

The Analysis of Cross-Classified Categorical Data

Stephen E. Fienberg

The Analysis of Cross-Classified Categorical Data

Second Edition

 Springer

Stephen E. Fienberg
Department of Statistics
Carnegie-Mellon University
Pittsburgh, PA 15213
fienberg@stat.cmu.edu

ISBN 978-0-387-72824-7

Library of Congress Control Number: 2007928366

© 2007 Springer Science+Business Media, LLC. This Springer edition is a reprint of the 1980 edition published by MIT Press.

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer Science+Business Media, LLC, 233 Spring Street, New York, NY 10013, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Printed on acid-free paper.

Printed in the United States of America.

9 8 7 6 5 4 3 2 1

springer.com

To Fred

Preface to the Second Edition

The reception given to the first edition of this book, especially by nonstatisticians, has been most pleasing. Yet several readers have written to me asking for further details on, or clarifications of, methods and examples, and suggesting the preparation of sets of problems at the end of each chapter so that the book would be more useful as a text. This second edition was prepared, in large part, as a response to these requests.

Methodological research on the analysis of categorical data based on the use of loglinear models has continued at a rapid pace over the last three years. In this new edition, I have attempted to expand the discussion of several topics, by drawing selectively from this new literature, while at the same time preserving the existing structure of chapters and sections.

While not a single chapter remains completely unchanged, the bulk of the new material consists of (1) problem sets at the end of Chapters 2 through 8, (2) expanded discussion of linear logistic response models and polytomous response models in Chapter 6, (3) a further discussion of retrospective epidemiological studies in Chapter 7, and (4) a new appendix on the small-sample behavior of goodness-of-fit statistics. I have added briefer materials and references elsewhere and corrected several minor errors from the first edition. A relatively major correction has been made in connection with the theorem on collapsing tables in Section 3.8.

I gave considerable thought to the preparation of an additional appendix on computer programs for the analysis of categorical data, but in the end I resisted the temptation to do so. Many programs for maximum-likelihood estimation in connection with loglinear models are now in widespread use. These include the GLIM package prepared in England under the guidance of John Nelder and the sponsorship of the Royal Statistical Society, and various adaptations of iterative scaling programs originally prepared by Yvonne Bishop and Shelby Haberman (e.g., BMDP3F in the BMDP Programs distributed by the UCLA Health Sciences Computing Facility). Most users are likely to find one or more suitable programs available at their own computer installation that can be used to work through the examples and problems in this book. My primary reason for not providing any further guidance to computer programs is that I believe there will be major changes in both their availability and in the numerical methods they will be using within the next two to three years. Thus any explicit advice I could offer now would be out of date soon after the publication of the second edition.

Many friends, colleagues, and students provided me with suggestions, comments, and corrections for this edition. These include John Duffy, O. Dudley Duncan, David Hoaglin, J. G. Kalbfleisch, Kinley Larntz, S. Keith Lee, William Mason, Michael Meyer, Doug Ratcliff and Stanley Wasserman. The

preparation of this edition was partially supported by Office of Naval Research Contract N00014-78-C-0600 to the University of Minnesota.

For the typing and organization of the final manuscript, as well as for the updating of the indexes, I am indebted to Linda D. Anderson.

New Brighton, Minnesota
November 1979

Stephen E. Fienberg

Preface to the First Edition

The analysis of cross-classified categorical data has occupied a prominent place in introductory and intermediate-level statistical methods courses for many years, but with a few exceptions the only techniques presented in such courses have been those associated with the analysis of two-dimensional contingency tables and the calculation of chi-square statistics. During the past 15 years, advances in statistical theory and the ready availability of high-speed computers have led to major advances in the analysis of multi-dimensional cross-classified categorical data. Bishop, Fienberg, and Holland [1975], Cox [1970a], Haberman [1974a], Lindsey [1973], and Plackett [1974] have all presented detailed expositions of these new techniques, but these books are not directed primarily to the nonstatistical reader, whose background may be limited to one or two semesters of statistical methods at a noncalculus level.

The present monograph is intended as an introduction to the recent work on the analysis of cross-classified categorical data using loglinear models. I have written primarily for nonstatisticians, and Appendix I contains a summary of theoretical statistical terminology for such readers. Most of the material should be accessible to those who are familiar with the analysis of two-dimensional contingency tables, regression analysis, and analysis-of-variance models. The monograph also includes a variety of new methods based on loglinear models that have entered the statistical literature subsequent to the preparation of my book with Yvonne Bishop and Paul Holland. In particular, Chapter 4 contains a discussion of contingency tables with ordered categories for one or more of the variables, and Chapter 8 presents several new applications of the methods associated with incomplete contingency tables (i.e., tables with structural zeros).

Versions of material in this monograph were prepared in the form of notes to accompany lectures delivered in July 1972 at the Advanced Institute on Statistical Ecology held at Pennsylvania State University and during 1973 through 1975 at a series of Training Sessions on the Multivariate Analysis of Qualitative Data held at the University of Chicago. Various participants at these lectures have provided me with comments and suggestions that have found their way into the presentation here. Most of the final version of the monograph was completed while I was on sabbatical leave from the University of Minnesota and under partial support from National Science Foundation Grant SOC72-05257 to the Department of Statistics, Harvard University, and grants from the Robert Wood Johnson Foundation and the Commonwealth Fund to the Center for the Analysis of Health Practices, Harvard School of Public Health.

I am grateful to Stephen S. Brier, Michael L. Brown, Ron Christensen, David R. Cox, William Fairley, S. Keith Lee, William Mason, and Roy E. Welsch for extremely valuable comments and suggestions. Many people have provided me with examples and other materials, from both published and unpublished works, that have found their way into the final manuscript, including Albert Beaton, Richard Campbell, O. Dudley Duncan, Leo Goodman, Shelby Haberman, David Hoaglin, Kinley Larntz, Marc Nerlove, S. James Press, Ira Reiss, Thomas Schoener, and Sanford Weisberg. Most of all, I am indebted to Yvonne Bishop, Paul Holland, and Frederick Mosteller, whose collaboration over a period of many years helped to stimulate the present work.

For the typing and organization of the final manuscript, I wish to thank Sue Hange, Pat Haswell, Susan Kaufman, and Laurie Pearlman.

New Brighton, Minnesota

Stephen E. Fienberg

Contents

| | |
|---|-----------|
| Preface to the Second Edition | vii |
| Preface to the First Edition | ix |
| 1 Introduction | 1 |
| 1.1 The Analysis of Categorical Data | 1 |
| 1.2 Forms of Multivariate Analysis | 2 |
| 1.3 Some Historical Background | 4 |
| 1.4 A Medical Example | 6 |
| 2 Two-Dimensional Tables | 8 |
| 2.1 Two Binomials | 8 |
| 2.2 The Model of Independence | 10 |
| 2.3 The Loglinear Model | 13 |
| 2.4 Sampling Models | 15 |
| 2.5 The Cross-Product Ratio and 2×2 Tables | 16 |
| 2.6 Interrelated Two-Dimensional Tables | 20 |
| 2.7 Correction for Continuity | 21 |
| 2.8 Other Scales for Analyzing Two-Dimensional Tables | 22 |
| Problems | 23 |
| 3 Three-Dimensional Tables | 27 |
| 3.1 The General Loglinear Model | 27 |
| 3.2 Sampling Models | 29 |
| 3.3 Estimated Expected Values | 32 |
| 3.4 Iterative Computation of Expected Values | 37 |
| 3.5 Goodness-of-Fit Statistics | 40 |
| 3.6 Hierarchical Models | 43 |
| 3.7 A Further Example | 44 |
| 3.8 Collapsing Tables | 48 |
| Problems | 51 |
| 4 Selection of a Model | 56 |
| 4.1 General Issues | 56 |
| 4.2 Conditional Test Statistics | 56 |
| 4.3 Partitioning Chi-Square | 57 |
| 4.4 Using Information about Ordered Categories | 61 |
| Problems | 68 |

| | | |
|--------------|---|-----|
| 5 | Four- and Higher-Dimensional Contingency Tables | 71 |
| 5.1 | The Loglinear Models and MLEs for Expected Values | 71 |
| 5.2 | Using Partitioning to Select a Model | 74 |
| 5.3 | Stepwise Selection Procedures | 77 |
| 5.4 | Looking at All Possible Effects | 80 |
| | Problems | 88 |
| 6 | Fixed Margins and Logit Models | 95 |
| 6.1 | A Three-Dimensional Example | 95 |
| 6.2 | Logit Models | 97 |
| 6.3 | Logit Models and Ordered Categories | 99 |
| 6.4 | Linear Logistic Response Models | 102 |
| 6.5 | Logistic Regression vs. Discriminant Analysis | 105 |
| 6.6 | Polytomous and Multivariate Response Variables | 110 |
| | Problems | 116 |
| 7 | Causal Analysis Involving Logit and Loglinear Models | 120 |
| 7.1 | Path Diagrams | 120 |
| 7.2 | Recursive Systems of Logit Models | 123 |
| 7.3 | Recursive Systems: A More Complex Example | 129 |
| 7.4 | Nonrecursive Systems of Logit Models | 133 |
| 7.5 | Retrospective Epidemiological Studies | 135 |
| 7.6 | Logistic Models for Retrospective Data | 137 |
| | Problems | 138 |
| 8 | Fixed and Random Zeros | 140 |
| 8.1 | Sampling Zeros and MLEs in Loglinear Models | 140 |
| 8.2 | Incomplete Two-Dimensional Contingency Tables | 142 |
| 8.3 | Incompleteness in Several Dimensions | 146 |
| 8.4 | Some Applications of Incomplete Table Methodology | 150 |
| | Problems | 159 |
| Appendix I | Statistical Terminology | 165 |
| Appendix II | Basic Estimation Results for Loglinear Models | 167 |
| Appendix III | Percentage Points of χ^2 Distribution | 171 |
| Appendix IV | Small-Sample Properties of χ^2 Statistics | 172 |

| | |
|---------------|------|
| Contents | xiii |
| References | 177 |
| Author Index | 191 |
| Subject Index | 195 |