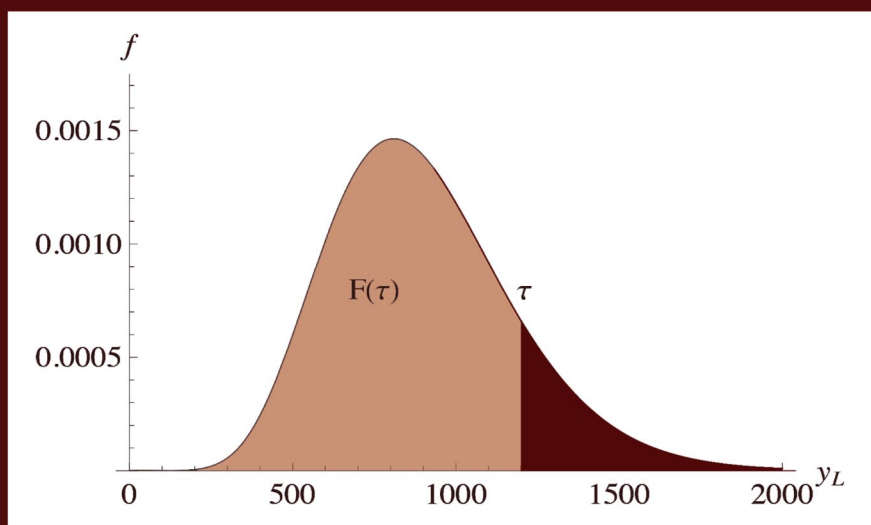


Chapman & Hall/CRC
Statistics in the Social and Behavioral Sciences Series

Generalized Linear Models for Categorical and Continuous Limited Dependent Variables



Michael Smithson
Edgar C. Merkle



CRC Press
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Generalized Linear Models for Categorical and Continuous Limited Dependent Variables

Chapman & Hall/CRC
Statistics in the Social and Behavioral Sciences Series

Series Editors

Jeff Gill

Washington University, USA

Steven Heeringa

University of Michigan, USA

Wim van der Linden

CTB/McGraw-Hill, USA

J. Scott Long

Indiana University, USA

Tom Snijders

Oxford University, UK
University of Groningen, UK

Aims and scope

Large and complex datasets are becoming prevalent in the social and behavioral sciences and statistical methods are crucial for the analysis and interpretation of such data. This series aims to capture new developments in statistical methodology with particular relevance to applications in the social and behavioral sciences. It seeks to promote appropriate use of statistical, econometric and psychometric methods in these applied sciences by publishing a broad range of reference works, textbooks and handbooks.

The scope of the series is wide, including applications of statistical methodology in sociology, psychology, economics, education, marketing research, political science, criminology, public policy, demography, survey methodology and official statistics. The titles included in the series are designed to appeal to applied statisticians, as well as students, researchers and practitioners from the above disciplines. The inclusion of real examples and case studies is therefore essential.

Published Titles

Analysis of Multivariate Social Science Data, Second Edition

David J. Bartholomew, Fiona Steele, Irini Moustaki, and Jane I. Galbraith

Applied Survey Data Analysis

Steven G. Heeringa, Brady T. West, and Patricia A. Berglund

Bayesian Methods: A Social and Behavioral Sciences Approach, Second Edition

Jeff Gill

Foundations of Factor Analysis, Second Edition

Stanley A. Mulaik

Incomplete Categorical Data Design: Non-Randomized Response Techniques for Sensitive Questions in Surveys

Guo-Liang Tian and Man-Lai Tang

Informative Hypotheses: Theory and Practice for Behavioral and Social Scientists

Herbert Hoijtink

Latent Markov Models for Longitudinal Data

Francesco Bartolucci, Alessio Farcomeni, and Fulvia Pennoni

Linear Causal Modeling with Structural Equations

Stanley A. Mulaik

Multiple Correspondence Analysis and Related Methods

Michael Greenacre and Jorg Blasius

Multivariable Modeling and Multivariate Analysis for the Behavioral Sciences

Brian S. Everitt

Statistical Test Theory for the Behavioral Sciences

Dato N. M. de Gruijter and Leo J. Th. van der Kamp

Generalized Linear Models for Categorical and Continuous Limited Dependent Variables

Michael Smithson and Edgar C. Merkle

This page intentionally left blank

Chapman & Hall/CRC
Statistics in the Social and Behavioral Sciences Series

Generalized Linear Models for Categorical and Continuous Limited Dependent Variables

Michael Smithson

The Australian National University, Canberra, Australia

Edgar C. Merkle

University of Missouri, Columbia, Missouri, USA



CRC Press

Taylor & Francis Group

Boca Raton London New York

CRC Press is an imprint of the
Taylor & Francis Group, an **informa** business
A CHAPMAN & HALL BOOK

CRC Press
Taylor & Francis Group
6000 Broken Sound Parkway NW, Suite 300
Boca Raton, FL 33487-2742

© 2014 by Taylor & Francis Group, LLC
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works
Version Date: 20130614

International Standard Book Number-13: 978-1-4665-5175-6 (eBook - PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access www.copyright.com (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

Trademark Notice: Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

Visit the Taylor & Francis Web site at
<http://www.taylorandfrancis.com>

and the CRC Press Web site at
<http://www.crcpress.com>

Dedication

Michael Smithson: To the memory of my father, James Edward Smithson (1921–1961). I think he would have liked this book.

Ed Merkle: To Lila and Grant.

This page intentionally left blank

Contents

Preface	xiii
List of Figures	xvii
List of Tables	xix
Notation	xxi
About the Authors	xxiii
1 Introduction and Overview	1
1.1 The Nature of Limited Dependent Variables	1
1.2 Overview of GLMs	2
1.2.1 Definition	2
1.2.2 Extensions	4
1.3 Estimation Methods and Model Evaluation	6
1.3.1 Model Evaluation and Diagnosis	6
1.3.2 Model Selection and Interpretation Issues	9
1.4 Organization of This Book	12
I Discrete Variables	15
2 Binary Variables	17
2.1 Logistic Regression	18
2.2 The Binomial GLM	20
2.2.1 Latent Variable Interpretation	22
2.2.2 Interpretation of Coefficients	22
2.2.3 Example	25
2.2.4 Extension to $n > 1$	27
2.2.5 Alternative Link Functions	30
2.3 Estimation Methods and Issues	32
2.3.1 Model Evaluation and Diagnostics	32
2.3.2 Overdispersion	38
2.3.3 Relationships to Other Models	39
2.4 Analyses in R and Stata	40

2.4.1	Analyses in R	40
2.4.2	Analyses in Stata	45
2.5	Exercises	49
3	Nominal Polytomous Variables	51
3.1	Multinomial Logit Model	51
3.2	Conditional Logit and Choice Models	58
3.3	Multinomial Processing Tree Models	61
3.4	Estimation Methods and Model Evaluation	66
3.4.1	Estimation Methods and Model Comparison	66
3.4.2	Model Evaluation and Diagnosis	67
3.5	Analyses in R and Stata	71
3.5.1	Analyses in R	71
3.5.2	Analyses in Stata	77
3.6	Exercises	80
4	Ordinal Categorical Variables	81
4.1	Modeling Ordinal Variables: Common Practice versus Best Practice	81
4.2	Ordinal Model Alternatives	82
4.2.1	The Proportional Odds Assumption	82
4.2.2	Modeling Relative Probabilities	83
4.3	Cumulative Models	84
4.3.1	The Proportional Odds Model	84
4.3.2	Example	86
4.4	Adjacent Models	89
4.4.1	The Adjacent Categories Model	89
4.4.2	Example	90
4.5	Stage Models	91
4.5.1	The Continuation Ratio Model	92
4.5.2	Example	94
4.6	Estimation Methods and Issues	97
4.6.1	Model Choice	98
4.6.2	Model Diagnostics	98
4.7	Analyses in R and Stata	102
4.7.1	Analyses in R	102
4.7.2	Analyses in Stata	107
4.8	Exercises	111
5	Count Variables	113
5.1	Distributions for Count Data	113
5.2	Poisson Regression Models	116
5.2.1	Model Definition	116
5.2.2	Example	118
5.2.3	Exposure	119
5.2.4	Overdispersion and Quasi-Poisson Models	120

CONTENTS	xi
----------	----

5.3	Negative Binomial Models	123
5.3.1	Model Definition	123
5.3.2	Example	125
5.4	Truncated and Censored Models	125
5.5	Zero-Inflated and Hurdle Models	127
5.5.1	Hurdle Models	127
5.5.2	Zero-Inflated Models	130
5.6	Estimation Methods and Issues	133
5.6.1	Negative Binomial Model Estimation	133
5.6.2	Model Diagnostics	133
5.7	Analyses in R and Stata	137
5.7.1	Analyses in R	138
5.7.2	Analyses in Stata	144
5.8	Exercises	148

II Continuous Variables 151

6	Doubly Bounded Continuous Variables	153
6.1	Doubly Bounded versus Censored	153
6.2	The beta GLM	153
6.3	Modeling Location and Dispersion	159
6.3.1	Judged Probability of Guilt	160
6.3.2	Reading Accuracy for Dyslexic and Non-Dyslexic Readers	162
6.3.3	Model Comparison	164
6.4	Estimation Methods and Issues	166
6.4.1	Estimator Bias	168
6.4.2	Model Diagnostics	171
6.5	Zero- and One-Inflated Models	176
6.6	Finite Mixture Models	178
6.6.1	Car Dealership Example	180
6.7	Analyses in R and Stata	182
6.7.1	Analyses in R	182
6.7.2	Analyses in Stata	186
6.8	Exercises	191
7	Censoring and Truncation	193
7.1	Models for Censored and Truncated Variables	193
7.1.1	Tobit Models	197
7.2	Non-Gaussian Censored Regression	203
7.3	Estimation Methods, Model Comparison, and Diagnostics	208
7.4	Extensions of Censored Regression Models	211
7.4.1	Proportional Hazard and Proportional Odds Models	212
7.4.2	Double and Interval Censoring	214
7.4.3	Censored Quantile Regression	220

7.5	Analyses in R and Stata	222
7.5.1	Analyses in R	222
7.5.2	Analyses in Stata	228
7.6	Exercises	232
8	Extensions	235
8.1	Extensions and Generalizations	235
8.2	Multilevel Models	236
8.2.1	Multilevel Binary Logistic Regression	236
8.2.2	Multilevel Count Models	239
8.2.3	Multilevel Beta Regression	241
8.3	Bayesian Estimation	245
8.3.1	Bayesian Binomial GLM	247
8.3.2	Bayesian Beta Regression	250
8.3.3	Modeling Random Sums	253
8.4	Evaluating Relative Importance of Predictors in GLMs	256
	References	261
	Author Index	275
	Subject Index	281

Preface

This book is devoted to the “other” kinds of dependent variables than those for which linear regression is appropriate. These include binary, polytomous nominal, categorical ordinal, counted, interval-valued, bounded continuous, censored, and truncated variables. We argue that these dependent variables are, if anything, more common throughout the human sciences than the kind that suit linear regression. Readers acquainted with the literature on such variables will have noticed the similarity between the titles of this book and of the pioneering textbook by Long (1997). Long’s book was eagerly acquired by the first author when it came out and proved an excellent source and guide over the years for both students and colleagues. Our book updates his book on topics they have in common, primarily regarding advances in special cases or extensions of models, estimation methods, model diagnostics, and of course software. Although the past two decades have seen many excellent books published on these topics, most of them are devoted to one or another specific subset of the topics. Ours is a broader but unified coverage in which we attempt to integrate the concepts and ideas shared across models and types of data, especially regarding conceptual links between discrete and continuous limited dependent variables.

At several points we bring together material that heretofore has been scattered across the literature in journal articles, book chapters, conference proceedings, software package documentation files, and blogs. Topics in our book not covered in Long’s include bounded continuous variables, a greater variety of boundary-inflated models, and methods for modeling heteroscedasticity. All of the dependent variables we consider have boundaries of some kind, be they due to categorical distinctions or bounds on a continuum. The distinctions among different kinds of bounds and how to incorporate them into statistical models are fairly challenging issues, and not much guidance is available in the literature. For example, we have observed that researchers can become confused about whether boundary observations on a variable should be regarded as accurate scores or censored values. Throughout the book we guide the reader to appropriate models on the basis of whether the bounds are inherent in a construct or variable, or imposed (e.g., by censoring or truncation). Likewise, although both the concepts and software are available for dealing with heteroscedasticity, it remains a relatively neglected topic in the applied statistical literature despite its considerable importance. Heteroscedasticity is especially relevant for the kinds of dependent variables we deal with here, both because it can frequently arise in the data and because some models for these variables are inherently heteroscedastic. We therefore treat both kinds of heteroscedasticity: Unconditional in the sense that it is

due to the bounded nature of the construct or variable, and conditional on values or states of independent variables.

Wherever possible, we have illustrated concepts, models, and techniques with real or realistic datasets and demonstrations in R (R Development Core Team, 2013) and Stata. Each substantive chapter also has several exercises at the end. Both illustrations and exercises are intended to help readers to build conceptual understanding and fluency in using these techniques.

Data and Software

We illustrate the models and methods in this book using both R (R Development Core Team, 2013) and Stata software. We elected to use these pieces of software through a combination of personal preference, popularity, and access. The data files used in this book are all freely available; to obtain them, R users can install the `smdata` package that is freely available on CRAN. The installation, followed by loading the package and loading datasets, can be completed with the following R commands.

```
## Install package
install.packages("smdata")

## Load package
library("smdata")

## Load, e.g., the email dataset within smdata
data("email")
```

On the topic of R, we also note that Thompson (2009) has written a valuable manual that describes the use of R for categorical data analysis; it can be freely obtained at <https://home.comcast.net/~lthompson221/Splushdiscrete2.pdf>.

Stata users and others can obtain the data from <https://dl.dropbox.com/u/1857674/SmiMerBook/SmiMerDoc.html>, which includes the data in both Stata format and csv format. The page also includes details about each data file and some extra code. Finally, Stata users will find the book on categorical and limited dependent variables by Long and Freese (2006) very useful.

Acknowledgments

We are indebted to numerous people who have helped to improve this book. We received substantial and very useful feedback on chapter drafts, illustrative examples, and supplementary material from Yiyun Shou, Jay Verkuilen, four anonymous reviewers, and the Australian National University psychology Honours and graduate students participating in statistical workshops during 2012. We also received useful feedback and discussion from graduate students in Ed's GLM courses at Wichita State University and from colleagues in the department of Psychological Sciences at the University of Missouri. Roger Koenker and Achim Zeileis provided prompt re-

sponses to our questions regarding the nuances of various packages and code, and the maintainers of R and CRAN aided us with the generous contribution of their time. Likewise, Bill Rising at Stata, as part of their author support program, made valuable suggestions for improving our attempts at Stata code. On the data side, Lukas Hulsey provided the medical treatment preference data that were used in Chapter 4, Justin Owens provided the eye-tracking data that were used in Chapter 5, Michael Gurr's Honours thesis data was used in Chapter 6, Ken Mavor contributed data used in Chapter 7, and Stephen Tang contributed data used in Chapter 8. At Chapman and Hall, John Kimmel gave us expert editorial advice, Shashi Kumar provided essential \LaTeX support, project coordinators Kathryn Everett and Kate Gallo guided us throughout the book's gestation, and project editor Marsha Hecht provided invaluable production assistance. Without them this book would not have been possible, let alone presentable. Of course, any errors or flaws remaining in this book are our responsibility.

The book was typeset by the authors in \LaTeX using the `chapmono` class. We also made extensive use of the `apacite` package for citations and indexes. We are indebted to the developers and maintainers of this free software.

Finally, we each owe ineradicable debts to our families, most especially to our wives, Susan and Victoria. Susan, yet again, tolerated the late nights, working weekends, obsessional preoccupation, and spousal neglect that unfortunately accompanies my book-writing (and she knows I can't even promise that this is the last one!). Victoria also tolerated these things, independently of book-writing to some degree. She additionally had a baby during the middle of book-writing, making for a very busy time. She continued to keep good spirits and flutter her eyelashes during this busy time, however, and I am grateful.

This page intentionally left blank

List of Figures

2.1	The binomial ($n = 10, p = .5$) distribution.	19
2.2	Relationship between $\text{logit}(p)$ and p .	21
2.3	Relationship between high school GPA and probability of passing calculus.	23
2.4	Scatter plots of “hours worked per week” and “year in school,” versus whether or not the individual skipped school.	26
2.5	Histogram of the number of days that students skipped school, out of the past 30.	29
2.6	Predictor variables versus Pearson residuals from the binary model fitted to the school-skipping data.	33
2.7	ROC curve for the binary logistic regression model fitted to the school-skipping data.	36
3.1	Age effect on odds ratios of three modes relative to injection.	54
3.2	Plot of the ratio of sex odds ratios coefficients from Table 3.5.	57
3.3	Choice MPT structure.	63
3.4	Simplified perceptual identification MPT structure.	65
3.5	Residuals for model of cocaine usage method with age as predictor.	69
4.1	Effects plot of the proportional odds model, as fit to the Hulsey data.	88
4.2	Box plots of weekly time spent emailing by marital status.	95
4.3	Stage model predicted probabilities of remaining at one’s current marital status.	97
5.1	Poisson distributions.	114
5.2	Negative binomial distributions.	115
5.3	Illustration of the logarithmic link function.	117
5.4	Box plots of the Owens et al. data by condition.	118
5.5	Histogram of days missed due to emotions, nerves, or mental health.	129
5.6	Comparison of observed and predicted proportions of Owens et al. fixation counts from the negative binomial model with offset.	134
5.7	Pearson residuals from the negative binomial model with offset, fitted to the Owens et al. data.	135
6.1	Beta distributions.	154
6.2	Reading accuracy scores for dyslexics and controls.	162

6.3	Reading accuracy scores for dyslexics and controls.	164
6.4	Pearson residuals for Models 1 and 2.	172
6.5	Four kinds of residuals.	173
6.6	Influence on location submodel coefficients.	174
6.7	Influence on precision submodel coefficients.	175
6.8	Generalized leverage.	176
6.9	Probability car bought from Carlos.	180
7.1	Latent response time distribution.	194
7.2	Truncated and censored pdfs.	195
7.3	Depression score predictions from linear regression and Tobit models.	199
7.4	Predicted values for y^* , $y y < \tau$, and censored y .	202
7.5	Weibull distributions with $\theta = 0.7, 1.3, 2, 3$, and $\lambda = 1$.	205
7.6	Example 7.2: Response times: Intuition-primed versus deliberation-primed conditions.	206
7.7	Response time Q-Q plots for the Weibull and log-normal models.	209
7.8	Example 7.2: Raw versus deviance residuals for response time Tobit and Weibull models.	210
7.9	Pearson residuals for the heteroscedastic Tobit model.	210
7.10	dfbetas residuals for Tobit location submodel.	211
7.11	dfbetas residuals for Tobit dispersion submodel.	212
7.12	Response time data empirical versus CPHM fitted survival curves.	214
7.13	Example 7.3: Pro-euthanasia attitudes predicted by Christian identity.	215
7.14	Example 7.4: Prediction lines for the marks-based, interval-censored, and midpoint-based regressions.	217
7.15	Lab scores: 25 th , 50 th , and 75 th quantile regressions compared with Tobit regression.	221
8.1	Multilevel binomial example, box plots of estimated subject effects by condition.	238
8.2	Multilevel binomial example, histogram of estimated item effects.	239
8.3	Joint posterior distributions of β_0 and β_2 , for models where school year is uncentered (top) and centered (bottom).	248
8.4	Posterior predictive distributions.	249
8.5	Chain mixing and density for δ_2 parameter.	251
8.6	QQ Plots for the three verdicts' posterior distributions.	252
8.7	95% credible region for the "guilty" verdicts QQ plot.	253
8.8	cdfs of duration sums versus posterior sums distribution.	255

List of Tables

2.1	Contingency table of school skipping versus year in school.	26
2.2	School skipping example, parameter estimates and standard errors of the binomial logistic model with $n = 30$ versus the binary logistic model.	29
2.3	Test accuracy example, parameter estimates and standard errors of the binomial logistic model with $n = 30$.	30
3.1	Sex by cocaine use method.	52
3.2	MNL model predicting cocaine usage method from sex.	53
3.3	MNL model predicting cocaine usage method from age.	53
3.4	Cocaine ingestion method by sex by race.	55
3.5	MNL model predicting cocaine ingestion method from sex by race.	56
3.6	Choice by initiated gaze.	60
3.7	CL model predicting choice from initiated gaze.	60
3.8	Transportation methods.	62
3.9	MPT coefficients.	63
3.10	Observed and fitted probabilities and summed residuals of CL models.	68
4.1	Proportion of participants in the Hulsey study making each rating, by condition. Cumulative probabilities are in parentheses.	86
4.2	Cumulative proportional odds model estimates for the Hulsey data, with two-tailed p -values.	87
4.3	Cumulative proportional odds model predictions for the Hulsey data.	88
4.4	Adjacent categories model estimates for the Hulsey data.	90
4.5	Adjacent categories model predictions for the Hulsey data.	91
4.6	Continuation ratio model estimates for the GSS data.	96
4.7	Stage model estimates (with no proportional odds assumption) for the GSS data.	96
4.8	Residuals for the cumulative proportional odds model fit to the Hulsey data.	99
5.1	Comparison of slope estimates from six count regression models.	125
5.2	Hurdle model estimates from fit to “work days missed” data, using a truncated Poisson distribution for the count part of the model.	130

5.3	Hurdle model estimates from a fit to the “work days missed” data, using a truncated negative binomial distribution for the count part of the model.	131
5.4	Zero-inflated negative binomial model estimates from a fit to the “work days missed” data.	132
5.5	Comparison of fitted zero and hurdle models’ log-likelihoods.	133
5.6	Overview of nested count regression models.	136
6.1	Probability-of-guilt model parameter estimates.	160
6.2	Dyslexic readers data: Parameter estimates for two models.	163
6.3	Variance–covariance matrix for probability-of-guilt model.	166
6.4	Bias-corrected estimates for the probability-of-guilt model.	169
6.5	Bias-corrected and bias-reduced estimates for the dyslexic readers model.	170
6.6	Dyslexic readers example: 1-inflated beta GLM.	178
6.7	Parameter estimates and confidence intervals for car dealership example.	181
6.8	Parameter estimates for the car dealership full model.	186
7.1	Reading accuracy scores Tobit model.	201
7.2	Reading accuracy scores heteroscedastic Tobit model.	203
7.3	Response times censored regression models.	207
7.4	Euthanasia attitude data doubly censored regression models.	216
7.5	Lab scores example <code>quantreg</code> output.	227
8.1	Test accuracy example, parameter estimates, and standard errors of the multilevel logistic model with crossed random effects.	238
8.2	Condom usage model fixed-effects parameter estimates.	240
8.3	Simple effects models: By condition.	241
8.4	Judged probability model parameter estimates.	244
8.5	Skipping school model parameter estimates.	247
8.6	Probability-of-guilt model MCMC estimates.	250
8.7	Duration sums MCMC estimates.	256
8.8	Dominance analysis applied to reading accuracy scores, Part 1.	259
8.9	Dominance analysis applied to reading accuracy scores, Part 2.	260

Notation

Symbol	Definition
\sim	is distributed as (e.g., $X \sim N(0, 1)$)
\approx	is approximately equal to (e.g., $\pi \approx 3.14$)
α_j	the intercept/cutpoint associated with the j^{th} logit of an ordinal regression model
$\text{Be}(\omega, v)$	a Beta distribution with parameters ω and v
$\text{B}(\cdot)$	Beta function
$\boldsymbol{\beta}$	a vector of regression coefficients
$\boldsymbol{\beta}_j$	a vector of regression coefficients associated with the j^{th} logit of an ordinal regression
β_k	the regression coefficient associated with x_k
χ_{df}^2	a chi-squared-distributed variate with df degrees of freedom
$D(M)$	The deviance associated with Model M
$\exp(x)$ or e^x	The exponential of x
$f(\cdot)$	usually a probability density function (pdf)
$F(\cdot)$	usually a cumulative density function (cdf)
$g(\cdot)$	usually a link function in the location submodel of a GLM
G^2	Likelihood ratio statistic
$\Gamma(\cdot)$	Gamma function
$h(\cdot)$	usually a link function in the dispersion submodel of a GLM
I	number of unique combinations of values that may be assumed by a vector of predictors \mathbf{x}
J	number of categories in a categorical or ordinal variable
$L(\boldsymbol{\theta} \mathbf{y}, \mathbf{X})$	A model likelihood function
$\ell(\boldsymbol{\theta} y_i, \mathbf{x}_i)$	Individual i 's contribution to the likelihood
μ	population mean
$N(\mu, \sigma)$	a normal distribution with mean μ and standard deviation σ
n	number of observations
φ	population precision in Beta and Negative Binomial Models
$\phi(\cdot), \Phi(\cdot)$	the pdf and cdf for the standard normal distribution $N(0, 1)$
π	mixture probability in zero-inflated models
$P(\cdot)$	Probability of an event or a proposition
p	Usually a probability in the role of a dependent variable in a GLM
q	the number of estimated model parameters
r	a model residual

σ	population standard deviation
Σ	a population variance–covariance matrix
t	exposure variable
τ	censoring or truncation threshold in probit, logit, and Tobit models
θ	a model parameter vector
ω, ν	parameters of a beta distribution
X^2	Pearson goodness-of-fit statistic
\mathbf{X}	a matrix of predictor variables
\mathbf{x}_i	a vector of predictor variables for observation i
x_k	the k^{th} predictor variable
y	a response variable
y^*	a latent, continuous variable underlying y
W	Wald statistic

About the Authors

Michael Smithson is a Professor in the Research School of Psychology at The Australian National University in Canberra. He received his PhD from the University of Oregon. His primary research interests are in judgment and decision making under uncertainty, statistical methods for the social sciences, and applications of fuzzy set theory to the social sciences. He has authored or co-authored six books and co-edited two volumes. His other publications include more than 140 refereed journal articles and book chapters.

Edgar C. Merkle is an Assistant Professor in the Department of Psychological Sciences at the University of Missouri. He received a PhD in Quantitative Psychology and an MS in Statistics, both from the Ohio State University. His research interests include latent variable models, subjective probabilities and forecasts, and statistical computing. He has authored numerous journal articles within these areas.

This page intentionally left blank

Introduction and Overview

1.1 The Nature of Limited Dependent Variables

Many variables in the social sciences are “limited” in the sense that their supports have boundaries. In fact, we claim that the vast majority of these variables are bounded. These bounds consist of two related kinds: Categorical boundaries and bounds on a continuum. The key distinction between these is that cases contained within categorical bounds are treated as identical in value or state, whereas cases falling in a bounded continuum may take different values within that range.

The primary rationale for this book and related books, including Long (1997), Agresti (2012), Powers and Xie (2008), and Bishop, Fienberg, and Holland (1975), is that the most popular distributions for model error (the normal and t distributions) assume that the dependent variables to which they are applied are unbounded (i.e., their support encompasses the entire real line). In contrast, bounded data require distributions that take their bounds into account.

Categorical bounds imply discreteness, and so categorical random variables require discrete distributions. These occupy the first part of this book. Bounds on continuous variables raise important measurement issues, primarily regarding cases at or near the boundaries. Are such cases accurately recorded or do their true scores lie “beyond the bounds?” Often the bounds are artifacts insofar as they are imposed by such considerations as constraints on the number of items comprising a scale or a test, a practical need to identify scale endpoints, or a decision to record only cases that exceed or lie below some threshold.

We distinguish three situations regarding bounds. First, the dependent variable data are completely known for all cases, so that cases on a boundary have been accurately measured and, in this sense, the boundary is “real.” Second, some of the cases are *censored* because they are only partially known (e.g., we know that a particular debtor owes less than \$1,000 but we do not know exactly how much she owes). Third, some cases are *truncated* because they have been excluded from the sample on the basis of some characteristic (e.g., a bank records only losses exceeding \$1,000). We will examine censoring and truncation in Chapters 5 and 7.

Chapter 6 deals with doubly bounded continuous dependent variables, i.e., those that have both a lower and an upper bound, where the boundary cases are real. The most obvious examples are proportions and percentages, but other examples are readily found such as borrowings on credit cards with upper limits, rating scales used by

judges in performance sports such as diving or gymnastics, prototypicality or degree of membership rating scales, and happiness or satisfaction ratings.

Bounds may be “unreal” in two ways. First, they may be imposed arbitrarily, such as the endpoints in the popular agree–disagree rating scales. If the endpoints are “strongly disagree” and “strongly agree” then all we know about cases occupying those endpoints is that they are that strong or stronger in their degree of (dis)agreement. Another influence on the veridicality of scale endpoints is the granularity of the scale. A confidence rating scale will have more accurate endpoints if it uses 20 bins than if it uses 4.

Second, the bounds may be real (e.g., a score of 0% or 100% on an exam) but the boundary cases may or may not be censored depending on the measurement purpose. A score of 0% or 100% on an exam is a true score if the exam includes all relevant test items, but not if it is considered to include only a sample of relevant items. Likewise, if a survey asking respondents to report the amount of time they spent eating during the last 24 hours turns up cases reporting 0, those cases may be taken as true scores if the purpose is simply to measure how much time people from the relevant population spent on eating during the past 24 hours, but not if it is to measure the proportion of 24 hours that people typically devote to eating.

Finally, limited dependent variables may be “boundary inflated” in the sense that the proportion of boundary cases exceeds that presupposed by the distribution model. A count variable such as number of cigarettes smoked in the past month might produce a large number of zeros in a population where most people do not smoke cigarettes. Chapters 5 and 6 include treatments of boundary-inflated models.

1.2 Overview of GLMs

Generalized linear models (GLMs), as originally described by Nelder and Wedderburn (1972) and expanded upon by McCullagh and Nelder (1989), form the foundation of this book. In this section, we define GLMs, relate them to simpler models, and discuss extensions.

1.2.1 Definition

GLMs are extensions of the standard linear regression model to situations where the dependent variable (or response variable) is limited. A linear predictor is common to all models, but the way in which the linear predictor relates to the data is different. To begin, consider the standard linear regression model, which we will generally call the *Gaussian GLM* for reasons that will become clear:

$$y = \mathbf{x}\boldsymbol{\beta} + e = \beta_0 + \beta_1 x_1 + \dots + \beta_K x_K + e, \quad (1.1)$$

where we assume a 1 in the first entry of \mathbf{x} and that $e \sim N(0, \sigma^2)$. With $K + 2$ or more observations, we can fit this model to the data and obtain estimates of $\boldsymbol{\beta}$ and σ^2 . Given estimates of $\boldsymbol{\beta}$ and some value of \mathbf{x} , we can then predict the mean of y . The regression weights $\boldsymbol{\beta}$ provide insight about the extent to which the predicted mean of y is impacted by each of the x_j . Thus, we typically wish to interpret the β s and test

whether they differ from zero. These tests and interpretations help us to summarize the general impact of the x_j on y .

To “generalize” the model in Equation (1.1), we must think about it in an alternative manner. We consider three features of the model: (i) the distribution associated with y , (ii) the parameter of the distribution on which we wish to focus, and (iii) the way in which we model the parameter via a linear predictor. For the Gaussian GLM, we have (i) a normal distribution for y , (ii) focus on μ , the normal distribution’s mean parameter, and (iii) a linear predictor placed directly on μ . The third feature is admittedly confusing for this model, because it is not clear why it is needed. However, it will help us soon.

In the spirit of the three features that we just noted, we could rewrite the model from (1.1) as

$$y|\mathbf{x}, \boldsymbol{\beta} \sim N(\mu, \sigma^2) \quad (1.2)$$

$$\mu = \mathbf{x}\boldsymbol{\beta}. \quad (1.3)$$

These equations no longer contain an error term that is added to a linear predictor. Instead, we assume a distribution for y and then place a linear predictor on the conditional mean of the distribution of the response variable. The distribution on y , the normal (also known as the Gaussian), leads us to call this model a Gaussian GLM. The parameter that we model, μ , is unbounded and can therefore be directly modeled via the linear predictor. In general, however, model parameters are not unbounded, so that a linear predictor may make nonsensical predictions. For example, we will see many situations where a model parameter can only assume values between 0 and 1. In this case, placing a linear predictor directly on the parameter may result in predictions that are less than zero or greater than 1. We need a function that “unbounds” the model parameter, making it sensible to use a linear predictor. This function is termed a *link function*, because it “links” a model parameter to a linear predictor. Given a distribution for y , there exists a special link function that implies some good statistical properties (involving the fact that the sufficient statistics for the model are of a simple form). This special link function is termed a *canonical link function*. However, choice of link function is more often guided by precedent than by whether or not it is canonical.

We will describe a variety of link functions throughout this book, though we illustrate only one here. In Chapter 2 and elsewhere, we model a probability parameter p that can only assume values in $(0, 1)$. To use a linear predictor in this model, the link function must transform values in $(0, 1)$ to values in $(-\infty, \infty)$. A common choice (which also happens to be the canonical choice) is the logit link, also known as the log-odds link. This is given as

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right). \quad (1.4)$$

This function is described in much more detail in Chapter 2. In a GLM context, we would place a linear predictor on $\text{logit}(p)$ to obtain

$$\text{logit}(p) = \mathbf{x}\boldsymbol{\beta}, \quad (1.5)$$

so that the regression weights lead to predictions on the transformed, logit scale. We typically wish to have predictions on the probability scale, and we can obtain these via an inverse link function. The inverse link transforms unbounded predictions back to the bounded scale, with the inverse logit given as

$$\text{logit}^{-1}(\mathbf{x}\boldsymbol{\beta}) = \frac{1}{1 + \exp(-\mathbf{x}\boldsymbol{\beta})}. \quad (1.6)$$

The main point that the reader should take from this discussion is the fact that, for many distributions, we must unbound a parameter before using a linear predictor. This is accomplished via a link function. Throughout this book, we typically describe GLMs via the three features above: the assumed distribution for y , the focal model parameter (which is usually the conditional mean of the distribution), and the link function that associates the focal parameter with a linear predictor. The near-exclusive focus on linear predictors largely follows the tradition of model developments; however, at least for unbounded data, a linear predictor generally approximates a function of interest via a first-order Taylor-series expansion (see, e.g., Venables, 2000). The accuracy of this approximation obviously varies.

1.2.2 Extensions

As defined by McCullagh and Nelder (1989), GLMs include only models whose distributions arise from the exponential family (see, e.g., Casella & Berger, 2002 for a formal definition of the exponential family). While this family includes many common distributions (including the normal, binomial, Poisson, beta, and gamma), many of the models that we describe in this book do not fall in the family. This makes the title of the book inaccurate in some respects, although the models not included in the exponential family still use many of the same concepts. We describe here some extensions of “plain” GLMs, some of which still fall in the GLM family and some of which do not. While we recognize the inaccuracy in terminology, we generally ignore the “GLM family or not” distinction because it does not have a major impact on the applied researcher.

We will often distinguish between *location* and *dispersion* parameters throughout the book. Informally, location parameters are those that influence the central tendency of a distribution. Dispersion parameters, on the other hand, influence only the variability of a distribution. When a parameter influences both central tendency and variability, it is generally regarded as a location parameter (the mean of the Poisson distribution is one example).

Applied researchers are quite familiar with modeling location: most popular statistical tests employ null hypotheses associated with mean parameters. Further, many of the models in this book will allow us to assess the impact of predictor variables on the mean of a distribution. Dispersion parameters, on the other hand, are often regarded as error parameters or nuisance parameters, and less attention is typically devoted to them. As we will see in this book, however, these parameters are often very important both for estimation and interpretation. Dispersion parameters can affect the interpretations that we make about location parameters, and they can also

be interesting to study in isolation. Expanding on the latter issue, there are some situations in which we would expect the predictor variables to directly influence the variability of a response variable. For bounded variables, this variability can additionally be interpreted as polarization. The associated statistical tests of polarization can inform many theories in the social sciences. Additionally, the mean and variance of a bounded variable are not independent of one another. Thus, dispersion parameters generally play an important role in models for bounded dependent variables.

For both location and dispersion parameters, there exist a few common scenarios where the attributes of y are not a good match with the associated model. In these scenarios, there are often relatively simple modifications that can be employed to salvage the original model. We describe the location scenarios separately from dispersion scenarios below.

Location. Focusing on location parameters, the scenarios primarily include situations where y is *censored* or *truncated*. Truncation occurs when we exclude some values of y above or below a certain point. For example, if we ask current smokers about the number of cigarettes that they smoke each month, we have excluded responses of “zero.” This situation may be called “truncation from below at 1” or, alternatively, “left truncation at 1.” In contrast, we would observe truncation from “above” (or from the “right”) if we excluded values above a certain point of the response variable y .

Censoring is similar to truncation, except that all values of y above or below a point τ have simply been relabeled as τ . For example, imagine a bathroom scale with a maximum weight of 400 lbs. While this will not be an issue for many people, anyone who weighs more than 400 lbs will be assigned a weight of 400. This is a case of “right-censoring at 400.” Censored observations are observed (e.g., a 450-lb individual is recorded as weighing 400 lbs), while truncated observations are unobserved (e.g., people who smoke zero cigarettes are not included in the sample).

Censoring and truncation are often observed in the context of count regression (Chapter 5), though we also encounter these issues in other situations. Chapter 7 describes the general handling of censoring and truncation, including situations where the associated threshold is unknown (for example, where we do not know the maximum weight that a scale can record).

Dispersion. Focusing on dispersion parameters, major scenarios are *overdispersion* and *heteroscedasticity*. When we have overdispersion, the model-predicted variance systematically underestimates the variance observed in y . That is, the model simply cannot account for the variance observed in y while simultaneously accounting for the mean (i.e., for the location) of y . The overdispersion issue is important because, when it occurs, the standard errors associated with the location estimates are too small. This can often lead the researcher to erroneously infer statistically significant effects of a predictor variable on y (i.e., Type I errors). To correct for overdispersion, it is possible to add an extra parameter that accounts for the variance that the original model could not. It is also sometimes possible to use a different distribution for y that better handles dispersion. The overdispersion issue is most relevant to the count