y*) that relates it to the linear predictor,

$$g[E(y)] = X\beta$$
.

Next we present more detail about each component of a GLM.

1.1.1 Random Component of a GLM

The random component of a GLM consists of a response variable y with independent observations (y_1, \ldots, y_n) having probability density or mass function for a distribution in the exponential family. In Chapter 4 we review this family of distributions, which has several appealing properties. For example, $\sum_i y_i$ is a sufficient statistic for its parameter, and regularity conditions (such as differentiation passing under an integral sign) are satisfied for derivations of properties such as optimal large-sample performance of maximum likelihood (ML) estimators.

By restricting GLMs to exponential family distributions, we obtain general expressions for the model likelihood equations, the asymptotic distributions of estimators for model parameters, and an algorithm for fitting the models. For now, it suffices to say that the distributions most commonly used in Statistics, such as the normal, binomial, and Poisson, are exponential family distributions.

¹The superscript T on a vector or matrix denotes the transpose; for example, here y is a column vector. Our notation makes no distinction between random variables and their observed values; this is generally clear from the context.