

# Complex Surveys

Parimal Mukhopadhyay

# Complex Surveys

Analysis of Categorical Data

 Springer

Parimal Mukhopadhyay  
Indian Statistical Institute  
Kolkata, West Bengal  
India

ISBN 978-981-10-0870-2      ISBN 978-981-10-0871-9 (eBook)  
DOI 10.1007/978-981-10-0871-9

Library of Congress Control Number: 2016936288

© Springer Science+Business Media Singapore 2016

This work is subject to copyright. All rights are reserved by the Publisher, whether the whole or part of the material is concerned, specifically the rights of translation, reprinting, reuse of illustrations, recitation, broadcasting, reproduction on microfilms or in any other physical way, and transmission or information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed.

The use of general descriptive names, registered names, trademarks, service marks, etc. in this publication does not imply, even in the absence of a specific statement, that such names are exempt from the relevant protective laws and regulations and therefore free for general use.

The publisher, the authors and the editors are safe to assume that the advice and information in this book are believed to be true and accurate at the date of publication. Neither the publisher nor the authors or the editors give a warranty, express or implied, with respect to the material contained herein or for any errors or omissions that may have been made.

Printed on acid-free paper

This Springer imprint is published by Springer Nature  
The registered company is Springer Science+Business Media Singapore Pte Ltd.

*To*  
*The Loving Memory of My Wife*  
Manju

# Preface

Most of the data a statistician uses is categorical in nature. In the realms of biomedicine and social sciences, ecology and demography, voting pattern and marketing, to name a few, categorical data dominate. They are the major raw materials for analysis for testing different hypotheses relating to the populations which generate such data. Such data are generally obtained from surveys carried under a complex survey design, generally a stratified multistage design. In analysis of data collected through sample surveys, standard statistical techniques are often routinely employed.

We recall the well-quoted phrase which we chanted in our undergraduate days: Let  $x_1, \dots, x_n$  be a random sample of size  $n$  drawn from a population with probability density function  $f(x)$  or probability mass function  $p_M(x)$ , etc. This means that the sampled units whose observations are  $x_1, x_2, \dots, x_n$ , are drawn by simple random sampling with replacement (*srswr*). This also implies that observed variables  $x_1, \dots, x_n$  are independently and identically distributed (IID). In fact, most of the results in theoretical statistics, including those in usual analysis of categorical data, are based on these assumptions.

However, survey populations are often complex with different cell probabilities in different subgroups of the population, and this implies a situation different from the IID setup. Longitudinal surveys—where sample subjects are observed over two or more time points—typically lead to dependent observations over time. Moreover, longitudinal surveys often have complex survey designs that involve clustering which results in cross-sectional dependence among samples.

In view of these facts, it is, therefore, necessary to modify the usual tests of goodness of fit, like Pearsonian chi-square, likelihood ratio and Wald statistic to make them suitable for use in the context of data from complex surveys. A host of ardent researchers have developed a number of such tests for this purpose over more than last four decades.

There are already a myriad number of textbooks and research monographs on analysis of categorical data. Then why is another book in the area required? My humble answer is that all those treatise provide excellent description of the

categorical data analysis under the classical setup (usual *srswr* or IID assumption), but none addresses the problem when the data are obtained through complex sample survey designs, which more often than not fail to satisfy the usual assumptions. The present endeavor tries to fill in the gap in the area.

The idea of writing this book is, therefore, to make a review of some of the ideas that have blown out in the field of analysis of categorical data from complex surveys. In doing so, I have tried to systematically arrange the results and provide relevant examples to illuminate the ideas. This research monograph is a review of the works already done in the area and does not offer any new investigation. As such I have unhesitatingly used a host of brilliant publications in this area. A brief outline of different chapters is indicated as under:

- (1) Chapter 1: Basic ideas of sampling; finite population; sampling design; estimator; different sampling strategies; design-based method of making inference; superpopulation model; model-based inference
- (2) Chapter 2: Effects of a true complex design on the variance of an estimator with reference to a *srswr* design or an IID-model setup; design effects; misspecification effects; multivariate design effect; nonparametric variance estimation
- (3) Chapter 3: Review of classical models of categorical data; tests of hypotheses for goodness of fit; log-linear model; logistic regression model
- (4) Analysis of categorical data from complex surveys under full or saturated models; different goodness-of-fit tests and their modifications
- (5) Analysis of categorical data from complex surveys under log-linear model; different goodness-of-fit tests and their modifications
- (6) Analysis of categorical data from complex surveys under binomial and polytomous logistic regression model; different goodness-of-fit tests and their modifications
- (7) Analysis of categorical data from complex surveys when misclassification errors are present; different goodness-of-fit tests and their modifications
- (8) Some procedures for obtaining approximate maximum likelihood estimators; pseudo-likelihood approach for estimation of finite population parameters; design-adjusted estimators; mixed model framework; principal component analysis
- (9) Appendix: Asymptotic properties of multinomial distribution; asymptotic distribution of different goodness-of-fit tests; Neyman's (1949) and Wald's (1943) procedures for testing general hypotheses relating to population proportions

I gratefully acknowledge my indebtedness to the authorities of PHI Learning, New Delhi, India, for kindly allowing me to use a part of my book, *Theory and Methods of Survey Sampling*, 2nd ed., 2009, in Chap. 2 of the present book. I am thankful to Mr. Shamin Ahmad, Senior Editor for Mathematical Sciences at

Springer, New Delhi, for his kind encouragement. The book was prepared at the Indian Statistical Institute, Kolkata, to the authorities of which I acknowledge my thankfulness. And last but not the least, I must acknowledge my indebtedness to my family for their silent encouragement and support throughout this project.

January 2016

Parimal Mukhopadhyay

# Contents

<b>1 Preliminaries</b> . . . . .	1
1.1 Introduction . . . . .	1
1.2 The Fixed Population Model. . . . .	2
1.3 Different Types of Sampling Designs. . . . .	8
1.4 The Estimators . . . . .	11
1.4.1 A Class of Estimators . . . . .	14
1.5 Model-Based Approach to Sampling . . . . .	17
1.5.1 Uses of Design and Model in Sampling . . . . .	23
1.6 Plan of the Book. . . . .	24
<b>2 The Design Effects and Misspecification Effects</b> . . . . .	27
2.1 Introduction . . . . .	28
2.2 Effect of a Complex Design on the Variance of an Estimator . . . . .	30
2.3 Effect of a Complex Design on Confidence Interval for $\theta$ . . . . .	37
2.4 Multivariate Design Effects. . . . .	40
2.5 Nonparametric Methods of Variance Estimation . . . . .	41
2.5.1 A Simple Method of Estimation of Variance of a Linear Statistic . . . . .	41
2.5.2 Linearization Method for Variance Estimation of a Nonlinear Estimator . . . . .	45
2.5.3 Random Group Method . . . . .	48
2.5.4 Balanced Repeated Replications . . . . .	52
2.5.5 The Jackknife Procedures. . . . .	56
2.5.6 The Jackknife Repeated Replication (JRR) . . . . .	58
2.5.7 The Bootstrap. . . . .	60
2.6 Effect of Survey Design on Inference About Covariance Matrix. . . . .	64
2.7 Exercises and Complements . . . . .	65



**3 Some Classical Models in Categorical Data Analysis . . . . . 67**

3.1 Introduction . . . . . 67

3.2 Statistical Models . . . . . 68

3.2.1 Fitting Statistical Models . . . . . 69

3.2.2 Large Sample Estimation Theory. . . . . 69

3.2.3 Asymptotic Properties of ML Estimates . . . . . 70

3.2.4 Testing of Parameters . . . . . 71

3.2.5 Transformation of the Central Limit Theorem. . . . . 73

3.3 Distribution Theory for Count Data . . . . . 73

3.3.1 Multinomial Models . . . . . 73

3.3.2 Poisson Models. . . . . 76

3.3.3 The Multinomial-Poisson Connection. . . . . 78

3.4 Goodness-of-Fit. . . . . 78

3.4.1 Likelihood Ratio Statistic. . . . . 78

3.4.2 Pearson’s Goodness-of-Fit Statistic . . . . . 79

3.5 Binomial Data. . . . . 81

3.5.1 Binomial Data and the Log-Likelihood Ratio . . . . . 81

3.6 Log-Linear Models . . . . . 82

3.6.1 Log-Linear Models for Two-Way Tables . . . . . 82

3.6.2 Log-Linear Models for Three-Way Tables . . . . . 84

3.7 Logistic Regression Analysis . . . . . 89

3.7.1 The Logistic Regression Model. . . . . 89

3.7.2 Fitting the Logistic Regression Model . . . . . 90

3.7.3 The Multiple Logistic Regression Model . . . . . 91

3.7.4 Fitting the Multiple Logistic Regression Model. . . . . 92

3.7.5 Polytomous Logistic Regression . . . . . 94

3.8 Fitting the Logistic Regression Models to Data  
from Complex Surveys. . . . . 95

**4 Analysis of Categorical Data Under a Full Model . . . . . 97**

4.1 Introduction . . . . . 97

4.2 Tests of Goodness-of-Fit . . . . . 98

4.2.1 Pearsonian Chi-Square Statistic. . . . . 98

4.2.2 Design-Based Wald Statistic. . . . . 99

4.2.3 Neyman’s (Multinomial Wald) Statistic . . . . . 101

4.2.4 Log-Likelihood Ratio Statistic . . . . . 101

4.2.5 Asymptotic Distribution of  $X_W^2$  and  $X_P^2$ . . . . . 101

4.2.6 Generalized Design Effect of  $\pi$ . . . . . 107

4.2.7 Modification to  $X_P^2$  . . . . . 109

4.2.8 Fay’s Jackknifed Statistic. . . . . 112

4.3 Testing for Homogeneity . . . . . 115

4.3.1 A Simple Method for Binary Data  
from Cluster Sampling. . . . . 120

4.4 Effect of Survey Design on Classical Tests  
of General Linear Hypotheses. . . . . 121

4.5 Tests of Independence in a Two-Way Table . . . . . 123

4.6 Some Evaluation of Tests Under Cluster Sampling . . . . . 128

4.7 Exercises and Complements . . . . . 130

**5 Analysis of Categorical Data Under Log-Linear Models . . . . . 135**

5.1 Introduction . . . . . 135

5.2 Log-Linear Models in Contingency Tables . . . . . 136

5.3 Tests for Goodness of Fit . . . . . 138

5.3.1 Other Standard Tests and Their First- and Second-Order Corrections . . . . . 139

5.3.2 Fay’s Jackknifed Tests . . . . . 143

5.4 Asymptotic Covariance Matrix of the Pseudo-MLE  $\hat{\pi}$  . . . . . 145

5.4.1 Residual Analysis . . . . . 146

5.5 Brier’s Model . . . . . 147

5.6 Nested Models . . . . . 149

5.6.1 Pearsonian Chi-Square and the Likelihood Ratio Statistic . . . . . 150

5.6.2 A Wald Statistic . . . . . 153

5.6.3 Modifications to Test Statistics . . . . . 154

5.6.4 Effects of Survey Design on  $X^2_p(2|1)$  . . . . . 155

**6 Analysis of Categorical Data Under Logistic Regression Model . . . . . 157**

6.1 Introduction . . . . . 157

6.2 Binary Logistic Regression . . . . . 158

6.2.1 Pseudo-MLE of  $\pi$  . . . . . 159

6.2.2 Asymptotic Covariance Matrix of the Estimators . . . . . 160

6.2.3 Goodness-of-Fit Tests . . . . . 163

6.2.4 Modifications of Tests . . . . . 164

6.3 Nested Model . . . . . 164

6.3.1 A Wald Statistic . . . . . 167

6.3.2 Modifications to Tests . . . . . 168

6.4 Choosing Appropriate Cell-Sample Sizes for Running Logistic Regression Program in a Standard Computer Package . . . . . 169

6.5 Model in the Polytomous Case . . . . . 170

6.6 Analysis Under Generalized Least Square Approach . . . . . 172

6.7 Exercises and Complements . . . . . 176

**7 Analysis in the Presence of Classification Errors . . . . . 179**

7.1 Introduction . . . . . 179

7.2 Tests for Goodness-of-Fit Under Misclassification . . . . . 180

7.2.1 Methods for Considering Misclassification Under SRS . . . . . 180

7.2.2 Methods for General Sampling Designs . . . . . 182

7.2.3 A Model-Free Approach . . . . . 183

- 7.3 Tests of Independence Under Misclassification . . . . . 185
  - 7.3.1 Methods for Considering Misclassification Under SRS. . . . . 186
  - 7.3.2 Methods for Arbitrary Survey Designs. . . . . 186
- 7.4 Test of Homogeneity . . . . . 188
- 7.5 Analysis Under Weighted Cluster Sample Design . . . . . 192
- 8 Approximate MLE from Survey Data. . . . . 195**
  - 8.1 Introduction . . . . . 195
  - 8.2 Exact MLE from Survey Data. . . . . 196
    - 8.2.1 Ignorable Sampling Designs . . . . . 196
    - 8.2.2 Exact MLE . . . . . 197
  - 8.3 MLE’s Derived from Weighted Distributions . . . . . 198
  - 8.4 Design-Adjusted Maximum Likelihood Estimation. . . . . 200
    - 8.4.1 Design-Adjusted Regression Estimation  
with Selectivity Bias . . . . . 205
  - 8.5 The Pseudo-Likelihood Approach to MLE  
from Complex Surveys. . . . . 208
    - 8.5.1 Analysis Based on Generalized Linear Model . . . . . 209
    - 8.5.2 Estimation for Linear Models . . . . . 212
  - 8.6 A Mixed (Design-Model) Framework. . . . . 216
  - 8.7 Effect of Survey Design on Multivariate Analysis  
of Principal Components . . . . . 220
    - 8.7.1 Estimation of Principal Components . . . . . 222
- Appendix A: Asymptotic Properties of Multinomial Distribution . . . . . 223**
- References . . . . . 241**

## About the Author

**Parimal Mukhopadhyay** is a former professor of statistics at the Indian Statistical Institute, Kolkata, India. He received his Ph.D. degree in statistics from the University of Calcutta in 1977. He also served as a faculty member at the University of Ife, Nigeria, Moi University, Kenya, University of South Pacific, Fiji Islands and held visiting positions at the University of Montreal, University of Windsor, Stockholm University and the University of Western Australia. He has written more than 75 research papers in survey sampling, some co-authored and eleven books on statistics. He was a member of the Institute of Mathematical Statistics and elected member of the International Statistical Institute.