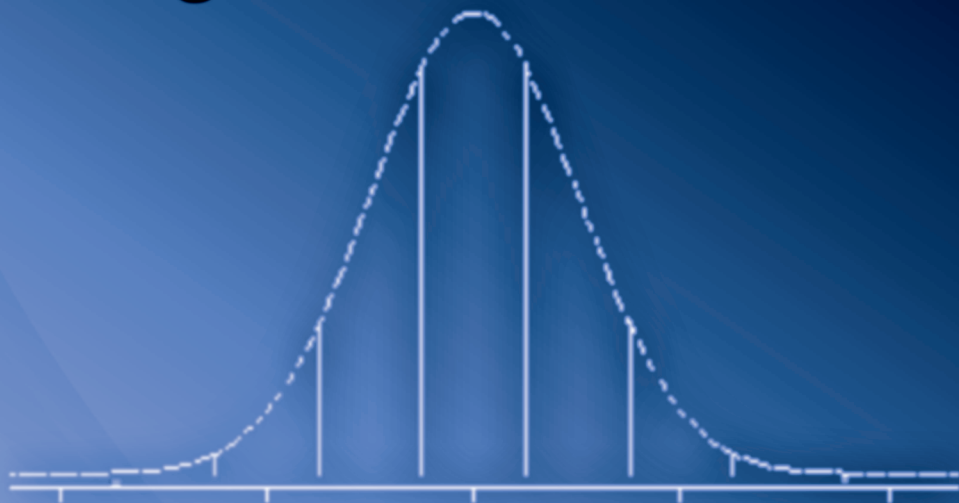


Multivariate Generalized Linear Mixed Models Using R



Damon M. Berridge
Robert Crouchley

 **CRC Press**
Taylor & Francis Group

A CHAPMAN & HALL BOOK

Multivariate Generalized Linear Mixed Models Using R

Damon M. Berridge
Robert Crouchley



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

© 2011 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: 20111012

International Standard Book Number-13: 978-1-4398-1327-0 (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>

Contents

List of Figures	x <i>i</i>
List of Tables	x <i>iii</i>
List of Applications	xv
List of Datasets	xv <i>ii</i>
Preface	xix
Acknowledgments	xx <i>iii</i>
1 Introduction	1
2 Generalized linear models for continuous/interval scale data	9
2.1 Introduction	9
2.2 Continuous/interval scale data	10
2.3 Simple and multiple linear regression models	11
2.4 Checking assumptions in linear regression models	12
2.5 Likelihood: multiple linear regression	13
2.6 Comparing model likelihoods	14
2.7 Application of a multiple linear regression model	15
2.8 Exercises on linear models	17
3 Generalized linear models for other types of data	21
3.1 Binary data	21
3.1.1 Introduction	21
3.1.2 Logistic regression	22
3.1.3 Logit and probit transformations	23
3.1.4 General logistic regression	24
3.1.5 Likelihood	24
3.1.6 Example with binary data	24
3.2 Ordinal data	26
3.2.1 Introduction	26

3.2.2	The ordered logit model	27
3.2.3	Dichotomization of ordered categories	29
3.2.4	Likelihood	29
3.2.5	Example with ordered data	30
3.3	Count data	32
3.3.1	Introduction	32
3.3.2	Poisson regression models	33
3.3.3	Likelihood	34
3.3.4	Example with count data	34
3.4	Exercises	37
4	Family of generalized linear models	43
4.1	Introduction	43
4.2	The linear model	44
4.3	The binary response model	44
4.4	The Poisson model	46
4.5	Likelihood	46
5	Mixed models for continuous/interval scale data	49
5.1	Introduction	49
5.2	Linear mixed model	49
5.3	The intraclass correlation coefficient	51
5.4	Parameter estimation by maximum likelihood	53
5.5	Regression with level-two effects	54
5.6	Two-level random intercept models	55
5.7	General two-level models including random intercepts	56
5.8	Likelihood	58
5.9	Residuals	58
5.10	Checking assumptions in mixed models	59
5.11	Comparing model likelihoods	60
5.12	Application of a two-level linear model	61
5.13	Two-level growth models	66
5.13.1	A two-level repeated measures model	66
5.13.2	A linear growth model	66
5.13.3	A quadratic growth model	67
5.14	Likelihood	67
5.15	Example using linear growth models	68
5.16	Exercises using mixed models for continuous/interval scale data	69
6	Mixed models for binary data	75
6.1	Introduction	75

6.2	The two-level logistic model	75
6.3	General two-level logistic models	77
6.4	Intraclass correlation coefficient	77
6.5	Likelihood	78
6.6	Example using binary data	78
6.7	Exercises using mixed models for binary data	81
7	Mixed models for ordinal data	85
7.1	Introduction	85
7.2	The two-level ordered logit model	85
7.3	Likelihood	86
7.4	Example using mixed models for ordered data	87
7.5	Exercises using mixed models for ordinal data	90
8	Mixed models for count data	93
8.1	Introduction	93
8.2	The two-level Poisson model	93
8.3	Likelihood	94
8.4	Example using mixed models for count data	95
8.5	Exercises using mixed models for count data	97
9	Family of two-level generalized linear models	99
9.1	Introduction	99
9.2	The mixed linear model	100
9.3	The mixed binary response model	100
9.4	The mixed Poisson model	102
9.5	Likelihood	102
10	Three-level generalized linear models	105
10.1	Introduction	105
10.2	Three-level random intercept models	105
10.3	Three-level generalized linear models	106
10.4	Linear models	107
10.5	Binary response models	108
10.6	Likelihood	108
10.7	Example using three-level generalized linear models	109
10.8	Exercises using three-level generalized linear mixed models	112
11	Models for multivariate data	115
11.1	Introduction	115
11.2	Multivariate two-level generalized linear model	116

11.3	Bivariate Poisson model: example	117
11.4	Bivariate ordered response model: example	121
11.5	Bivariate linear-probit model: example	126
11.6	Multivariate two-level generalized linear model likelihood	131
11.7	Exercises using multivariate generalized linear mixed models	131
12	Models for duration and event history data	135
12.1	Introduction	135
12.1.1	Left censoring	135
12.1.2	Right censoring	135
12.1.3	Time-varying explanatory variables	136
12.1.4	Competing risks	136
12.2	Duration data in discrete time	137
12.2.1	Single-level models for duration data	137
12.2.2	Two-level models for duration data	139
12.2.3	Three-level models for duration data	140
12.3	Renewal data	143
12.3.1	Introduction	143
12.3.2	Example: renewal models	145
12.4	Competing risk data	147
12.4.1	Introduction	147
12.4.2	Likelihood	148
12.4.3	Example: competing risk data	150
12.5	Exercises using renewal and competing risks models	153
13	Stayers, non-susceptibles and endpoints	157
13.1	Introduction	157
13.2	Mover-stayer model	157
13.3	Likelihood incorporating the mover-stayer model	160
13.4	Example 1: stayers within count data	161
13.5	Example 2: stayers within binary data	164
13.6	Exercises: stayers	166
14	Handling initial conditions/state dependence in binary data	169
14.1	Introduction to key issues: heterogeneity, state dependence and non-stationarity	169
14.2	Example	170
14.3	Random effects models	171
14.4	Initial conditions problem	172

14.5	Initial treatment	173
14.6	Example: depression data	174
14.7	Classical conditional analysis	174
14.8	Classical conditional model: example	175
14.9	Conditioning on initial response but allowing random effect u_{0j} to be dependent on z_j	176
14.10	Wooldridge conditional model: example	177
14.11	Modelling the initial conditions	178
14.12	Same random effect in the initial response and subsequent response models with a common scale parameter	179
14.13	Joint analysis with a common random effect: example	180
14.14	Same random effect in models of the initial response and subsequent responses but with different scale parameters	181
14.15	Joint analysis with a common random effect (different scale parameters): example	182
14.16	Different random effects in models of the initial response and subsequent responses	183
14.17	Different random effects: example	184
14.18	Embedding the Wooldridge approach in joint models for the initial response and subsequent responses	185
14.19	Joint model incorporating the Wooldridge approach: example	187
14.20	Other link functions	187
14.21	Exercises using models incorporating initial conditions/state dependence in binary data	188
15	Incidental parameters: an empirical comparison of fixed effects and random effects models	195
15.1	Introduction	195
15.2	Fixed effects treatment of the two-level linear model	197
15.3	Dummy variable specification of the fixed effects model	199
15.4	Empirical comparison of two-level fixed effects and random effects estimators	200
15.5	Implicit fixed effects estimator	204
15.6	Random effects models	204
15.7	Comparing two-level fixed effects and random effects models	208
15.8	Fixed effects treatment of the three-level linear model	208

15.9	Exercises comparing fixed effects and random effects . . .	209
A	SabreR installation, SabreR commands, quadrature, estimation, endogenous effects	215
A.1	SabreR installation	215
A.2	SabreR commands	215
	A.2.1 The arguments of the SabreR object	215
	A.2.2 The anatomy of a SabreR command file	216
A.3	Quadrature	218
	A.3.1 Standard Gaussian quadrature	218
	A.3.2 Performance of Gaussian quadrature	219
	A.3.3 Adaptive quadrature	221
A.4	Estimation	223
	A.4.1 Maximizing the log likelihood of random effects models	223
A.5	Fixed effects linear models	225
A.6	Endogenous and exogenous variables	226
B	Introduction to R for Sabre	229
B.1	Getting started with R	229
	B.1.1 Preliminaries	229
	B.1.1.1 Working with R in interactive mode	229
	B.1.1.2 Basic functions	231
	B.1.1.3 Getting help	232
	B.1.1.4 Stopping R	232
	B.1.2 Creating and manipulating data	232
	B.1.2.1 Vectors and lists	232
	B.1.2.2 Vectors	233
	B.1.2.3 Vector operations	234
	B.1.2.4 Lists	235
	B.1.2.5 Data frames	236
	B.1.3 Session management	237
	B.1.3.1 Managing objects	237
	B.1.3.2 Attaching and detaching objects	237
	B.1.3.3 Serialization	238
	B.1.3.4 R scripts	238
	B.1.3.5 Batch processing	239
	B.1.4 R packages	239
	B.1.4.1 Loading a package into R	239
	B.1.4.2 Installing a package for use in R	239
	B.1.4.3 R and Statistics	240
B.2	Data preparation for SabreR	240

B.2.1	Creation of dummy variables	240
B.2.2	Missing values	243
B.2.3	Creating lagged response covariate data	245
References		249
Author Index		259
Subject Index		263

This page intentionally left blank

List of Figures

11.1	The relationship between wages and trade union membership: I	127
11.2	The relationship between wages and trade union membership: II	127
11.3	The relationship between wages and trade union membership: III	128
12.1	Duration data	136
12.2	Diagrammatic representation of renewal data	143
12.3	Example of competing risk data: failure due to two failure mechanisms	148
12.4	Data required to model failure due to mechanism A	148
12.5	Data required to model failure due to mechanism B	149
13.1	The normal distribution	158
13.2	Quadrature points approximating the normal distribution	158
13.3	Quadrature with left and right endpoints	159
13.4	Quadrature with left endpoint only	159
B.1	First few lines of <code>essays.tab</code>	241
B.2	First few lines of new dataset <code>essays2.tab</code>	242
B.3	First few lines of <code>thaieduc.tab</code>	243
B.4	Ungrouped depression data (<code>depression0.tab</code>)	246
B.5	First few lines of <code>depression.tab</code>	247
B.6	First few lines of new dataset <code>depression2.tab</code>	248

This page intentionally left blank

List of Tables

11.1	Crosstabulation of <code>dvisits</code> by <code>prescrib</code>	119
12.1	Sample of duration data in continuous time	138
12.2	Sample of duration data, reconfigured in discrete time .	138
12.3	Sample of renewal data in continuous time	144
12.4	Sample of renewal data, reconfigured in discrete time . .	144
12.5	Sample of competing risk data in continuous time	149
12.6	Sample of competing risk data, reconfigured in discrete time	149
13.1	Observed migration frequencies	162
14.1	Depression data (1 = depressed, 0 = not depressed) . .	171

This page intentionally left blank

List of Applications

Angina pectoris (renewal data), 153
Attitudes to abortion, 6, 39, 91
Attitudes to gender roles (bivariate ordered data), 115, 121

Choosing teaching as a profession, 30, 87

Demand for health care, 34, 95, 216
Demand for health care (bivariate count data), 115–117
Depression, 170, 172, 174, 175, 177, 180, 182, 184, 187, 245

Educational attainment, 3, 18, 70
Effect of education on log wages, 210
Effect of job training on firm scrap rates, 209
Epileptic seizures, 7, 40, 97
Essay grading, 240
Essay grading (binary response), 4, 37, 81
Essay grading (continuous response), 2, 17, 70
Essay grading (ordered response), 5, 38, 90
Expiratory flow rates (bivariate data), 131

Female employment participation (stayers in binary data), 167
Female UK labour force participation, 191
Filled and lapsed vacancies (competing risk data), 150
Filling job vacancies (competing risk data), 116
Filling vacancies (three-level data), 140
Fish caught by US National Park visitors (stayers in count data), 168

German unemployment (competing risk data), 155

Headaches, 7
Headaches (count data), 40, 97

Immunization of Guatemalan children, 5, 38, 83
Immunization of Guatemalan children (binary response), 113

Log wages (three-level data), 109

- Mathematics achievement, 15, 61, 105
- Migration moves (binary data), 164
- Migration moves (count data), 161

- Patents and R&D expenditure, 192
- Psychological distress, 2, 10, 11, 14, 17, 49, 52, 55, 56, 60, 69
- Pupil rating of school managers, 4, 19, 72
- Pupil rating of school managers (three-level data), 208

- Repeating a grade, 24, 78, 243
- Residential mobility, 145
- Respiratory status, 6, 39, 92

- Skin cancer deaths, 7, 41, 98, 113
- Student evaluation of teachers, 68

- Tower of London, 5, 37, 82
- Tower of London (binary response), 112
- Trade union membership, 4, 37, 81
- Trade union membership (stayers in binary data), 166
- Trade union membership of females, 189
- Trade union membership of young males, 188

- Unemployment claims, 3, 18, 71

- Wage determinants, 3, 18, 72
- Wages and trade union membership, 116, 121
- Wages and trade union membership (bivariate data), 126, 132
- Wages of young women, 200

List of Datasets

abortion2.tab, 6, 39, 91
angina.tab, 154
deaths.tab, 7, 41, 98, 114
depression.tab, 174, 179, 180, 182, 184, 187, 247
depression0.tab, 245, 246
depression1.tab, 248
depression2.tab, 174, 175, 177, 247, 248
epilep.tab, 7, 40, 97
essays.tab, 240, 242
essays2.tab, 4, 37, 81, 242
essays_ordered.tab, 6, 38, 90
ezunem2.tab, 3, 18, 71
fish.tab, 168
ghq2.tab, 2, 10, 70
grader1.tab, 2, 70
grader2.tab, 2, 17
guatemala_immun.tab, 5, 38, 83, 113
headache2.tab, 7, 40, 97
hsb.tab, 15, 61
jtrain.tab, 209
labour.tab, 167
manager.tab, 4, 19, 72, 110
neighbourhood.tab, 3, 18, 70
nls.tab, 126, 128, 166, 200
nls wage-union.tab, 126
opfama.tab, 121, 122
opfamaf.tab, 124
opfamf.tab, 122, 123
patents.tab, 192
pefr.tab, 132
racd.tab, 35, 95, 216, 217
respiratory2.tab, 7, 39, 92
roch.tab, 145
rochmig.tab, 164
rochmigx.tab, 161

teacher1.tab, 30
teacher2.tab, 30, 87
thaieduc.tab, 243, 244
thaieduc1.tab, 24, 78, 244, 245
thaieduc2.tab, 78, 245
tower1.tab, 5, 37, 82, 112
unemployedR.tab, 155
unionjmw1.tab, 188, 189
unionjmw2.tab, 188, 189
unionred1.tab, 189, 190
unionred2.tab, 189, 190
vacancies.tab, 151
visit-prescribe.tab, 118
vwks4_30k.tab, 140
wagepan.tab, 3, 5, 19, 37, 72, 81, 132
wagepan2.tab, 210
wemp-base1.tab, 191
wemp-base2.tab, 191

Preface

The main aims of this book are to provide an introduction to the principles of modelling as applied to longitudinal data from panel and related studies with the necessary statistical theory, and to describe the application of these principles to the analysis of a wide range of examples using the Sabre software (<http://sabre.lancs.ac.uk/>) from within R.

This material on multivariate generalized linear mixed models arises from the activities at the Economic and Social Research Council (ESRC)-funded Colaboratory for Quantitative e-Social Science (CQeSS) at Lancaster University from 2003 to 2008. Sabre is a program for the statistical analysis of multi-process event/response sequences. These responses can take the form of binary, ordinal, count and linear recurrent events. The response sequences can also be of different types, for example, a linear response (wages) and a binary one (trade union membership). Such multi-process data are common in many research areas, for example, in the analysis of work and life histories from the British Household Panel Survey or the German Socio-Economic Panel Study where researchers often want to disentangle state dependence (the effect of previous responses or related outcomes) from any omitted effects that might be present in recurrent behaviour (for example, unemployment). Understanding the need to disentangle these generic substantive issues dates back to the study of accident proneness in the 1950s and has since been discussed in many applied areas, including consumer behaviour and voting behaviour. These issues, and others relating to the analysis of longitudinal or event history data, are discussed in more detail in the following text:

- Shahtahmasebi, S. and Berridge, D. (2010) *Conceptualizing Human Behaviour in Health and Social Research: A Practical Guide to Data Analysis*, New York: Nova

Some key contributions in the References, including a number of Heckman's seminal works, have been reprinted in the following series:

1. Penn, R. and Berridge, D. (2010) *Social Statistics Volume 1: The Fundamentals of Descriptive Social Statistics*, London: Sage
2. Penn, R. and Berridge, D. (2010) *Social Statistics Volume 2: The Development of Statistical Modelling*, London: Sage

3. Penn, R. and Berridge, D. (2010) *Social Statistics Volume 3: Statistical Modelling of Longitudinal Data*, London: Sage
4. Penn, R. and Berridge, D. (2010) *Social Statistics Volume 4: Statistical Modelling of Ordinal Categorical Data*, London: Sage

Those contributions appearing in this series are indicated by asterisks in the References. One asterisk indicates Volume 1, two asterisks indicate Volume 2, and so on.

Sabre can also be used to model collections of single sequences such as may occur in medical trials on the number of headaches experienced over a sequence of weeks, or in single-equation descriptions of cross-sectional clustered data such as the educational attainment of children in schools.

Sabre is available in three forms: (1) stand-alone (as discussed in Shahtahmasebi and Berridge, 2010), (2) the R plugin (as discussed in the current text), and (3) the Stata plugin (as discussed on the Sabre web page — see above).

The class of models that can be estimated by Sabre may be termed Multivariate Generalized Linear Mixed Models (MGLMMs). These models have special features to help them disentangle state dependence from the incidental parameters (omitted or unobserved effects). The incidental parameters can be treated as random or fixed. The random effects models can be estimated with standard Gaussian quadrature or adaptive Gaussian quadrature. Quadrature methods (and particularly adaptive Gaussian quadrature) are the most reliable way of handling random effects in MGLMMs, as the adequacy of the numerical integration can be improved by adding more quadrature points. The number of quadrature points required will depend on the model being estimated. If additional quadrature points fail to improve the log likelihood, then we have found an accurate evaluation of the integral. Even though the linear model integral has a closed form solution, we do not use it as it cannot easily be used in multivariate models when some of the joint sequences do not have interval level responses. Also current computational facilities on many desktop computers often make the delay involved in using numerical integration for the linear model negligible for many small to medium-sized data sets. ‘End effects’ can also be added to the models to accommodate ‘stayers’ or ‘non-susceptibles’. The fixed effects algorithm we have developed uses code for large sparse matrices from the Harwell Subroutine Library; see <http://www.cse.scitech.ac.uk/nag/hs1/>.

Also included in Sabre is the option to undertake all the calculations using increased accuracy. Numerical underflow and overflow often occur in the estimation process for models with incidental parameters. We suppose that many of the alternative software systems truncate their

calculations without informing the user when this happens as there is little discussion of this in their respective user manuals.

This book is written in a way that we have found appropriate for some of our short courses. The book starts by discussing members of the family of generalized linear models and gradually adds complexity to the modelling framework by incorporating random effects. We then review the generalized linear model notation before illustrating a range of more substantively appropriate random effects models, for example, the three-level model, multivariate (in particular, bivariate and trivariate) models, endpoint, event history and state dependence models. The MGLMMs are estimated using either standard Gaussian quadrature or adaptive Gaussian quadrature. The book compares two-level fixed and random effects linear models. Additional information on quadrature, model estimation and endogenous variables is included in Appendix A. Appendix B contains an introduction to R and some examples of using R to pre-process the data for Sabre.

There are two other related SabreR booklets available from the Sabre web page:

- Exercises for SabreR
- Solutions Manual for SabreR Exercises

These booklets contain the exercises and solutions on small data sets that have been written to accompany this book. These exercises will run quickly on a desktop PC.

Drafts of the chapters of this book were developed and revised in the process of preparing and delivering short courses in ‘Statistical Modelling using Sabre’, ‘Multilevel Modelling’ and ‘Event History Analysis’ given at CQeSS and the Department of Mathematics and Statistics at Lancaster University and elsewhere. We are grateful to many of the students of these courses who are from a range of backgrounds (for example, computational science and the social sciences) and whose comments and criticisms improved these early drafts. We think that the book should serve as a training manual for postgraduate Masters and research students, and as a self-teaching manual for data analysts.

If you have any suggestions as to how this book could be improved—for instance by the addition of other material—please let us know via the Sabre mailing list, sabre@lancaster.ac.uk.

We accept no liability for anything that might happen as a consequence of your use of Sabre, though we are happy to accept recognition of its successful use.

Dr. Damon M. Berridge and Professor Robert Crouchley
Lancaster University
February 2011

This page intentionally left blank

Acknowledgments

Many thanks to Dr. Iraj Kazemi for helping to draft the material in the first ten chapters of this book. Thanks to Professor Richard B. Davies for inspiring the early development of Sabre (Poisson and logit models with endpoints).

Many thanks to Daniel Grose for writing the R side of the SabreR library. Dan also wrote much of the introductory material on R in Appendix B. David Stott and John Pritchard undertook all the recent development work on Sabre. Dave wrote the standard Gaussian and adaptive Gaussian quadrature algorithms. John wrote the algorithm for manipulating the large sparse matrices used by the fixed effects estimator.

This work was supported by the following ESRC research grants:

- RES-149-28-1003: The Colaboratory for Quantitative e-Social Science (E-Social Science Centre Lancaster Node), principal investigator: Professor Robert Crouchley
- RES-149-25-0010: An OGSA Component-Based Approach to Middleware for Statistical Modelling, principal investigator: Professor Robert Crouchley
- RES-576-25-0019: The Lancaster-Warwick-Stirling Node: Developing Statistical Modelling in the Social Sciences (National Centre for Research Methods (NCRM) Phase 2), principal investigator: Professor Brian Francis

The NCRM Phase 2 grant was particularly important for the development of the bivariate ordered response model reported in Chapter 11.

Finally, we wish to express our gratitude to Professor Roger Penn for his assistance in proofreading the final draft of this book.

This page intentionally left blank

1

Introduction

A major objective of this book is to provide data analysts with the tools to analyze large and complex datasets using methodologically sound models, thereby enabling them to answer increasingly complex research questions. The statistical software used in this book is SabreR. This is a version of the package Sabre, for the statistical analysis of multi-process event/response sequences, which has been implemented within the R environment.

These responses can take the form of binary, ordinal, count and linear recurrent events. The response sequences can also be of different types, for example, a linear response (wages) and a binary response (trade union membership). Such multi-process data are common in many research areas, for example, in the analysis of work and life histories from the British Household Panel Survey or the German Socio-Economic Panel Study where researchers often want to disentangle state dependence (the effect of previous responses or related outcomes) from any omitted effects that might be present in recurrent behaviour (unemployment).

Understanding of the need to disentangle these generic substantive issues dates back to the study of accident proneness [14] and has been discussed in many applied areas, including consumer behaviour [75] and voting behaviour [34].

SabreR can also be used to model collections of single sequences such as those that may occur in medical trials, for example, headaches and epileptic seizures [29,30], or in single-equation descriptions of cross-sectional clustered data such as the educational attainment of children in schools.

The class of models that can be estimated by SabreR may be called multivariate generalized linear mixed models. These models have special features added to standard models to help us disentangle state dependence from the incidental parameters (omitted or unobserved effects). The incidental parameters can be treated as random or fixed, the random effects models being estimated using standard Gaussian quadrature or adaptive Gaussian quadrature. ‘End effects’ can also be added to the models to accommodate ‘stayers’ or ‘non-susceptibles’, resulting in a more parsimonious model which provides a better fit to the data with

fewer parameters than a non-parametric specification of the random effects. The fixed effects algorithm we have developed uses code for large sparse matrices from the Harwell Subroutine Library [49].

SabreR also includes the option to undertake all of the calculations using increased accuracy. This is important because numerical underflow and overflow often occur in the estimation process for models with incidental parameters.

Chapters 2 and 3 cover the analysis of single-level data of various types: continuous, binary, ordinal and count data using univariate generalized linear models. The material covered in these chapters is summarized in Chapter 4. Chapters 5 to 8 extend these models to handle multi-level, specifically two-level, data of various types: continuous, binary, ordinal and count data, using univariate generalized linear mixed models. The models considered in Chapters 5 to 8 are summarized in Chapter 9, and are generalized to handle three-level data in Chapter 10.

A key feature of this book is the emphasis on the application of statistical models to real-life examples. At the heart of each chapter will be a fully worked example. In addition, readers will have the opportunity to apply these statistical models and to interpret the resulting output through a large number of exercises spanning a wide variety of areas of application. The exercises illustrating the use of models for continuous/interval scale data in Chapters 2 and 5 are based on the following examples:

Example 1.1. Psychological distress

Twelve students completed the twelve-item version of Goldberg's General Health Questionnaire (GHQ) [42]. The questionnaire was completed by each student on two different occasions, separated by three days. A psychological distress score was computed, on the basis of the twelve GHQ items, for each student on each of the two occasions [39]. These student-occasion-specific scores are saved in the file `ghq2.tab`.

Example 1.2. Essay grading (continuous response)

Johnson and Albert [66] analyzed data on the grading of essays by several experts. Essays were graded on a scale between 1 and 10, with a score of 10 corresponding to 'excellent'. In this example, we consider a subset of the data limited to the grades given to 198 essays by markers 1 and 4. This subset of data is stored in the data file `grader1.tab` which may be found on the Sabre web page. The grades given by markers 1 and 4 are stacked in a single column `grade` in the file `grader2.tab`. This file also includes an identifier which distinguishes between the two graders, in other words, the variable `dg4` which takes value 1 if the grader is number 4, and value 0 otherwise. Alternative treatments of the response are considered in *Examples 1.7* and *1.11*.

Example 1.3. Educational attainment

Garner and Raudenbush [41] and Raudenbush and Bryk [93] studied the role of school and/or neighbourhood effects on the educational attainment of young people, from one Scottish Local Education Authority, who left school between 1984 and 1986. The primary outcome of interest is a young person's combined end-of-school educational attainment as measured by his/her grades.

Explanatory variables are available at two levels: (i) the individual young person level and (ii) the school and/or neighbourhood level. Most explanatory variables present in the dataset are specific to each young person. These variables include: young person's gender; verbal reasoning quotient and reading ability as measured by tests in primary school at age 11–12; father's occupation and education. The one school/neighbourhood-specific explanatory variable is an index of social deprivation for the local community within which the young person lived. The data are stored in the file `neighbourhood.tab` on the Sabre web page.

Example 1.4. Unemployment claims

Indiana's enterprise zone programme provided tax credits for cities with high poverty and unemployment levels. In a bid to establish whether those cities targetted by the programme had significantly lower unemployment claims than those cities lying outside enterprise zones, Papke [85] analyzed annual data from 1980 to 1988. The dataset (`ezunem2.tab`) comprises the number of unemployment claims in 22 cities, and whether each city was located within an enterprise zone, in each of the nine years 1980 to 1988.

Example 1.5. Wage determinants

Vella and Verbeek [103] analyzed annual data on 545 males from the Youth Sample of the US National Longitudinal Survey for the period 1980 to 1987. The version of the data used in this book (`wagepan.tab`) was obtained from Wooldridge [106]. We wish to relate the outcome of primary interest, log hourly wage (in US dollars), to a time-invariant factor (ethnicity) and a variety of time-dependent explanatory variables. Those variables allowed to vary over time include respondent demographics (marital status, region of US lived in, rural/urban area lived in), education (years of schooling), labour market experience and trade union membership. These data are re-considered in *Example 1.8*, where trade union membership is regarded as the binary response of interest.

Having analyzed these data in Chapters 2 and 5, we will return to this dataset on further occasions in this book. In Chapter 11, in the context of bivariate models, we will estimate a joint model for wages and

trade union membership. We will allow trade union membership to be endogenous in the wage equation. In Chapter 14, we will use the data on trade union membership to illustrate Wooldridge's [107] treatment of the initial conditions problem in first-order Markov models. In Chapter 15, we compare and contrast the inferences made when we first assume fixed effects and then proceed under the assumption of random effects. We will use these data to relate log wages to time-varying explanatory variables such as number of years of labour market experience, marital status and trade union membership, and to time-invariant factors including race and education.

Example 1.6. Pupil rating of school managers

856 pupils in 94 schools were asked to rate the performance of their school managers/directors on the basis of six questions, each response recorded on a four-point scale [64]. The response to each item given by each pupil is presented in the dataset `manager.tab`. Pupil-specific explanatory variables are gender and school year. School-specific factors are gender of the school manager/director and type of school which is classified into the following three categories: 'general (AVO)', 'professional (MBO&T)' and 'day/evening'.

The exercises illustrating the use of models for binary data in Chapters 3 and 6 are based on the following examples:

Example 1.7. Essay grading (binary response)

In an extension to *Example 1.1*, we use data on the grades given to 198 essays by markers 1 to 5. Essays were graded on a scale from 1 to 10, with 10 classified as 'excellent'. For the purposes of the current example, the original essay grading variable is converted into a binary response variable, labelled as `pass` in the dataset `essays2.tab`. The variable `pass` takes the value 1 for grades 5 to 10, and value 0 for grades 1 to 4. The primary objective in this example is to test for significant differences in this binary response between markers, whilst adjusting for six explanatory variables which characterize the 198 essays.

Four of these factors are lexical in nature: average word length (`wordlength`), square root of the number of words (`sqrtwords`), average sentence length (`sentlength`) and proportion of words in the essay which are prepositions (`prepos`). A fifth explanatory variable is related to punctuation: number of commas, multiplied by 100 and divided by the total number of words in the essay (`commas`). The sixth factor is the percentage of words in the essay which are spelt incorrectly (`errors`).

Example 1.8. Trade union membership

In *Example 1.5*, we related data from the Youth Sample of the US