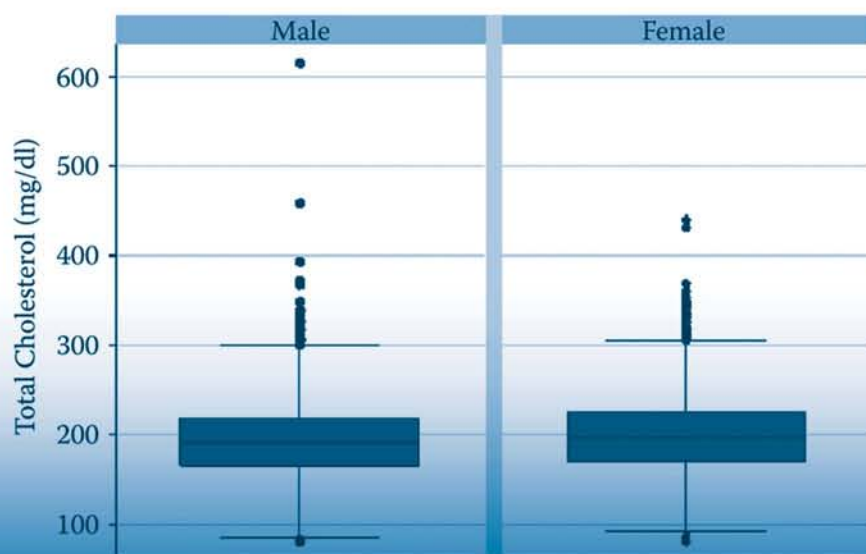


Chapman & Hall/CRC  
Statistics in the Social and Behavioral Sciences Series

# Applied Survey Data Analysis



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



CRC Press  
Taylor & Francis Group

A CHAPMAN & HALL BOOK

# **Applied Survey Data Analysis**

# Chapman & Hall/CRC

## Statistics in the Social and Behavioral Sciences Series

### Series Editors

**A. Colin Cameron**  
University of California, Davis, USA

**J. Scott Long**  
Indiana University, USA

**Andrew Gelman**  
Columbia University, USA

**Sophia Rabe-Hesketh**  
University of California, Berkeley, USA

**Anders Skrondal**  
Norwegian Institute of Public Health, Norway

### 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*

**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*

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

# Applied Survey Data Analysis

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



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

© 2010 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: 20140514

International Standard Book Number-13: 978-1-4200-8067-4 (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) (<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

---

Preface.....xv

**1. Applied Survey Data Analysis: Overview ..... 1**

    1.1 Introduction ..... 1

    1.2 A Brief History of Applied Survey Data Analysis ..... 3

        1.2.1 Key Theoretical Developments ..... 3

        1.2.2 Key Software Developments ..... 5

    1.3 Example Data Sets and Exercises ..... 6

        1.3.1 The National Comorbidity Survey Replication (NCS-R)..... 6

        1.3.2 The Health and Retirement Study (HRS)—2006 ..... 7

        1.3.3 The National Health and Nutrition Examination Survey (NHANES)—2005, 2006..... 7

        1.3.4 Steps in Applied Survey Data Analysis..... 8

            1.3.4.1 Step 1: Definition of the Problem and Statement of the Objectives..... 8

            1.3.4.2 Step 2: Understanding the Sample Design ..... 9

            1.3.4.3 Step 3: Understanding Design Variables, Underlying Constructs, and Missing Data..... 10

            1.3.4.4 Step 4: Analyzing the Data ..... 11

            1.3.4.5 Step 5: Interpreting and Evaluating the Results of the Analysis ..... 11

            1.3.4.6 Step 6: Reporting of Estimates and Inferences from the Survey Data ..... 12

**2. Getting to Know the Complex Sample Design ..... 13**

    2.1 Introduction ..... 13

        2.1.1 Technical Documentation and Supplemental Literature Review..... 13

    2.2 Classification of Sample Designs ..... 14

        2.2.1 Sampling Plans..... 15

        2.2.2 Inference from Survey Data ..... 16

    2.3 Target Populations and Survey Populations..... 16

    2.4 Simple Random Sampling: A Simple Model for Design-Based Inference..... 18

        2.4.1 Relevance of SRS to Complex Sample Survey Data Analysis..... 18

        2.4.2 SRS Fundamentals: A Framework for Design-Based Inference..... 19

        2.4.3 An Example of Design-Based Inference under SRS ..... 21

- 2.5 Complex Sample Design Effects .....23
  - 2.5.1 Design Effect Ratio .....23
  - 2.5.2 Generalized Design Effects and Effective Sample Sizes .....25
- 2.6 Complex Samples: Clustering and Stratification .....27
  - 2.6.1 Clustered Sampling Plans .....28
  - 2.6.2 Stratification.....31
  - 2.6.3 Joint Effects of Sample Stratification and Clustering.....34
- 2.7 Weighting in Analysis of Survey Data.....35
  - 2.7.1 Introduction to Weighted Analysis of Survey Data.....35
  - 2.7.2 Weighting for Probabilities of Selection .....37
  - 2.7.3 Nonresponse Adjustment Weights .....39
    - 2.7.3.1 Weighting Class Approach .....40
    - 2.7.3.2 Propensity Cell Adjustment Approach.....40
  - 2.7.4 Poststratification Weight Factors .....42
  - 2.7.5 Design Effects Due to Weighted Analysis .....44
- 2.8 Multistage Area Probability Sample Designs.....46
  - 2.8.1 Primary Stage Sampling .....47
  - 2.8.2 Secondary Stage Sampling .....48
  - 2.8.3 Third and Fourth Stage Sampling of Housing Units and Eligible Respondents .....49
- 2.9 Special Types of Sampling Plans Encountered in Surveys.....50
- 3. Foundations and Techniques for Design-Based Estimation and Inference.....53**
  - 3.1 Introduction .....53
  - 3.2 Finite Populations and Superpopulation Models .....54
  - 3.3 Confidence Intervals for Population Parameters .....56
  - 3.4 Weighted Estimation of Population Parameters.....56
  - 3.5 Probability Distributions and Design-Based Inference .....60
    - 3.5.1 Sampling Distributions of Survey Estimates.....60
    - 3.5.2 Degrees of Freedom for  $t$  under Complex Sample Designs.....63
  - 3.6 Variance Estimation.....65
    - 3.6.1 Simplifying Assumptions Employed in Complex Sample Variance Estimation.....66
    - 3.6.2 The Taylor Series Linearization Method .....68
      - 3.6.2.1 TSL Step 1 .....69
      - 3.6.2.2 TSL Step 2 .....70
      - 3.6.2.3 TSL Step 3 .....71
      - 3.6.2.4 TSL Step 4 .....71
      - 3.6.2.5 TSL Step 5 .....73
    - 3.6.3 Replication Methods for Variance Estimation.....74
      - 3.6.3.1 Jackknife Repeated Replication.....75

- 3.6.3.2 Balanced Repeated Replication ..... 78
      - 3.6.3.3 The Bootstrap ..... 82
    - 3.6.4 An Example Comparing the Results from the TSL,  
JRR, and BRR Methods ..... 82
  - 3.7 Hypothesis Testing in Survey Data Analysis ..... 83
  - 3.8 Total Survey Error and Its Impact on Survey Estimation and  
Inference ..... 85
    - 3.8.1 Variable Errors ..... 86
    - 3.8.2 Biases in Survey Data ..... 87
- 4. Preparation for Complex Sample Survey Data Analysis ..... 91**
  - 4.1 Introduction ..... 91
  - 4.2 Analysis Weights: Review by the Data User ..... 92
    - 4.2.1 Identification of the Correct Weight Variables for the  
Analysis ..... 93
    - 4.2.2 Determining the Distribution and Scaling of the  
Weight Variables ..... 94
    - 4.2.3 Weighting Applications: Sensitivity of Survey  
Estimates to the Weights ..... 96
  - 4.3 Understanding and Checking the Sampling Error  
Calculation Model ..... 98
    - 4.3.1 Stratum and Cluster Codes in Complex Sample  
Survey Data Sets ..... 99
    - 4.3.2 Building the NCS-R Sampling Error Calculation  
Model ..... 100
    - 4.3.3 Combining Strata, Randomly Grouping PSUs, and  
Collapsing Strata ..... 103
    - 4.3.4 Checking the Sampling Error Calculation Model for  
the Survey Data Set ..... 105
  - 4.4 Addressing Item Missing Data in Analysis Variables ..... 108
    - 4.4.1 Potential Bias Due to Ignoring Missing Data ..... 108
    - 4.4.2 Exploring Rates and Patterns of Missing Data Prior  
to Analysis ..... 109
  - 4.5 Preparing to Analyze Data for Sample Subpopulations ..... 110
    - 4.5.1 Subpopulation Distributions across Sample Design  
Units ..... 111
    - 4.5.2 The Unconditional Approach for Subclass Analysis ..... 114
    - 4.5.3 Preparation for Subclass Analyses ..... 114
  - 4.6 A Final Checklist for Data Users ..... 115
- 5. Descriptive Analysis for Continuous Variables ..... 117**
  - 5.1 Introduction ..... 117
  - 5.2 Special Considerations in Descriptive Analysis of Complex  
Sample Survey Data ..... 118
    - 5.2.1 Weighted Estimation ..... 118



- 5.2.2 Design Effects for Descriptive Statistics..... 119
    - 5.2.3 Matching the Method to the Variable Type ..... 119
  - 5.3 Simple Statistics for Univariate Continuous Distributions..... 120
    - 5.3.1 Graphical Tools for Descriptive Analysis of Survey Data ..... 120
    - 5.3.2 Estimation of Population Totals..... 123
    - 5.3.3 Means of Continuous, Binary, or Interval Scale Data..... 128
    - 5.3.4 Standard Deviations of Continuous Variables ..... 130
    - 5.3.5 Estimation of Percentiles and Medians of Population Distributions..... 131
  - 5.4 Bivariate Relationships between Two Continuous Variables .... 134
    - 5.4.1 X-Y Scatterplots..... 134
    - 5.4.2 Product Moment Correlation Statistic ( $r$ )..... 135
    - 5.4.3 Ratios of Two Continuous Variables ..... 136
  - 5.5 Descriptive Statistics for Subpopulations..... 137
  - 5.6 Linear Functions of Descriptive Estimates and Differences of Means ..... 139
    - 5.6.1 Differences of Means for Two Subpopulations ..... 141
    - 5.6.2 Comparing Means over Time ..... 143
  - 5.7 Exercises ..... 144
- 6. Categorical Data Analysis ..... 149**
  - 6.1 Introduction ..... 149
  - 6.2 A Framework for Analysis of Categorical Survey Data ..... 150
    - 6.2.1 Incorporating the Complex Design and Pseudo-Maximum Likelihood..... 150
    - 6.2.2 Proportions and Percentages..... 150
    - 6.2.3 Cross-Tabulations, Contingency Tables, and Weighted Frequencies ..... 151
  - 6.3 Univariate Analysis of Categorical Data ..... 152
    - 6.3.1 Estimation of Proportions for Binary Variables ..... 152
    - 6.3.2 Estimation of Category Proportions for Multinomial Variables ..... 156
    - 6.3.3 Testing Hypotheses Concerning a Vector of Population Proportions..... 158
    - 6.3.4 Graphical Display for a Single Categorical Variable..... 159
  - 6.4 Bivariate Analysis of Categorical Data ..... 160
    - 6.4.1 Response and Factor Variables ..... 160
    - 6.4.2 Estimation of Total, Row, and Column Proportions for Two-Way Tables..... 162
    - 6.4.3 Estimating and Testing Differences in Subpopulation Proportions ..... 163
    - 6.4.4 Chi-Square Tests of Independence of Rows and Columns ..... 164
    - 6.4.5 Odds Ratios and Relative Risks..... 170

- 6.4.6 Simple Logistic Regression to Estimate the Odds Ratio ..... 171
    - 6.4.7 Bivariate Graphical Analysis..... 173
  - 6.5 Analysis of Multivariate Categorical Data ..... 174
    - 6.5.1 The Cochran–Mantel–Haenszel Test ..... 174
    - 6.5.2 Log-Linear Models for Contingency Tables..... 176
  - 6.6 Exercises ..... 177
- 7. Linear Regression Models ..... 179**
  - 7.1 Introduction ..... 179
  - 7.2 The Linear Regression Model ..... 180
    - 7.2.1 The Standard Linear Regression Model ..... 182
    - 7.2.2 Survey Treatment of the Regression Model..... 183
  - 7.3 Four Steps in Linear Regression Analysis..... 185
    - 7.3.1 Step 1: Specifying and Refining the Model..... 186
    - 7.3.2 Step 2: Estimation of Model Parameters..... 187
      - 7.3.2.1 Estimation for the Standard Linear Regression Model ..... 187
      - 7.3.2.2 Linear Regression Estimation for Complex Sample Survey Data ..... 188
    - 7.3.3 Step 3: Model Evaluation ..... 193
      - 7.3.3.1 Explained Variance and Goodness of Fit..... 193
      - 7.3.3.2 Residual Diagnostics..... 194
      - 7.3.3.3 Model Specification and Homogeneity of Variance ..... 194
      - 7.3.3.4 Normality of the Residual Errors..... 195
      - 7.3.3.5 Outliers and Influence Statistics ..... 196
    - 7.3.4 Step 4: Inference ..... 196
      - 7.3.4.1 Inference Concerning Model Parameters ..... 199
      - 7.3.4.2 Prediction Intervals..... 202
  - 7.4 Some Practical Considerations and Tools..... 204
    - 7.4.1 Distribution of the Dependent Variable ..... 204
    - 7.4.2 Parameterization and Scaling for Independent Variables ..... 205
    - 7.4.3 Standardization of the Dependent and Independent Variables ..... 208
    - 7.4.4 Specification and Interpretation of Interactions and Nonlinear Relationships ..... 208
    - 7.4.5 Model-Building Strategies ..... 210
  - 7.5 Application: Modeling Diastolic Blood Pressure with the NHANES Data ..... 211
    - 7.5.1 Exploring the Bivariate Relationships ..... 212
    - 7.5.2 Naïve Analysis: Ignoring Sample Design Features ..... 215
    - 7.5.3 Weighted Regression Analysis ..... 216

7.5.4	Appropriate Analysis: Incorporating All Sample Design Features.....	218
7.6	Exercises .....	224
<b>8.</b>	<b>Logistic Regression and Generalized Linear Models for Binary Survey Variables .....</b>	<b>229</b>
8.1	Introduction .....	229
8.2	Generalized Linear Models for Binary Survey Responses.....	230
8.2.1	The Logistic Regression Model.....	231
8.2.2	The Probit Regression Model .....	234
8.2.3	The Complementary Log–Log Model.....	234
8.3	Building the Logistic Regression Model: Stage 1, Model Specification .....	235
8.4	Building the Logistic Regression Model: Stage 2, Estimation of Model Parameters and Standard Errors.....	236
8.5	Building the Logistic Regression Model: Stage 3, Evaluation of the Fitted Model.....	239
8.5.1	Wald Tests of Model Parameters .....	239
8.5.2	Goodness of Fit and Logistic Regression Diagnostics.....	243
8.6	Building the Logistic Regression Model: Stage 4, Interpretation and Inference .....	245
8.7	Analysis Application .....	251
8.7.1	Stage 1: Model Specification .....	252
8.7.2	Stage 2: Model Estimation .....	253
8.7.3	Stage 3: Model Evaluation.....	255
8.7.4	Stage 4: Model Interpretation/Inference .....	256
8.8	Comparing the Logistic, Probit, and Complementary Log–Log GLMs for Binary Dependent Variables .....	259
8.9	Exercises .....	262
<b>9.</b>	<b>Generalized Linear Models for Multinomial, Ordinal, and Count Variables.....</b>	<b>265</b>
9.1	Introduction .....	265
9.2	Analyzing Survey Data Using Multinomial Logit Regression Models.....	265
9.2.1	The Multinomial Logit Regression Model .....	265
9.2.2	Multinomial Logit Regression Model: Specification Stage.....	267
9.2.3	Multinomial Logit Regression Model: Estimation Stage.....	268
9.2.4	Multinomial Logit Regression Model: Evaluation Stage.....	268

- 9.2.5 Multinomial Logit Regression Model: Interpretation Stage..... 270
    - 9.2.6 Example: Fitting a Multinomial Logit Regression Model to Complex Sample Survey Data..... 271
  - 9.3 Logistic Regression Models for Ordinal Survey Data ..... 277
    - 9.3.1 Cumulative Logit Regression Model ..... 278
    - 9.3.2 Cumulative Logit Regression Model: Specification Stage ..... 279
    - 9.3.3 Cumulative Logit Regression Model: Estimation Stage ..... 279
    - 9.3.4 Cumulative Logit Regression Model: Evaluation Stage ..... 280
    - 9.3.5 Cumulative Logit Regression Model: Interpretation Stage ..... 281
    - 9.3.6 Example: Fitting a Cumulative Logit Regression Model to Complex Sample Survey Data..... 282
  - 9.4 Regression Models for Count Outcomes ..... 286
    - 9.4.1 Survey Count Variables and Regression Modeling Alternatives..... 286
    - 9.4.2 Generalized Linear Models for Count Variables..... 288
      - 9.4.2.1 The Poisson Regression Model..... 288
      - 9.4.2.2 The Negative Binomial Regression Model ..... 289
      - 9.4.2.3 Two-Part Models: Zero-Inflated Poisson and Negative Binomial Regression Models ..... 290
    - 9.4.3 Regression Models for Count Data: Specification Stage..... 291
    - 9.4.4 Regression Models for Count Data: Estimation Stage..... 292
    - 9.4.5 Regression Models for Count Data: Evaluation Stage..... 292
    - 9.4.6 Regression Models for Count Data: Interpretation Stage..... 293
    - 9.4.7 Example: Fitting Poisson and Negative Binomial Regression Models to Complex Sample Survey Data..... 294
  - 9.5 Exercises ..... 298
- 10. Survival Analysis of Event History Survey Data ..... 303**
  - 10.1 Introduction ..... 303
  - 10.2 Basic Theory of Survival Analysis..... 303
    - 10.2.1 Survey Measurement of Event History Data ..... 303
    - 10.2.2 Data for Event History Models ..... 305
    - 10.2.3 Important Notation and Definitions..... 306
    - 10.2.4 Models for Survival Analysis..... 307

- 10.3 (Nonparametric) Kaplan–Meier Estimation of the Survivor Function.....308
  - 10.3.1 K–M Model Specification and Estimation.....309
  - 10.3.2 K–M Estimator—Evaluation and Interpretation ..... 310
  - 10.3.3 K–M Survival Analysis Example..... 311
- 10.4 Cox Proportional Hazards Model ..... 315
  - 10.4.1 Cox Proportional Hazards Model: Specification..... 315
  - 10.4.2 Cox Proportional Hazards Model: Estimation Stage ..... 316
  - 10.4.3 Cox Proportional Hazards Model: Evaluation and Diagnostics..... 317
  - 10.4.4 Cox Proportional Hazards Model: Interpretation and Presentation of Results..... 319
  - 10.4.5 Example: Fitting a Cox Proportional Hazards Model to Complex Sample Survey Data ..... 319
- 10.5 Discrete Time Survival Models.....322
  - 10.5.1 The Discrete Time Logistic Model .....323
  - 10.5.2 Data Preparation for Discrete Time Survival Models ..... 324
  - 10.5.3 Discrete Time Models: Estimation Stage.....327
  - 10.5.4 Discrete Time Models: Evaluation and Interpretation.....328
  - 10.5.5 Fitting a Discrete Time Model to Complex Sample Survey Data.....329
- 10.6 Exercises .....333
- 11. Multiple Imputation: Methods and Applications for Survey Analysts .....335**
  - 11.1 Introduction .....335
  - 11.2 Important Missing Data Concepts .....336
    - 11.2.1 Sources and Patterns of Item-Missing Data in Surveys .....336
    - 11.2.2 Item-Missing Data Mechanisms.....338
    - 11.2.3 Implications of Item-Missing Data for Survey Data Analysis.....341
    - 11.2.4 Review of Strategies to Address Item-Missing Data in Surveys.....342
  - 11.3 An Introduction to Imputation and the Multiple Imputation Method.....345
    - 11.3.1 A Brief History of Imputation Procedures.....345
    - 11.3.2 Why the Multiple Imputation Method?.....346
    - 11.3.3 Overview of Multiple Imputation and MI Phases .....348
  - 11.4 Models for Multiply Imputing Missing Data.....350
    - 11.4.1 Choosing the Variables to Include in the Imputation Model.....350

- 11.4.2 Distributional Assumptions for the Imputation Model.....352
  - 11.5 Creating the Imputations.....353
    - 11.5.1 Transforming the Imputation Problem to Monotonic Missing Data.....353
    - 11.5.2 Specifying an Explicit Multivariate Model and Applying Exact Bayesian Posterior Simulation Methods.....354
    - 11.5.3 Sequential Regression or “Chained Regressions” .....354
  - 11.6 Estimation and Inference for Multiply Imputed Data.....355
    - 11.6.1 Estimators for Population Parameters and Associated Variance Estimators .....356
    - 11.6.2 Model Evaluation and Inference.....357
  - 11.7 Applications to Survey Data.....359
    - 11.7.1 Problem Definition .....359
    - 11.7.2 The Imputation Model for the NHANES Blood Pressure Example.....360
    - 11.7.3 Imputation of the Item-Missing Data.....361
    - 11.7.4 Multiple Imputation Estimation and Inference.....363
      - 11.7.4.1 Multiple Imputation Analysis 1: Estimation of Mean Diastolic Blood Pressure .....364
      - 11.7.4.2 Multiple Imputation Analysis 2: Estimation of the Linear Regression Model for Diastolic Blood Pressure .....365
  - 11.8 Exercises .....368
- 12. Advanced Topics in the Analysis of Survey Data.....371**
  - 12.1 Introduction .....371
  - 12.2 Bayesian Analysis of Complex Sample Survey Data .....372
  - 12.3 Generalized Linear Mixed Models (GLMMs) in Survey Data Analysis.....375
    - 12.3.1 Overview of Generalized Linear Mixed Models .....375
    - 12.3.2 Generalized Linear Mixed Models and Complex Sample Survey Data .....379
    - 12.3.3 GLMM Approaches to Analyzing Longitudinal Survey Data.....382
    - 12.3.4 Example: Longitudinal Analysis of the HRS Data .....389
    - 12.3.5 Directions for Future Research.....395
  - 12.4 Fitting Structural Equation Models to Complex Sample Survey Data.....395
  - 12.5 Small Area Estimation and Complex Sample Survey Data .....396
  - 12.6 Nonparametric Methods for Complex Sample Survey Data .....397
- Appendix A: Software Overview .....399**
  - A.1 Introduction .....399

- A.1.1 Historical Perspective.....400
  - A.1.2 Software for Sampling Error Estimation..... 401
- A.2 Overview of Stata® Version 10+ ..... 407
- A.3 Overview of SAS® Version 9.2 ..... 410
  - A.3.1 The SAS SURVEY Procedures..... 411
- A.4 Overview of SUDAAN® Version 9.0..... 414
  - A.4.1 The SUDAAN Procedures..... 415
- A.5 Overview of SPSS® ..... 421
  - A.5.1 The SPSS Complex Samples Commands..... 422
- A.6 Overview of Additional Software ..... 427
  - A.6.1 WesVar® ..... 427
  - A.6.2 IVEware (Imputation and Variance Estimation Software) ..... 428
  - A.6.3 Mplus ..... 429
  - A.6.4 The R survey Package ..... 429
- A.7 Summary ..... 430
- References ..... 431**
- Index ..... 443**

---

## Preface

---

This book is written as a guide to the applied statistical analysis and interpretation of survey data. The motivation for this text lies in years of teaching graduate courses in applied methods for survey data analysis and extensive consultation with social and physical scientists, educators, medical researchers, and public health professionals on best methods for approaching specific analysis questions using survey data. The general outline for this text is based on the syllabus for a course titled “Analysis of Complex Sample Survey Data” that we have taught for over 10 years in the Joint Program in Survey Methodology (JPSM) based at the University of Maryland (College Park) and in the University of Michigan’s Program in Survey Methodology (MPSM) and Summer Institute in Survey Research Techniques.

Readers may initially find the topical outline and content choices a bit unorthodox, but our instructional experience has shown it to be effective for teaching this complex subject to students and professionals who have a minimum of a two-semester graduate level course in applied statistics. The practical, everyday relevance of the chosen topics and the emphasis each receives in this text has also been informed by over 60 years of combined experience in consulting on survey data analysis with research colleagues and students under the auspices of the Survey Methodology Program of the Institute for Social Research (ISR) and the University of Michigan Center for Statistical Consultation and Research (CSCAR). For example, the emphasis placed on topics as varied as weighted estimation of population quantities, sampling error calculation models, coding of indicator variables in regression models, and interpretation of results from generalized linear models derives directly from our long-term observation of how often naïve users make critical mistakes in these areas.

This text, like our courses that it will serve, is designed to provide an intermediate-level statistical overview of the analysis of complex sample survey data—emphasizing methods and worked examples while reinforcing the principles and theory that underly those methods. The intended audience includes graduate students, survey practitioners, and research scientists from the wide array of disciplines that use survey data in their work. Students and practitioners in the statistical sciences should also find that this text provides a useful framework for integrating their further, more in-depth studies of the theory and methods for survey data analysis.

Balancing theory and application in any text is no simple matter. The distinguished statistician D. R. Cox begins the outline of his view of applied statistical work by stating, “Any simple recommendation along the lines *in applications one should do so and so* is virtually bound to be wrong in some or, indeed, possibly many contexts. On the other hand, descent into yawning



### THEORY BOX P.1 AN EXAMPLE THEORY BOX

Theory boxes are used in this volume to develop or explain a fundamental theoretical concept underlying statistical methods. The content of these “gray-shaded” boxes is intended to stand alone, supplementing the interested reader’s knowledge, but not necessary for understanding the general discussion of applied statistical approaches to the analysis of survey data.

vacuous generalities is all too possible” (Cox, 2007). Since the ingredients of each applied survey data analysis problem vary—the aims, the sampling design, the available survey variables—there is no single set of recipes that each analyst can simply follow without additional thought and evaluation on his or her part. On the other hand, a text on applied methods should not leave survey analysts alone, fending for themselves, with only abstract theoretical explanations to guide their way through an applied statistical analysis of survey data.

On balance, the discussion in this book will tilt toward proven recipes where theory and practice have demonstrated the value of a specific approach. In cases where theoretical guidance is less clear, we identify the uncertainty but still aim to provide advice and recommendations based on experience and current thinking on best practices.

The chapters of this book are organized to be read in sequence, each chapter building on material covered in the preceding chapters. Chapter 1 provides important context for the remaining chapters, briefly reviewing historical developments and laying out a step-by-step process for approaching a survey analysis problem. Chapters 2 through 4 will introduce the reader to the fundamental features of **complex sample designs** and demonstrate how design characteristics such as stratification, clustering, and weighting are easily incorporated into the statistical methods and software for survey estimation and inference. Treatment of statistical methods for survey data analysis begins in Chapters 5 and 6 with coverage of univariate (i.e., single-variable) descriptive and simple bivariate (i.e., two-variable) analyses of continuous and categorical variables. Chapter 7 presents the linear regression model for continuous dependent variables. Generalized linear regression modeling methods for survey data are treated in Chapters 8 and 9. Chapter 10 pertains to methods for event-history analysis of survey data, including models such as the Cox proportional hazards model and discrete time models. Chapter 11 introduces methods for handling missing data problems in survey data sets. Finally, the coverage of statistical methods for survey data analysis concludes in Chapter 12 with a discussion of new developments in the area of survey applications of advanced statistical techniques, such as multilevel analysis.

To avoid repetition in the coverage of more general topics such as the recommended steps in a regression analysis or testing hypotheses concerning regression parameters, topics will be introduced as they become relevant to the specific discussion. For example, the iterative series of steps that we recommend analysts follow in regression modeling of survey data is introduced in Chapter 7 (linear regression models for continuous outcomes), but the series applies equally to model specification, estimation, evaluation, and inference for generalized linear regression models (Chapters 8 and 9). By the same token, specific details of the appropriate procedures for each step (e.g., regression model diagnostics) are covered in the chapter on a specific technique. Readers who use this book primarily as a reference volume will find cross-references to earlier chapters useful in locating important background for discussion of specific analysis topics.

There are many quality software choices out there for survey data analysts. We selected Stata® for all book examples due to its ease of use and flexibility for survey data analysis, but examples have been replicated to the greatest extent possible using the SAS®, SPSS®, IVEware, SUDAAN®, R, WesVar®, and Mplus software packages on the book Web site (<http://www.isr.umich.edu/src/smp/asda/>). Appendix A reviews software procedures that are currently available for the analysis of complex sample survey data in these other major software systems.

Examples based on the analysis of major survey data sets are routinely used in this book to demonstrate statistical methods and software applications. To ensure diversity in sample design and substantive content, example exercises and illustrations are drawn from three major U.S. survey data sets: the 2005–2006 National Health and Nutrition Examination Survey (NHANES); the 2006 Health and Retirement Study (HRS); and the National Comorbidity Survey-Replication (NCS-R). A description of each of these survey data sets is provided in Section 1.3. A series of practical exercises based on these three data sets are included at the end of each chapter on an analysis topic to provide readers and students with examples enabling practice with using statistical software for applied survey data analysis.

Clear and consistent use of statistical notation is important. Table P.1 provides a summary of the general notational conventions used in this book. Special notation and symbol representation will be defined as needed for discussion of specific topics.

The materials and examples presented in the chapters of this book (which we refer to in subsequent chapters as ASDA) are supplemented through a companion Web site (<http://www.isr.umich.edu/src/smp/asda/>). This Web site provides survey analysts and instructors with additional resources in the following areas: links to new publications and an updated bibliography for the survey analysis topics covered in Chapters 5–12; links to sites for example survey data sets; replication of the command setups and output for the analysis examples in the SAS, SUDAAN, R, SPSS, and Mplus software systems; answers to frequently asked questions (FAQs); short technical

TABLE P.1  
Notational Conventions for Applied Survey Data Analysis

Notation	Properties	Explanation of Usage
Indices and Limits		
$N, n$	Standard usage	Population size, sample size
$M, m$	Standard usage	Subpopulation size, subpopulation sample size
$h$	Subscript	Stratum index (e.g., $\bar{y}_h$ )
$\alpha$	Subscript	Cluster or primary stage unit (PSU) index (e.g., $\bar{y}_{ha}$ )
$i$	Subscript	Element (respondent) index (e.g., $y_{ia}$ )
$j, k, l$	Subscripts	Used to index vector or matrix elements (e.g., $\beta_j$ )
Survey Variables and Variable Values		
$y, x$	Roman, lowercase, italicized, end of alphabet	Survey variables (e.g., systolic blood pressure, mmHg; weight, kg)
$Y_i, X_i$	Roman, uppercase, end of alphabet, subscript	True population values of $y, x$ for individual $i$ , with $i = 1, \dots, N$ comprising the population
$y_i, x_i$	Roman, lowercase, end of alphabet, subscript	Sample survey observation for individual $i$ (e.g., $y_i = 124.5$ mmHg, $x_i = 80.2$ kg)
$y, x, Y, X$	As above, <b>bold</b>	Vectors (or matrices) of variables or variable values (e.g., $y = \{y_1, y_2, \dots, y_n\}$ )
Model Parameters and Estimates		
$\beta_j, \gamma_j$	Greek, lowercase	Regression model parameters, subscripts
$\hat{\beta}_j, \hat{\gamma}_j$	Greek, lowercase, “^” hat	Estimates of regression model parameters
$\boldsymbol{\beta}, \boldsymbol{\gamma}, \hat{\boldsymbol{\beta}}, \hat{\boldsymbol{\gamma}}$	As above, <b>bold</b>	Vectors (or matrices) of parameters or estimates (e.g., $\boldsymbol{\beta} = \{\beta_0, \beta_1, \dots, \beta_p\}$ )
$B_j, b_j, \boldsymbol{B}, \boldsymbol{b}$	Roman, otherwise as above	As above but used to distinguish finite population regression coefficients from probability model parameters and estimates
Statistics and Estimates		
$\bar{Y}, P, \sigma_y^2, S_y^2, \bar{y}, p, s_y^2$	Standard usage	Population mean, proportion and variance; sample estimates as used in Cochran (1977)
$\boldsymbol{\Sigma}, \hat{\boldsymbol{\Sigma}}$	Standard usage	Variance–covariance matrix; sample estimate of variance–covariance matrix
$R^2, r, \psi$	Standard usage	Multiple-coefficient of determination ( $R$ -squared), Pearson product moment correlation, odds ratio
$\rho_y$	Greek, lowercase	Intraclass correlation for variable $y$
$Z, t, \chi^2, F$	Standard usage	Probability distributions

reports related to special topics in applied survey data analysis; and reviews of statistical software system updates and any resulting changes to the software commands or output for the analysis examples.

In closing, we must certainly acknowledge the many individuals who contributed directly or indirectly in the production of this book. Gail Arnold provided invaluable technical and organizational assistance throughout the production and review of the manuscript. Rod Perkins provided exceptional support in the final stages of manuscript review and preparation. Deborah Kloska and Lingling Zhang generously gave of their time and statistical expertise to systematically review each chapter as it was prepared. Joe Kazemi and two anonymous reviewers offered helpful comments on earlier versions of the introductory chapters, and SunWoong Kim and Azam Khan also reviewed the more technical material in our chapters for accuracy. We owe a debt to our many students in the JPSM and MPSM programs who over the years have studied with us—we only hope that you learned as much from us as we did from working with you. As lifelong students ourselves, we owe a debt to our mentors and colleagues who over the years have instilled in us a passion for statistical teaching and consultation: Leslie Kish, Irene Hess, Graham Kalton, Morton Brown, Edward Rothman, and Rod Little. Finally, we wish to thank the support staff at Chapman Hall/CRC Press, especially Rob Calver and Sarah Morris, for their continued guidance.

**Steven G. Heeringa**  
**Brady T. West**  
**Patricia A. Berglund**  
*Ann Arbor, Michigan*



# 1

---

## *Applied Survey Data Analysis: Overview*

---

---

### 1.1 Introduction

Modern society has adopted the **survey method** as a principal tool for looking at itself—“a telescope on society” in the words of House et al. (2004). The most common application takes the form of the periodic media surveys that measure population attitudes and beliefs on current social and political issues:

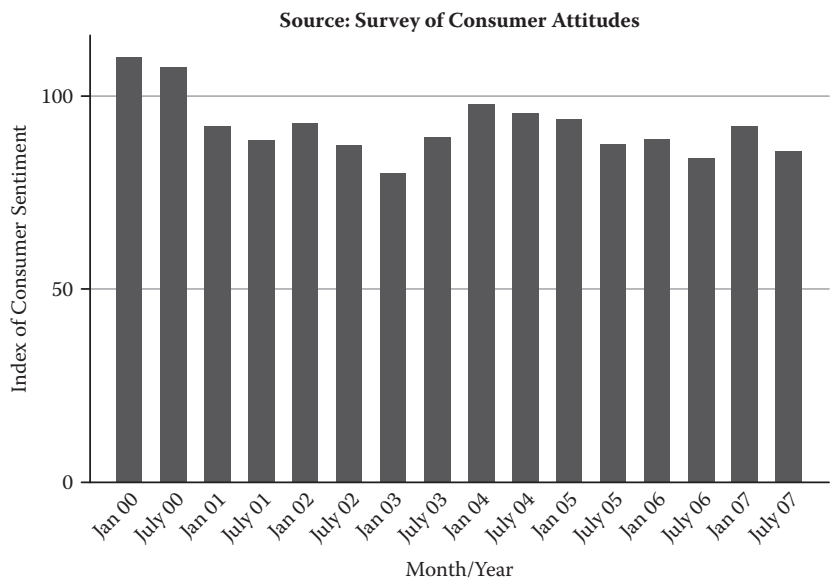
Recent international reports have said with near certainty that human activities are the main cause of global warming since 1950. The poll found that 84 percent of Americans see human activity as at least contributing to warming. (*New York Times*, April 27, 2007).

One step removed from the media limelight is the use of the survey method in the realms of marketing and consumer research to measure the preferences, needs, expectations, and experiences of consumers and to translate these to indices and other statistics that may influence financial markets or determine quality, reliability, or volume ratings for products as diverse as automobiles, hotel services, or TV programming:

CBS won the overall title with an 8.8 rating/14 share in primetime, ABC finished second at 7.7/12.... (<http://www.zap2it.com>, January 11, 2008)

The Index of Consumer Sentiment (see Figure 1.1) fell to 88.4 in the March 2007 survey from 91.3 in February and 96.9 in January, but it was nearly identical with the 88.9 recorded last March. (Reuters, University of Michigan, April 2007)

Also outside the view of most of society is the use of large-scale scientific surveys to measure labor force participation, earnings and expenditures, health and health care, commodity stocks and flows, and many other topics. These larger and longer-term programs of survey research are critically important to social scientists, health professionals, policy makers, and administrators and thus indirectly to society itself.



**FIGURE 1.1**  
Index of Consumer Sentiment, January 2000–July 2007.

Real median household income in the United States rose between 2005 and 2006, for the second consecutive year. Household income increased 0.7 percent, from \$47,845 to \$48,201. (DeNavas-Walt, Proctor, and Smith, 2007)

In a series of logistic models that included age and one additional variable (i.e., education, gender, race, or APOE genotype), older age was consistently associated with increased risk of dementia ( $p < 0.0001$ ). In these trivariate models, more years of education was associated with lower risk of dementia ( $p < 0.0001$ ). There was no significant difference in dementia risk between males and females ( $p = 0.26$ ). African Americans were at greater risk for dementia ( $p = 0.008$ ). As expected, the presence of one (Odds Ratio = 2.1; 95% C.I. = 1.45 – 3.07) or two (O.R. = 7.1; 95% C.I. = 2.92 – 17.07) APOE e4 alleles was significantly associated with increased risk of dementia. (Plassman et al., 2007)

The focus of this book will be on analysis of **complex sample survey data** typically seen in large-scale scientific surveys, but the general approach to survey data analysis and specific statistical methods described here should apply to all forms of survey data.

To set the historical context for contemporary methodology, Section 1.2 briefly reviews the history of developments in theory and methods for applied survey data analysis. Section 1.3 provides some needed background on the data sets that will be used for the analysis examples in Chapters 2–12. This short overview chapter concludes in Section 1.4 with

a general review of the sequence of steps required in any applied analysis of survey data.

---

## 1.2 A Brief History of Applied Survey Data Analysis

Today's survey data analysts approach a problem armed with substantial background in statistical survey theory, a literature filled with empirical results and high-quality software tools for the task at hand. However, before turning to the best methods currently available for the analysis of survey data, it is useful to look back at how we arrived at where we are today. The brief history described here is certainly a selected interpretation, chosen to emphasize the evolution of probability sampling design and related statistical analysis techniques that are most directly relevant to the material in this book. Readers interested in a comprehensive review of the history and development of survey research in the United States should see Converse (1987). Bulmer (2001) provides a more international perspective on the history of survey research in the social sciences. For the more statistically inclined, Skinner, Holt, and Smith (1989) provide an excellent review of the development of methods for descriptive and analytical treatment of survey data. A comprehensive history of the impacts of sampling theory on survey practice can be found in O'Muircheartaigh and Wong (1981).

### 1.2.1 Key Theoretical Developments

The science of survey sampling, survey data collection methodology, and the analysis of survey data date back a little more than 100 years. By the end of the 19th century, an open and international debate established the **representative sampling method** as a statistically acceptable basis for the collection of observational data on populations (Kaier, 1895). Over the next 30 years, work by Bowley (1906), Fisher (1925), and other statisticians developed the role of randomization in sample selection and large-sample methods for estimation and statistical inference for simple random sample (SRS) designs.

The early work on the representative method and inference for simple random and stratified random samples culminated in a landmark paper by Jerzy Neyman (1934), which outlined a cohesive framework for estimation and inference based on estimated confidence intervals for population quantities that would be derived from the probability distribution for selected samples over repeated sampling. Following the publication of Neyman's paper, there was a major proliferation of new work on survey sample designs, estimation of population statistics, and variance estimation required to develop confidence intervals for sample-based inference, or what in more recent times has been labeled **design-based inference** (Cochran, 1977; Deming, 1950;



Hansen, Hurwitz, and Madow, 1953; Kish, 1965; Sukatme, 1954; Yates, 1949). House et al. (2004) credit J. Steven Stock (U.S. Department of Agriculture) and Lester Frankel (U.S. Bureau of the Census) with the first applications of area probability sampling methods for household survey data collections. Even today, the primary techniques for sample design, population estimation, and inference developed by these pioneers and published during the period 1945–1975 remain the basis for almost all descriptive analysis of survey data.

The developments of the World War II years firmly established the probability sample survey as a tool for describing population characteristics, beliefs, and attitudes. Based on Neyman's (1934) theory of inference, survey sampling pioneers in the United States, Britain, and India developed optimal methods for sample design, estimators of survey population characteristics, and confidence intervals for population statistics. As early as the late 1940s, social scientists led by sociologist Paul Lazarsfeld of Columbia University began to move beyond using survey data to simply describe populations to using these data to explore relationships among the measured variables (see Kendall and Lazarsfeld, 1950; Klein and Morgan, 1951). Skinner et al. (1989) and others before them labeled these two distinct uses of survey data as *descriptive* and *analytical*. Hyman (1955) used the term *explanatory* to describe scientific surveys whose primary purpose was the analytical investigation of relationships among variables.

During the period 1950–1990, analytical treatments of survey data expanded as new developments in statistical theory and methods were introduced, empirically tested, and refined. Important classes of methods that were introduced during this period included log-linear models and related methods for contingency tables, generalized linear models (e.g., logistic regression), survival analysis models, general linear mixed models (e.g., hierarchical linear models), structural equation models, and latent variable models. Many of these new statistical techniques applied the method of maximum likelihood to estimate model parameters and standard errors of the estimates, assuming that the survey observations were *independent* observations from a known probability distribution (e.g., binomial, multinomial, Poisson, product multinomial, normal). As discussed in Chapter 2, data collected under most contemporary survey designs do not conform to the key assumptions of these methods.

As Skinner et al. (1989) point out, survey statisticians were aware that straightforward applications of these new methods to complex sample survey data could result in underestimates of variances and therefore could result in biased estimates of confidence intervals and test statistics. However, except in limited situations of relatively simple designs, exact determination of the size and nature of the bias (or a potential correction) were difficult to express analytically. Early investigations of such “design effects” were primarily empirical studies, comparing design-adjusted variances for estimates with the variances that would be obtained if the

data were truly identically and independently distributed (equivalent to a simple random sample of equal size). Over time, survey statisticians developed special approaches to estimating these models that enabled the survey analyst to take into account the complex characteristics of the survey sample design (e.g., Binder, 1983; Kish and Frankel, 1974; Koch and Lemeshow, 1972; Pfeffermann et al., 1998; Rao and Scott, 1981). These approaches (and related developments) are described in Chapters 5–12 of this text.

### 1.2.2 Key Software Developments

Development of the underlying statistical theory and empirical testing of new methods were obviously important, but the survey data analyst needed computational tools to apply these techniques. We can have nothing but respect for the pioneers who in the 1950s fitted multivariate regression models to survey data using only hand computations (e.g., sums, sums of squares, sums of cross-products, matrix inversions) performed on a rotary calculator and possibly a tabulating machine (Klein and Morgan, 1951). The origin of statistical software as we know it today dates back to the 1960s, with the advent of the first mainframe computer systems. Software systems such as BMDP and OSIRIS and later SPSS, SAS, GLIM, S, and GAUSS were developed for mainframe users; however, with limited exceptions, these major software packages did not include programs that were adapted to complex sample survey data.

To fill this void during the 1970s and early 1980s, a number of stand-alone programs, often written in the Fortran language and distributed as compiled objects, were developed by survey statisticians (e.g., OSIRIS: PSALMS and REPERR, CLUSTERS, CARP, SUDAAN, WesVar). By today's standards, these programs had a steep "learning curve," limited data management flexibility, and typically supported only descriptive analysis (means, proportions, totals, ratios, and functions of descriptive statistics) and linear regression modeling of multivariate relationships. A review of the social science literature of this period shows that only a minority of researchers actually employed these special software programs when analyzing complex sample survey data, resorting instead to standard analysis programs with their default assumption that the data originated with a simple random sample of the survey population.

The appearance of microcomputers in the mid-1980s was quickly followed by a transition to personal computer versions of the major statistical software (BMDP, SAS, SPSS) as well as the advent of new statistical analysis software packages (e.g., SYSTAT, Stata, S-Plus). However, with the exception of specialized software systems (WesVar PC, PC CARP, PC SUDAAN, Micro-OSIRIS, CLUSTERS for PC, IVEware) that were often designed to read data sets stored in the formats of the larger commercial software packages, the microcomputing revolution still did not put tools for the analysis of complex

sample survey data in the hands of most survey data analysts. Nevertheless, throughout the late 1980s and early 1990s, the scientific and commercial pressures to incorporate programs of this type into the major software systems were building. Beginning with Version 6.12, SAS users had access to PROC SURVEYMEANS and PROC SURVEYREG, two new SAS procedures that permitted simple descriptive analysis and linear regression analysis for complex sample survey data. At about the same time, the Stata system for statistical analysis appeared on the scene, providing complex sample survey data analysts with the “svy” versions of the more important analysis programs. SPSS’s entry into the world of complex sample survey data analysis came later with the introduction of the Complex Samples add-on module in Version 13. Appendix A of this text covers the capabilities of these different systems in detail.

The survey researcher who sits down today at his or her personal computing work station has access to powerful software systems, high-speed processing, and high-density data storage capabilities that the analysts in the 1970s, 1980s, and even the 1990s could not have visualized. All of these recent advances have brought us to a point at which today’s survey analyst can approach both simple and complex problems with the confidence gained through a fundamental understanding of the theory, empirically tested methods for design-based estimation and inference, and software tools that are sophisticated, accurate, and easy to use.

Now that we have had a glimpse at our history, let’s begin our study of applied survey data analysis.

---

## 1.3 Example Data Sets and Exercises

Examples based on the analysis of major survey data sets are routinely used in this book to demonstrate statistical methods and software applications. To ensure diversity in sample design and substantive content, example exercises and illustrations are drawn from three major U.S. survey data sets.

### 1.3.1 The National Comorbidity Survey Replication (NCS-R)

The NCS-R is a 2002 study of mental illness in the U.S. household population ages 18 and over. The core content of the NCS-R is based on a lay-administered interview using the World Health Organization (WHO) CIDI (Composite International Diagnostic Interview) diagnostic tool, which is designed to measure primary mental health diagnostic symptoms, symptom severity, and use of mental health services (Kessler et al., 2004). The NCS-R was based on interviews with randomly chosen adults in an equal probability, multistage sample of households selected from the University of Michigan

National Sample master frame. The survey response rate was 70.9%. The survey was administered in two parts: a Part I core diagnostic assessment of all respondents ( $n = 9,282$ ), followed by a Part II in-depth interview with 5,692 of the 9,282 Part I respondents, including all Part I respondents who reported a lifetime mental health disorder and a probability subsample of the disorder-free respondents in the Part I screening.

The NCS-R was chosen as an example data set for the following reasons: (1) the scientific content and, in particular, its binary measures of mental health status; (2) the multistage design with primary stage stratification and clustering typical of many large-scale public-use survey data sets; and (3) the two-phase aspect of the data collection.

### **1.3.2 The Health and Retirement Study (HRS)—2006**

The Health and Retirement Study (HRS) is a longitudinal study of the American population 50 years of age and older. Beginning in 1992, the HRS has collected data every two years on a longitudinal panel of sample respondents born between the years of 1931 and 1941. Originally, the HRS was designed to follow this probability sample of age-eligible individuals and their spouses or partners as they transitioned from active working status to retirement, measuring aging-related changes in labor force participation, financial status, physical and mental health, and retirement planning. The HRS observation units are age-eligible individuals and “financial units” (couples in which at least one spouse or partner is HRS eligible). Beginning in 1993 and again in 1998 and 2004, the original HRS 1931–1941 birth cohort panel sample was augmented with probability samples of U.S. adults and spouses/partners from (1) pre-1924 (added in 1993); (2) 1924–1930 and 1942–1947 (added in 1998); and (3) 1948–1953 (added in 2004). In 2006, the HRS interviewed over 22,000 eligible sample adults in the composite panel.

The HRS samples were primarily identified through in-person screening of large, multistage area probability samples of U.S. households. For the pre-1931 birth cohorts, the core area probability sample screening was supplemented through sampling of age-eligible individuals from the U.S. Medicare Enrollment Database. Sample inclusion probabilities for HRS respondents vary slightly across birth cohorts and are approximately two times higher for African Americans and Hispanics. Data from the 2006 wave of the HRS panel are used for most of the examples in this text, and we consider a longitudinal analysis of multiple waves of HRS data in Chapter 12.

### **1.3.3 The National Health and Nutrition Examination Survey (NHANES)—2005, 2006**

Sponsored by the National Center for Health Statistics (NCHS) of the Centers for Disease Control and Prevention (CDC), the NHANES is a survey of the adult, noninstitutionalized population of the United States. The NHANES

is designed to study the prevalence of major disease in the U.S. population and to monitor the change in prevalence over time as well as trends in treatment and major disease risk factors including personal behaviors, environmental exposure, diet, and nutrition. The NHANES survey includes both an in-home medical history interview with sample respondents and a detailed medical examination at a local mobile examination center (MEC). The NHANES surveys were conducted on a periodic basis between 1971 and 1994 (NHANES I, II, III), but beginning in 1999, the study transitioned to a continuous interviewing design. Since 1999, yearly NHANES data collections have been performed in a multistage sample that includes 15 primary stage unit (PSU) locations with new sample PSUs added in each data collection year. Approximately 7,000 probability sample respondents complete the NHANES in-home interview phase each year and roughly 5,000 of these individuals also consent to the detailed MEC examination. To meet specific analysis objectives, the NHANES oversamples low-income persons, adolescents between the ages of 12 and 19, persons age 60 and older, African Americans, and Hispanics of Mexican ancestry. To ensure adequate precision for sample estimates, NCHS recommends pooling data for two or more consecutive years of NHANES data collection. The NHANES example analyses provided in this text are based on the combined data collected in 2005 and 2006. The unweighted response rate for the interview phase of the 2005–2006 NHANES was approximately 81%.

Public use versions of each of these three major survey data sets are available online. The companion Web site for this book provides the most current links to the official public use data archives for each of these example survey data sets.

### **1.3.4 Steps in Applied Survey Data Analysis**

Applied survey data analysis—both in daily practice and here in this book—is a process that requires more of the analyst than simple familiarity and proficiency with statistical software tools. It requires a deeper understanding of the sample design, the survey data, and the interpretation of the results of the statistical methods. Following a more general outline for applied statistical analysis presented by Cox (2007), Figure 1.2 outlines a sequence of six steps that are fundamental to applied survey data analysis, and we describe these steps in more detail in the following sections.

#### ***1.3.4.1 Step 1: Definition of the Problem and Statement of the Objectives***

The first of the six steps involves a clear specification of the problem to be addressed and formulation of objectives for the analysis exercise. For example, the “problem” may be ambiguity among physicians over whether there should be a lower threshold for prostate biopsy following prostate specific antigen (PSA) screening in African American men (Cooney et al., 2001). The

Step	Activity
1	Definition of the problem and statement of the objectives.
2	Understanding the sample design.
3	Understanding design variables, underlying constructs, and missing data.
4	Analyzing the data.
5	Interpreting and evaluating the results of the analysis.
6	Reporting of estimates and inferences from the survey data.

**FIGURE 1.2**

Steps in applied survey data analysis.

corresponding objective would be to estimate the 95th percentile and the 95% confidence bounds for this quantity ( $\pm .2$  ng/ml PSA) in a population of African American men. The estimated 95% confidence bounds can in turn be used by medical experts to determine if the biopsy threshold for African American men should be different than for men of other race and ethnic groups.

As previously described, the problems to which survey data analyses may be applied span many disciplines and real-world settings. Likewise, the statistical objectives may vary. Historically, the objectives of most survey data analyses were to describe characteristics of a target population: its average household income, the median blood pressure of men, or the proportion of eligible voters who favor candidate X. But survey data analyses can also be used for decision making. For example, should a pharmaceutical company recall its current products from store shelves due to a perceived threat of contamination? In a population case-control study, does the presence of silicone breast implants significantly increase the odds that a woman will contract a connective tissue disease such as scleroderma (Burns et al., 1996)? In recent decades, the objective of many sample survey data analyses has been to explore and extend the understanding of multivariate relationships among variables in the target population. Sometimes multivariate modeling of survey data is seen simply as a descriptive tool, defining the form of a functional relationship as it exists in a finite population. But it is increasingly common for researchers to use observational data from complex sample surveys to probe causality in the relationships among variables.

#### **1.3.4.2 Step 2: Understanding the Sample Design**

The survey data analyst must understand the sample design that was used to collect the data he or she is about to analyze. Without an understanding of key properties of the survey sample design, the analysis may be inefficient,



biased, or otherwise lead to incorrect inference. An experienced researcher who designs and conducts a randomized block experimental design to test the relative effectiveness of new instructional methods should not proceed to analyze the data as a simple factorial design, ignoring the blocking that was built into his or her experiment. Likewise, an economics graduate student who elects to work with the longitudinal HRS data should understand that the nationally representative sample of older adults includes stratification, clustering, and disproportionate sampling (i.e., compensatory population weighting) and that these design features may require special approaches to population estimation and inference.

At this point, we may have discouraged the reader into thinking that an in-depth knowledge of survey sample design is required to work with survey data or that he or she may need to relearn what was studied in general courses on applied statistical methods. This is not the case. Chapters 2 through 4 will introduce the reader to the fundamental features of **complex sample designs** and will demonstrate how design characteristics such as stratification, clustering, and weighting are easily incorporated into the statistical methods and software for survey estimation and inference. Chapters 5–12 will show the reader that relatively simple extensions of his or her current knowledge of applied statistical analysis methods provide the necessary foundation for efficient and accurate analysis of data collected in sample surveys.

#### ***1.3.4.3 Step 3: Understanding Design Variables, Underlying Constructs, and Missing Data***

The typical scientific survey data set is **multipurpose**, with the final data sets often including hundreds of variables that span many domains of study— income, education, health, family. The sheer volume of available data and the ease by which it can be accessed can cause survey data analysts to become complacent in their attempts to fully understand the properties of the data that are important to their choice of statistical methods and the conclusions that they will ultimately draw from their analysis. Step 2 described the importance of understanding the sample design. In the survey data, the key features of the sample design will be encoded in a series of **design variables**. Before analysis begins, some simple questions need to be put to the candidate data set: What are the empirical distributions of these design variables, and do they conform to the design characteristics outlined in the technical reports and online study documentation? Does the original survey question that generated a variable of interest truly capture the underlying construct of interest? Are the response scales and empirical distributions of responses and independent variables suitable for the intended analysis? What is the distribution of missing data across the cases and variables, and is there a potential impact on the analysis and the conclusions that will be drawn?

Chapter 4 discusses techniques for answering these and other questions before proceeding to statistical analysis of the survey data.

#### ***1.3.4.4 Step 4: Analyzing the Data***

Finally we arrive at the step to which many researchers rush to enter the process. We are all guilty of wanting to jump ahead. Identifying the problem and objectives seems intuitive. We tell ourselves that formalizing that step wastes time. Understanding the design and performing data management and exploratory analysis to better understand the data structure is boring. After all, the statistical analysis step is where we obtain the results that enable us to describe populations (through confidence intervals), to extend our understanding of relationships (through statistical modeling), and possibly even to test scientific hypotheses.

In fact, the statistical analysis step lies at the heart of the process. Analytic techniques must be carefully chosen to conform to the analysis objectives and the properties of the survey data. Specific methodology and software choices must accommodate the design features that influence estimation and inference. Treatment of statistical methods for survey data analysis begins in Chapters 5 and 6 with coverage of univariate (i.e., single-variable) descriptive and simple bivariate (i.e., two-variable) analyses of continuous and categorical variables. Chapter 7 presents the linear regression model for continuous dependent variables, and generalized linear regression modeling methods for survey data are treated in Chapters 8 and 9. Chapter 10 pertains to methods for event-history analysis of survey data, including models such as the Cox proportional hazard model and discrete time logistic models. Chapter 11 introduces methods for handling missing data problems in survey data sets. Finally, the coverage of statistical methods for survey data analysis concludes with a discussion of new developments in the area of survey applications of advanced statistical techniques, such as multilevel analysis, in Chapter 12.

#### ***1.3.4.5 Step 5: Interpreting and Evaluating the Results of the Analysis***

Knowledge of statistical methods and software tools is fundamental to success as an applied survey data analyst. However, setting up the data, running the programs, and printing the results are not sufficient to constitute a thorough treatment of the analysis problem. Likewise, scanning a column of  $p$ -values in a table of regression model output does not inform us concerning the form of the “final model” or even the pure effect of a single predictor. As described in Step 3, interpretation of the results from an analysis of survey data requires a consideration of the error properties of the data. Variability of sample estimates will be reflected in the **sampling errors** (i.e., confidence intervals, test statistics) estimated in the course of the statistical analysis. **Nonsampling errors**, including potential bias due to survey nonresponse



and item missing data, cannot be estimated from the survey data (Lessler and Kalsbeek, 1992). However, it may be possible to use ancillary data to explore the potential direction and magnitude of such errors. For example, an analyst working for a survey organization may statistically compare survey respondents with nonrespondents in terms of known correlates of key survey variables that are readily available on the sampling frame to assess the possibility of nonresponse bias.

As survey data analysts have pushed further into the realm of multivariate modeling of survey data, care is required in interpreting fitted models. Is the model reasonably identified, and do the data meet the underlying assumptions of the model estimation technique? Are there alternative models that explain the observed data equally well? Is there scientific support for the relationship implied in the modeling results? Are interpretations that imply causality in the modeled relationships supported (Rothman, 1988)?

#### ***1.3.4.6 Step 6: Reporting of Estimates and Inferences from the Survey Data***

The end products of applied survey data analyses are reports, papers, or presentations designed to communicate the findings to fellow scientists, policy analysts and administrators and decision makers. This text includes discussion of standards and proven methods for effectively presenting the results of applied survey data analyses, including table formatting, statistical contents, and the use of statistical graphics.

With these six steps in mind, we now can begin our walk through the process of planning, formulating, and conducting analysis of survey data.

---

## References

---

- Agresti, A., *Categorical Data Analysis*, John Wiley & Sons, New York, 2002.
- Allison, P.D., Discrete-time methods for the analysis of event histories, *Sociological Methodology*, 13, 61–98, 1982.
- Allison, P.D., *Survival Analysis Using the SAS System: A Practical Guide*, SAS Institute, Cary, NC, 1995.
- Allison, P.D., *Logistic Regression Using the SAS® System: Theory and Application*, Cary, NC, 1999.
- Archer, K.J. and Lemeshow, S., Goodness-of-fit test for a logistic regression model estimated using survey sample data, *Stata Journal*, 6(1), 97–105, 2006.
- Archer, K.J., Lemeshow, S., and Hosmer, D.W., Goodness-of-fit tests for logistic regression models when data are collected using a complex sample design, *Computational Statistics and Data Analysis*, 51, 4450–4464, 2007.
- Barnard, J. and Rubin, D.B., Small-sample degrees of freedom with multiple imputation, *Biometrika*, 86(4), 948–955, 1999.
- Belli, R.F., Computerized event history calendar methods: Facilitating autobiographical recall, *American Statistical Association, Proceedings of the Section on Survey Research Methods*, 471–475, 2000.
- Biemer, P.P., Groves, R.M., Lyberg, L.E., Mathiowetz, N.A., and Sudman, S. (Eds.), *Measurement Errors in Surveys*, John Wiley & Sons, New York, 1991.
- Binder, D.A., On the variances of asymptotically normal estimators from complex surveys, *Survey Methodology*, 7, 157–170, 1981.
- Binder, D.A., On the variances of asymptotically normal estimators from complex surveys, *International Statistical Review*, 51, 279–292, 1983.
- Binder, D.A., Use of estimating functions for interval estimation from complex surveys, presented at *International Statistical Institute Meetings in Cairo*, 1991.
- Binder, D.A., Fitting Cox's proportional hazards model from survey data, *Biometrika*, 79, 139–147, 1992.
- Binder, D.A., Longitudinal surveys: Why are these surveys different from all other surveys? *Survey Methodology*, 24(2), 101–108, 1998.
- Bishop, Y.M., Feinberg, S.E., and Holland, P.W., *Discrete Multivariate Analysis*, MIT Press, Cambridge, MA, 1975.
- Bollen, K.A., *Structural Equations with Latent Variables*, Wiley-Interscience, New York, 1989.
- Bowley, A.L., Address to the Economic Science and Statistics Section of the British Association for the Advancement of Science, *Journal of the Royal Statistical Society*, 69, 548–557, 1906.
- Breidt, F.J. and Opsomer, J.D., Nonparametric and semiparametric estimation in complex surveys, in C.R. Rao and D. Pfeffermann (Eds.), *Sample Surveys: Theory, Methods and Inference, Handbook of Statistics*, Vol. 29, Elsevier, North Holland, 2009.
- Brewer, K.R.W. and Mellor, R.W., The effects of sample structure on analytic surveys, *Australian Journal of Statistics*, 15, 145–152, 1973.

- Bulmer, M., History of social survey, in N.J. Smeltser and P.B. Baltes, *International Encyclopedia of the Social and Behavioral Sciences*, vol. 21, 14469–14473, Elsevier, Oxford, 2001.
- Burns, C.J., Laing, T.J., Gillespie, B.W., Heeringa, S.G., Alcser, K.H., Mayes, M.D., et al., The epidemiology of scleroderma among women: Assessment of risk from exposure to silicone and silica, *Journal of Rheumatology*, 23(11), 1904–1912, 1996.
- Buskirk, T. and Lohr, S., Asymptotic properties of kernel density estimation with complex survey data, *Journal of Statistical Planning and Inference*, 128, 165–190, 2005.
- Cameron, A. and Trivedi, P., *Regression Analysis of Count Data*, Cambridge University Press, Cambridge, 1998.
- Carlin, J.B., Galati, J.C., and Royston, P., A new framework for managing and analyzing multiply imputed data in Stata, *Stata Journal*, 8(1), 49–67, 2008.
- Chambers, R.L., Dorfman, A.H., and Sverchkov, M. Yu., Nonparametric regression with complex sample survey data, in R.L. Chambers and C.J. Skinner, (Eds.), *Analysis of Survey Data*, John Wiley and Sons, London, 2003.
- Chambers, R.L. and Skinner, C.J. (Eds.), *Analysis of Survey Data*, John Wiley & Sons, New York, 2003.
- Cleveland, W.S., *Visualizing Data*, Hobart Press, Summit, NJ, 1993.
- Cleves, M., Gould, W.W., Gutierrez, R.G., and Marchenko, Y., *An Introduction to Survival Analysis using Stata*, 2d ed., Stata Press, College Station, TX, 2008.
- Cochran, W.G., *Sampling Techniques*, 3d ed., John Wiley & Sons, New York, 1977.
- Converse, J.M., *Survey Research in the United States: Roots and Emergence*, University of California Press, Berkeley, 1987.
- Cooney, K.A., Strawderman, M.S., Wojno, K.J., Doerr, K.M., Taylor, A., Alcser, K.H., et al., Age-specific distribution of serum prostate-specific antigen in a community-based study of African-American men, *Urology*, 57, 91–96, 2001.
- Cox, D.R., Regression models and life tables, *Journal of the Royal Statistical Society-B*, 34, 187–220, 1972.
- Cox, D.R., Applied Statistics: A Review, *Annals of Applied Statistics*, 1(1), 1–16, 2007, 1.
- Cox, D.R. and Snell, E.J., *The Analysis of Binary Data*, 2d ed., Chapman and Hall, London, 1989.
- DeMaris, A., *Regression with Social Data*, John Wiley & Sons, New York, 2004.
- Deming, W.E., *Some Theory of Sampling*, John Wiley & Sons, New York, 1950.
- DeNavas-Walt, C., Proctor, B.D., and Smith, J., Current population reports, P60-233, *Income, Poverty and Health Insurance Coverage in the United States: 2006*, U.S. Government Printing Office, Washington, DC, 2007.
- Deville, J.-C. and Särndal, C.-E., Calibration estimators in survey sampling, *Journal of the American Statistical Association*, 87, 376–382, 1992.
- Diggle, P.J., Heagerty, P., Liang, K.-Y., and Zeger, S.L., *Analysis of Longitudinal Data*, 2d ed., Clarendon Press, Oxford, 2002.
- Draper, N.R. and Smith, H., *Applied Regression Analysis*, 2d ed., John Wiley & Sons, New York, 1981.
- DuMouchel, W.H. and Duncan, G.S., Using sample survey weights in multiple regression analyses of stratified samples, *Journal of the American Statistical Association*, 78, 535–543, 1983.
- Elliott, M.R., Bayesian weight trimming for generalized linear regression models, *Survey Methodology*, 33(1), 23–34, 2007.
- Elliott, M.R. and Little, R.J.A., Model-based approaches to weight trimming, *Journal of Official Statistics*, 16, 191–210, 2000.

- Ezzatti-Rice, T.M., Khare, M., Rubin, D.B., Little, R.J.A., and Schafer, J.L., A comparison of imputation techniques in the Third National Health and Nutrition Examination Survey, *Proceedings of the American Statistical Association, Survey Research Methods Section*, 303–308, 1993.
- Faraway, J.J., *Linear Models with R*, Chapman & Hall, CRC, London, 2005.
- Faraway, J.J., *Extending the Linear Model with R: Generalized Linear, Mixed Effects and Nonparametric Regression Models*, Chapman & Hall/CRC, New York, 2006.
- Fellegi, I.P., Approximate tests of independence and goodness of fit based on stratified multistage samples, *Journal of the American Statistical Association*, 75, 261–268, 1980.
- Fisher, R.A., *Statistical Methods for Research Workers*, Oliver and Boyd, Edinburgh, 1925.
- Fitzmaurice, G.M., Davidian, M., Verbeke, G., and Molenberghs, G. (Eds.), *Longitudinal Data Analysis*, John Wiley & Sons, Hoboken, NJ, 2009.
- Fitzmaurice, G.M., Laird, N.M., and Ware, J.H., *Applied Longitudinal Analysis*, John Wiley & Sons, Hoboken, NJ, 2004.
- Fox, J., *Applied Regression Analysis and Generalized Linear Model*, 2d ed., Sage, Thousand Oaks, CA, 2008.
- Freedman, D.A., On the so-called “Huber Sandwich Estimator” and “robust standard errors,” *American Statistician*, 60(4), 299–302, 2006.
- Freedman, D.A., Survival analysis: A primer, *American Statistician*, 62, 110–119, 2008.
- Fuller, W.A., Regression analysis for sample survey, *Sankhya, Series C*, 37, 117–132, 1975.
- Fuller, W.A., *Measurement Error Models*, John Wiley & Sons, New York, 1987.
- Fuller, W.A., Regression estimation for survey samples (with discussion), *Survey Methodology*, 28(1), 5–23, 2002.
- Fuller, W.A., Kennedy, W., Schnell, D., Sullivan, G., and Park, H.J., *PC CARP*, Iowa State University, Statistical Laboratory, Ames, 1989.
- Gelfand, A.E., Hills, S.E., Racine-Poon, A., and Smith, A.F.M., Illustration of Bayesian inference in normal data models using Gibbs sampling, *Journal of the American Statistical Association*, 85, 972–985, 1990.
- Gelfand, A.E. and Smith, A.F.M., Sampling-based approaches to calculating marginal densities, *Journal of the American Statistical Association*, 85, 398–409, 1990.
- Gelman, A., Struggles with survey weighting and regression modeling, *Statistical Science*, 22(2), 153–164, 2007.
- Gelman, A., Carlin, J.B., Stern, H.S., and Rubin, D.B., *Bayesian Data Analysis*, 2d ed., Chapman & Hall / CRC Press, Boca Raton, FL, 2003.
- Gelman, A. and Hill, J., *Data Analysis Using Regression and Multilevel / Hierarchical Models*, Cambridge University Press, New York, 2006.
- Goldstein, H., *Multilevel Statistical Models*, 3d ed., Arnold, London, 2003.
- Greenwood, M., The “error of sampling” of the Survivorship Tables. Reports on public health and medical subjects, No. 33, Appendix 1, H.M. Stationery Office, London, 1926.
- Grizzle, J., Starmer, F., and Koch, G., Analysis of categorical data by linear models, *Biometrics*, 25, 489–504, 1969.
- Groves, R.M., *Survey Errors and Survey Costs*, 2d ed., John Wiley & Sons, New York, 2004.
- Groves, R.M. and Couper, M., *Nonresponse in Household Interview Surveys*, John Wiley & Sons, New York, 1998.

- Groves, R.M. and Heeringa, S.G., Responsive design for household surveys: Tools for actively controlling survey errors and costs, *Journal of the Royal Statistical Society Series A: Statistics in Society*, 169(3), 439–457, 2006.
- Groves, R.M., Fowler, F.J., Couper, M.P., Lepkowski, J.M., Singer, E., and Tourangeau, R., *Survey Methodology*, John Wiley & Sons, New York, 2004.
- Hansen, M.H., Hurwitz, W.N., and Madow, W.G., *Sample Survey Methods and Theory, Volumes I and II*, John Wiley & Sons, New York, 1953.
- Hansen, M.H., Madow, W.G., and Tepping, B.J., An evaluation of model-dependent and probability-sampling inferences in sample surveys, *Journal of the American Statistical Association*, 78, 776–793, 1983.
- Harms, T. and Duchesne, P., On kernel nonparametric regression designed for complex survey data, *Metrika*, published online March 12, 2009 at <http://www.springerlink.com/content/b61n11736222pn4/fulltext.pdf>.
- Harrell, F.E. Jr., *Regression Modeling Strategies, with Applications to Linear Models, Logistic Regression, and Survival Analysis*, Springer-Verlag, New York, 2001.
- Heeringa, S. and O'Muircheartaigh, C., Sample design for cross-national, cross-cultural survey programs, in J. Harkness, M. Braun, B. Edwards, T. Johnson, L. Lyberg, P. Mohler, et al. (Eds.), *Survey Methods in Multinational, Multiregional, and Multicultural Contexts*, John Wiley & Sons, Hoboken, NJ (in press).
- Heeringa, S. and O'Muircheartaigh, C., *Survey Methods in Multinational, Multiregional, and Multicultural Contexts*, 247–263.
- Heeringa, S.G., Alcser, K.H., Doerr, K., Strawderman, M., Cooney, K., Medberry, B., et al., Potential selection bias in a community-based study of PSA Levels in African-American men, *Journal of Clinical Epidemiology*, 54(2), 142–148, 2001.
- Heeringa, S.G. and Connor, J., 1980 SRC National Sample: Design and Development, Technical report, Survey Research Center, University of Michigan, Ann Arbor, 1986.
- Heeringa, S.G. and Connor, J., Technical Description of the Health and Retirement Survey Sample Design, Technical report, Survey Research Center, University of Michigan, Ann Arbor, 1995, accessed June 2009 at <http://hrsonline.isr.umich.edu/sitedocs/userg/HRSSAMP.pdf>.
- Heeringa, S., Little, R.J.A., and Raghunathan, T., Multivariate imputation of coarsened survey data on household wealth, in R.M. Groves, D.A. Dillman, J.L. Eltinge, and R.J.A. Little (Eds.), *Survey Nonresponse*, John Wiley & Sons, New York, 2002.
- Heeringa, S., Wagner, J., Torres, M., Duan, N., Adams, T. and Berglund, P., Sample designs and sampling methods for the Collaborative Psychiatric Epidemiology Studies (CPES), *International Journal of Methods in Psychiatric Research*, 13(4), 221–239, 2004.
- Herzog, T. and Rubin, D.B., Using multiple imputations to handle nonresponse in sample surveys, in W.G. Madow, I. Olkin, and D.B. Rubin (Eds.), *Incomplete Data in Sample Surveys, Volume 2: Theory and Bibliography*, Academic Press, New York, 1983.
- Hilbe, J.M., *Negative Binomial Regression*, Cambridge University Press, Cambridge, 2007.
- Hill, M.S., *The Panel Study of Income Dynamics: A User's Guide*, Sage, Beverly Hills, CA, 1992.
- Hoerl, A.E. and Kennard, R.W., Ridge regression: Biased estimation for nonorthogonal problems, *Technometrics*, 12, 55–67, 1970.

- Holt, D. and Smith, T.M.F., Post stratification, *Journal of the Royal Statistical Society, Series A (General)*, 142(1), 33–46, 1979.
- Horvitz, D.G. and Thompson, D.J., A generalization of sampling without replacement from a finite universe, *Journal of the American Statistical Association*, 47, 663–685, 1952.
- Hosmer, D.W. and Lemeshow, S., *Applied Logistic Regression*, Wiley, New York, 1989.
- Hosmer, D.W. and Lemeshow, S., *Applied Logistic Regression*, 2d ed., John Wiley & Sons, New York, 2000.
- Hosmer, D.W., Lemeshow, S., and May, S., *Applied Survival Analysis: Regression Modeling of Time to Event Data*, 2d ed., John Wiley & Sons, Hoboken, NJ, 2008.
- House, J.S., Juster, F.T., Kahn, R.L., Schuman, H., and Singer, E., *A Telescope on Society: Survey Research and Social Science at the University of Michigan and Beyond*, University of Michigan Press, Ann Arbor, 2004.
- Rao, J.N.K. and Rust, K.F., Variance estimation for complex surveys using replication techniques, *Statistical Methods in Medical Research*, 5, 283–310, 1996.
- Hyman, H.H., *Survey Design and Analysis*, Free Press, New York, 1955.
- Jann, B., Multinomial goodness of fit: Large-sample tests with survey design correction and exact tests for small samples, *Stata Journal*, 8(2), 147–169, 2008.
- Jans, M., Heeringa, S.G., and Charest, A.-S., Imputation for missing physiological and health measurement data: Tests and applications, *Proceedings of the Survey Research Methods Section of the American Statistical Association*, 2450–2457, 2008.
- Judge, G.G., Griffiths, W.E., Hill, R.C., and Lee, T.-C., *The Theory and Practice of Econometrics*, 2d ed., John Wiley & Sons, New York, 1985.
- Judkins, D.R., Fay's method for variance estimation, *Journal of Official Statistics*, 6, 223–239, 1990.
- Juster, F.T. and Suzman, R., The Health and Retirement Study: An overview, *Journal of Human Resources*, 1995(30 Suppl.), S7–S6, 1995.
- Kaier, A.N., Observations et expériences concernant des denombrements représentatives. Discussion appears in Liv. 1, XCIII–XCVII, *Bulletin of the International Statistical Institute*, 9, Liv. 2, 176–183, 1895.
- Kalbfleisch, J.D. and Prentice, R.L., *The Statistical Analysis of Failure Time Data*, Wiley, New York, 1980.
- Kalbfleisch, J.D. and Prentice, R.L., *The Statistical Analysis of Failure Time Data*, 2d ed., John Wiley & Sons, New York, 2002.
- Kalton, G., *Introduction to Survey Sampling*, Sage, Beverly Hills, CA, 1983.
- Kalton, G., Handling wave nonresponse in panel surveys, *Journal of Official Statistics*, 2(3), 303–314, 1986.
- Kalton, G. and Citro, C., Panel surveys: Adding the fourth dimension, *Survey Methodology*, 19, 205–215, 1993.
- Kalton, G. and Kasprzyk, D., The treatment of missing survey data, *Survey Methodology*, 12(1), 1–16, 1986.
- Kavoussi, S.K., West, B.T., Taylor, G.W., and Lebovic, D.I., Periodontal disease and endometriosis: Analysis of the National Health and Nutrition Examination Survey, *Fertility & Sterility*, 91(2), 335–342, 2009.
- Keeter, S., Miller, C., Kohut, A., Groves, R.M., and Presser, S., Consequences of reducing nonresponse in a national telephone survey, *Public Opinion Quarterly*, 64, 125–148, 2000.



- Kendall, P.L. and Lazarsfeld, P.F., Problems of survey analysis, in R.K. Merton and P.F. Lazarsfeld (Eds.), *Continuities in Social Research: Studies in the Scope and Method of "The American Soldier,"* Free Press, Chicago, 1950.
- Kennickell, A.B., Multiple imputation in the Survey of Consumer Finances, Federal Reserve Board, Paper 78, Washington, DC, September 1998.
- Kessler, R.C., Berglund, P., Chiu, W.T., Demler, O., Heeringa, S., Hiripi, E., et al., The US National Comorbidity Survey Replication (NCS-R): Design and field procedures, *International Journal of Methods in Psychiatric Research*, 13(2), 69–92, 2004.
- Kish, L., A procedure for objective respondent selection within the household, *Journal of the American Statistical Association*, 44, 380–387, 1949.
- Kish, L., *Survey Sampling*, John Wiley & Sons, New York, 1965.
- Kish, L., *Statistical Design for Research*, New York: John Wiley & Sons, 1987.
- Kish, L. and Frankel, M.R., Inference from complex samples, *Journal of the Royal Statistical Society, Series B*, 36, 1–37, 1974.
- Kish, L. and Hess, I., On variances of ratios and their differences in multi-stage samples, *Journal of the American Statistical Association*, 54, 416–446, 1959.
- Kish, L., Groves, R.M., and Krotki, K., Sampling errors for fertility surveys, *Occasional Papers*, No. 17, World Fertility Survey, 1976.
- Klein, L.R. and Morgan, J.N., Results of alternative statistical treatment of sample survey data, *Journal of the American Statistical Association*, 46, 442–460, 1951.
- Kleinbaum D., Kupper L., and Muller K., *Applied Regression Analysis and Other Multivariable Methods*, 2d ed., Duxbury Press, Belmont, CA, 1988.
- Kline, R.B., *Principles and Practice of Structural Equation Modeling*, 2d ed., Guilford Press, New York, 2004.
- Koch, G.G. and Lemeshow, S., An application of multivariate analysis to complex sample survey data, *Journal of the American Statistical Association*, 54, 59–78, 1972.
- Kolenikov, S., Resampling variance estimation for complex survey data, *Stata Journal* (in press).
- Korn, E.L. and Graubard, B.I., Simultaneous testing of regression coefficients with complex survey data: Use of Bonferroni t statistics, *American Statistician*, 44, 270–276, 1990.
- Korn, E.L. and Graubard, B.I., Scatterplots with survey data, *American Statistician*, 52(1), 58–69, 1998.
- Korn, E.L. and Graubard, B.I., *Analysis of Health Surveys*, John Wiley & Sons, New York, 1999.
- Kott, P.S., A model-based look at linear regression with survey data, *American Statistician*, 45, 107–112, 1991.
- Kott, P.S. and Carr, D.A., Developing an estimation strategy for a pesticide data program, *Journal of Official Statistics*, 13(4), 367–383, 1997.
- Kovar, J.G., Rao, J.N.K., and Wu, C.F.J., Bootstrap and other methods to measure errors in survey estimates, *Canadian Journal of Statistics*, 16 Suppl., 25–45, 1988.
- Landis, R.J., Stanish, W.M., Freeman, J.L., and Koch, G.G., A computer program for the generalized chi-square analysis of categorical data using weighted least squares (GENCAT), *Computer Programs in Biomedicine*, 6, 196–231, 1976.
- Lee, E.S. and Forthofer, R.N., *Analyzing Complex Survey Data*, 2d ed., Sage, Thousand Oaks, CA, 2006.
- Lee, E.T., *Statistical Methods for Survival Analysis*, John Wiley & Sons, New York, 1992.

- Lepkowski, J.M. and Couper, M.P., Nonresponse in the second wave of longitudinal household surveys, pp. 259–271 in R.M. Groves, D.A. Dillman, J.L. Eltinge, and R.J.A. Little (Eds.), *Survey Nonresponse*, John Wiley & Sons, New York, 2002.
- Lessler, J.T. and Kalsbeek, W.D., *Nonsampling Errors in Surveys*, John Wiley & Sons, New York, 1992.
- Levy, P.S. and Lemeshow, S., *Sampling of Populations: Methods and Applications*, 4th ed., John Wiley & Sons, New York, 2007.
- Li, J., Linear regression diagnostics in cluster samples, *Proceedings of the Survey Research Methods Section of the American Statistical Association*, Joint Statistical Meetings, 2007.
- Li, J. and Valliant, R., Influence analysis in linear regression with sampling weights, *Proceedings of the Section on Survey Methods Research, American Statistical Association*, 3330, 2006.
- Li, J. and Valliant, R., Survey weighted hat matrix and leverages, *Survey Methodology*, 35(1), 15–24, 2009.
- Li, K.H., Raghunathan, T.E., and Rubin, D.B., Large sample significance levels from multiply-imputed data using moment-based statistics and an F reference distribution, *Journal of the American Statistical Association*, 86, 1065–1073, 1991.
- Little, R.J., The Bayesian approach to sample survey inference, chapter 4 in R. Chambers and C.J. Skinner (Eds.), *Analysis of Survey Data*, John Wiley & Sons, Hoboken, NJ, 2003.
- Little, R.J.A., Inference with survey weights, *Journal of Official Statistics*, 7, 405–424, 1991.
- Little, R.J.A. and Rubin, D.B., *Statistical Analysis with Missing Data*, 2nd ed., John Wiley & Sons, New York, 2002.
- Little, R.J. and Vartivarian, S., Does weighting for nonresponse increase the variance of survey means? *Survey Methodology*, 31(2), 161–168, 2005.
- Lohr, S.L., *Sampling: Design and Analysis*, Duxbury Press, Pacific Grove, CA, 1999.
- Long, J.S. and Freese, J., *Regression Models for Categorical Dependent Variables Using Stata*, 2nd ed., Stat Press, College Station, Texas, 2006.
- Loomis, D., Richardson, D.B., and Elliott, L., Poisson regression analysis of ungrouped data, *Occupational and Environmental Medicine* 62, 325–329, 2005.
- Lumley, T., R software from the R Project, <http://www.r-project.org/>, V2.7 Analysis of complex survey samples, maintained by Thomas Lumley, University of Washington, 2005.
- Lynn, P. (Ed.), *Methodology of Longitudinal Surveys*, John Wiley & Sons, New York, 2009.
- Madow, W.G. and Olkin, I. (Eds.), *Incomplete Data in Sample Surveys, Volume 3: Proceedings of the Symposium*, Academic Press, New York, 1983.
- Mahalanobis, P.C., Recent experiments in statistical sampling in the Indian Statistical Institute, *Journal of the Royal Statistical Society*, 109, 325–370, 1946.
- Maindonald, J.H. and Braun, W.J., *Data analysis and graphics using R: an example-based approach*, 2d ed., Cambridge University Press, New York, 2007.
- McCabe, S.E., West, B.T., Morales, M., Cranford, J.A., and Boyd, C.J., Does early onset of non-medical use of prescription drugs predict subsequent prescription drug abuse and dependence? Results from a national study, *Addiction*, 102(12), 1920–1930, 2007.
- McCarthy, P.J., Pseudoreplication: Half samples, *Review of the International Statistical Institute*, 37, 239–264, 1969.



- McCulloch, C.E. and Searle, S.R., *Generalized, Linear and Mixed Models*, John Wiley & Sons, New York, 2001.
- McCullagh, P. and Nelder, J.A., *Generalized Linear Models*, 2d ed., Chapman and Hall, London, 1989.
- McFadden, D., Conditional logit analysis of qualitative choice behavior, in P. Zarembka (Ed.), *Frontiers in Economics*, Academic Press, New York, 1974.
- Menard, S.W. (Ed.), *Handbook of Longitudinal Research*, Academic Press, New York, 2008.
- Miller, R., *Survival Analysis*, John Wiley & Sons, New York, 1981.
- Miller, R.G., The jackknife—a review, *Biometrika*, 61, 1–15, 1974.
- Mitchell, M.N., *A Visual Guide to Stata Graphics*, 2d ed., Stata Press, College Station, TX, 2008.
- Molenberghs, G. and Verbeke, G., *Models for Discrete Longitudinal Data*, Springer, New York, 2005.
- Mohadjer, L. and Curtin, L.R., NHANES, Balancing sample design goals for the National Health and Nutrition Examination Survey, *Survey Methodology*, 34(1), 119–126, 2008.
- Morel, G., Logistic regression under complex survey designs, *Survey Methodology*, 15, 202–223, 1989.
- Muthén, B.O. and Satorra, A., Complex sample data in structural equation modeling, *Sociological Methodology*, 25, 267–316, 1995.
- Muthén, L.K. and Muthén, B.O., *Mplus User's Guide*, 5th ed., Muthén and Muthén, Los Angeles, CA, 1998–2007.
- Nagelkerke, N.J.D., A note on the general definition of the coefficient of determination, *Biometrika*, 78(3), 691–692, 1981.
- Neter, J., Kutner, M.H., Wasserman, W., and Nachtsheim, C.J., *Applied Linear Statistical Models*, 4th ed., McGraw-Hill/Irwin, Boston, 1996.
- Neyman, J., On the two different aspects of the representative method: The method of stratified sampling and the method of purposive selection, *Journal of the Royal Statistical Society*, 97, 558–606, 1934.
- O'Muircheartaigh, C. and Wong, S.T., The impact of sampling theory on survey practices: A review, *Bulletin of the International Statistical Institute*, 465–493, 1981.
- Opsomer, J.D. and Miller, C.P., Selecting the amount of smoothing in nonparametric regression estimation for complex surveys, *Journal of Nonparametric Statistics*, 17(5), 593–611, 2005.
- Peterson, B. and Harrell, F., Partial proportional odds models for ordinal response variables, *Applied Statistics*, 39, 205–217, 1990.
- Pfeffermann, D. and Holmes, D.J., Robustness considerations in the choice of method of inference for regression analysis of survey data, *Journal of the Royal Statistical Society, Series A*, 148, 268–278, 1985.
- Pfefferman, D., Skinner, C.J., Holmes, D.J., Goldstein, H., and Rasbash, J., Weighting for unequal selection probabilities in multilevel models, *Journal of the Royal Statistical Society, Series B* 60(1), 23–40, 1998.
- Plassman, B.L., Langa, K.M., Fisher, G.G., Heeringa, S.G., Weir, D.R., Ofstedal, M.B., et al., Prevalence of dementia in the United States: The Aging, Demographics, and Memory Study, *Neuroepidemiology*, 29, 125–132, 2007.
- Potter, F., A study of procedures to identify and trim extreme sample weights, *Proceedings of the Survey Research Methods Section*, American Statistical Association, 225–230, 1990.

- Rabe-Hesketh, S., Skrondal, A., and Pickles, A., GLLAMM Manual, U.C. Berkeley Division of Biostatistics Working Paper Series, Working Paper 160, 2004.
- Rabe-Hesketh, S. and Skrondal, A., Multilevel modelling of complex survey data, *Journal of the Royal Statistical Society-A*, 169, 805–827, 2006.
- Rabe-Hesketh, S. and Skrondal, A., *Multilevel and longitudinal modeling using Stata*, 2d ed., Stata Press, College Station, TX, 2008.
- Raghunathan, T.E. and Grizzle, J.E., A split questionnaire survey design, *Journal of the American Statistical Association*, 90(429), 54–63, 1995.
- Raghunathan, T.E., Lepkowski, J.M., Van Hoewyk, J., and Solenberger, P., A multivariate technique for multiply imputing missing values using a sequence of regression models, *Survey Methodology*, 27(1), 85–95, 2001.
- Rao, J.N.K., *Small Area Estimation*, Wiley Series in Survey Methodology, John Wiley & Sons, New York, 2003.
- Rao, J.N.K. and Scott, A.J., The analysis of categorical data from complex sample surveys: Chi-squared tests for goodness of fit and independence in two-way tables, *Journal of the American Statistical Association*, 76, 221–230, 1981.
- Rao, J.N.K. and Scott, A.J., On chi-squared test for multiway contingency tables with cell proportions estimated from survey data, *Annals of Statistics*, 12, 46–60, 1984.
- Rao, J.N.K. and Shao, J., Jackknife variance estimation with survey data under hot deck imputation, *Biometrika*, 79, 811–822, 1992.
- Rao, J.N.K. and Thomas, D.R., The analysis of cross-classified categorical data from complex sample surveys, *Sociological Methodology*, 18, 213–269, 1988.
- Rao, J.N.K. and Wu, C.F.J., Inference from stratified samples: Second order analysis of three methods for nonlinear statistics, *Journal of the American Statistical Association*, 80, 620–630, 1985.
- Rao, J.N.K. and Wu, C.F.J., Resampling inference with complex survey data, *Journal of the American Statistical Association*, 83, 231–241, 1988.
- Raudenbush, S.W., Synthesizing results for NAEP trial state assessment, in D.W. Grissmer and M. Ross (Eds.), *Analytic Issues in the Assessment of Student Achievement*, National Center for Educational Statistics, Washington, DC, 2000.
- Raudenbush, S.W. and Bryk, A.S., *Hierarchical Linear Models: Applications and Data Analysis Methods*, 2d ed., Sage, Newbury Park, CA, 2002.
- Reiter, J.P., Raghunathan, T.E., and Kinney, S.K., The importance of modeling the sampling design in multiple imputation for missing data, *Survey Methodology*, 32(2), 143–149, 2006.
- Research Triangle Institute (RTI), *SUDAAN 9.0 User's Manual: Software for Statistical Analysis of Correlated Data*, RTI, Research Triangle Park, NC, 2004.
- Roberts, G., Rao, J.N.K., and Kumar, S., Logistic regression analysis of sample survey data, *Biometrika*, 74, 1–12, 1987.
- Rothman, K.J., *Causal Inference*, Epidemiology Resources, MA, 1988, out of print.
- Royston, P., Multiple imputation of missing values, *Stata Technical Journal*, 5(4), 527–536, 2005.
- Rubin, D.B., Inference and missing data, *Biometrika*, 63(3), 581–592, 1976.
- Rubin, D.B., Basic ideas of multiple imputation for nonresponse, *Survey Methodology*, 12(1), 37–47, 1986.
- Rubin, D.B., *Multiple Imputation for Nonresponse in Surveys*, John Wiley & Sons, New York, 1987.
- Rubin, D.B., Multiple imputation after 18+ years, *Journal of the American Statistical Association*, 91(434), 473–489, 1996.

- Rubin, D.B. and Schenker, N., Multiple imputation for interval estimation from simple random samples with ignorable nonresponse, *Journal of the American Statistical Association*, 81, 366–374, 1986.
- Rueters/University of Michigan Surveys of Consumers, Accessed May 1, 2008 at [http://thomsonreuters.com/products\\_services/financial/UMichigan\\_Surveys\\_of\\_Consumers](http://thomsonreuters.com/products_services/financial/UMichigan_Surveys_of_Consumers), April 2007 Report, 1.
- Rust, K., Variance estimation for complex estimators in sample surveys, *Journal of Official Statistics*, 1, 381–397, 1985.
- Rust, K. and Hsu, V., Confidence intervals for statistics for categorical variables from complex samples, *Proceedings of the 2007 Joint Statistical Meetings*, Salt Lake City, UT, 2007.
- SAS Institute, Inc., *SAS/STAT® User's Guide, Version 9*, SAS Institute, Cary, NC, 2003.
- SAS Institute Inc., *SAS/GRAPH® 9.2: Statistical Graphics Procedures Guide*, SAS Institute, Cary, NC, 2009.
- Satterthwaite, F.E., An approximate distribution of estimates of variance components, *Biometrics*, 110–114, 1946.
- Schafer, J.L., *MIX: Multiple Imputation for Mixed Continuous and Categorical Data*, software library for S-PLUS, 1996, Written in S-PLUS and Fortran-77, at <http://www.stat.psu.edu/~jls/>.
- Schafer, J.L., *Analysis of Incomplete Multivariate Data*, Chapman & Hall, London, 1997.
- Schafer, J.L., *NORM: Multiple Imputation of Incomplete Multivariate Data under a Normal Model, Version 2*, 1999, Software for Windows 95/98/NT, at <http://www.stat.psu.edu/~jls/misoftwa.html>.
- Schafer, J.L., Ezatti-Rice, T.M., Johnson, W., Khare, M., Little, R.J.A., and Rubin, D.B., The NHANES III multiple imputation project, *Proceedings of the Survey Research Methods Section*, American Statistical Association, 696–701, 1996.
- Schoenfeld, D., Residuals for the proportional hazards regression model, *Biometrika*, 239–241, 1982.
- Schumacker, R.E. and Lomax, R.G., *A Beginner's Guide to Structural Equation Modeling*, 2d ed., Lawrence Erlbaum, 2004, Hillsdale, NJ.
- Shah, B.V., Holt, M.M., and Folsom, R.F., Inference about regression models from sample survey data, *Bulletin of the International Statistical Institute*, 41(3), 43–57, 1977.
- Shao, J. and Tu, D., *The Jackknife and Bootstrap*, Springer-Verlag, New York, 1995.
- Shao, J. and Wu, C.F.J., A general theory for jackknife variance estimation, *Annals of Statistics*, 17, 1176–1197, 1989.
- Singer, J.D. and Willett, J.B., It's about time: Using discrete-time survival analysis to study duration and the timing of events, *Journal of Educational and Behavioral Statistics*, 18, 155–195, 1993.
- Skinner, C. and Vieira, M. de T., Variance estimation in the analysis of clustered longitudinal survey data, *Survey Methodology*, 33(1), 3–12, 2007.
- Skinner, C.J. and Holmes, D.J., Random effects models for longitudinal survey data, chapter 14 in R.L. Chambers and C.J. Skinner (Eds.), *Analysis of Survey Data*, John Wiley & Sons, London, 2003.
- Skinner, C.J., Holt, D., and Smith, T.M.F., *Analysis of Complex Surveys*, John Wiley & Sons, New York, 1989.

- Skrondal, A. and Rabe-Hesketh, S., *Generalized Latent Variable Modeling: Multilevel, Longitudinal, and Structural Equation Models*, Chapman & Hall / CRC Press, Boca Raton, FL, 2004.
- Sribney, W.M., Two-way contingency tables for survey or clustered data, *Stata Technical Bulletin*, 45, 33–49, 1998.
- Stapleton, L.M., Variance estimation using replication methods in structural equation modeling with complex sample data, *Structural Equation Modeling: A Multidisciplinary Journal*, 15(2), 183–210, 2008.
- STATA Corp., *Release 10, P Manual, STATA Survey Data Manual*, College Station, TX, 2008.
- Statistical Solutions, Solas 3.0, at [http://www.statsol.ie/html/solas/solas\\_home.html](http://www.statsol.ie/html/solas/solas_home.html).
- Stiller, J.G. and Dalzell, D.R., Hot-deck imputation with SAS arrays and macros for large surveys, *Proceedings of the Twenty-Third Annual AS Users Group International Conference*, 1378–1383, 1998.
- Stokes, M.E., Davis, C.S., and Koch G.G., *Categorical Data Analysis Using the SAS System, Second edition*, SAS Institute Inc., Cary, NC, 2002.
- Striegel-Moore, R.H., Franko, D.L., Thompson, D., Affenito, S., and May, A., Exploring the typology of night eating syndrome, *International Journal of Eating Disorders*, 41(5), 411–418, 2008.
- Sukatme, P.V., *Sampling Theory of Surveys, with Applications*, Iowa State College Press, Ames, 1954.
- Tanner, M. and Wong, W., The calculation of posterior distributions by data augmentation, *Journal of the American Statistical Association*, 82, 528–550, 1997.
- Therneau, T.M., Grambsch, P.M., and Fleming, T.R., Martingale-based residuals for survival models, *Biometrika*, 77(1), 147–160, 1990.
- Thomas, D.R. and Rao, J.N.K., Small-sample comparisons of level and power for simple goodness-of-fit statistics under cluster sampling, *Journal of the American Statistical Association*, 82, 630–636, 1987.
- Thomas, N., Raghunathan, T.E., Schenker, N., Katzoff, M.J., and Johnson, C.L., An evaluation of matrix sampling methods using data from the National Health and Nutrition Examination Survey, *Survey Methodology*, 32(2), 217–231, 2006.
- Thompson, S.K. and Seber, G.A.F., *Adaptive Sampling*, John Wiley & Sons, New York, 1996.
- Tufte, E.R., *The Visual Display of Information*, Graphics Press, Cheshire, CT, 1983.
- Tukey, J.W., *Exploratory Data Analysis*, Addison-Wesley, Reading, MA, 1977.
- University of Michigan, Computer Support Group, OSIRIS VI: Statistical Analysis and Data Management Software System, Survey Research Center, Institute for Social Research, 1982.
- Valliant, R., Comparisons of variance estimators in stratified random and systematic sampling, *Journal of Official Statistics*, 6(2), 115–131, 1990.
- Valliant, R., The effect of multiple weighting steps on variance estimation, *Journal of Official Statistics*, 20(1), 1–18, 2004.
- Valliant, R., Dorfman, A.H., and Royall, R.M., *Finite Population Sampling and Inference: A Prediction Approach*, John Wiley & Sons, New York, 2000.
- Van Buuren, S. and Oudshoorn, C.G.M., *Flexible multivariate imputation by MICE*, Leiden: TNO Preventie en Gezondheid, TNO/VGZ/PG 99.054, 1999.
- Verbeke, G. and Molenberghs, G., *Linear Mixed Models for Longitudinal Data*, Springer, New York, 2005.

- Vieira, M.D.T. and Skinner, C.J., Estimating models for panel survey data under complex sampling, *Journal of Official Statistics*, 24, 343–364, 2008.
- West, B.T., Berglund, P., and Heeringa, S.G., A closer examination of subpopulation analysis of complex-sample survey data, *Stata Journal*, 8(4), 520–531, 2008.
- West, B.T., Welch, K.B., and Galecki, A.T., *Linear Mixed Models: A Practical Guide Using Statistical Software*, Chapman & Hall / CRC Press, Boca Raton, FL, 2007.
- Westat, Inc., *WesVar 4.0 User's Guide*, Westat, Rockville, MD, 2000.
- Wolter, K.M., *Introduction to Variance Estimation*, 2d ed., Springer-Verlag, New York, 2007.
- Woodruff, R.S., A simple method for approximating the variance of a complicated estimate, *Journal of the American Statistical Association*, 66, 411–414, 1971.
- Yamaguchi, K., *Event History Analysis*, Sage Publications, Newbury Park, CA, 1991.
- Yates, F., *Sampling Methods for Censuses and Surveys*, Griffin, London, 1949 (2d ed., 1953; 3d ed., 1960).
- Zajacova, A., Dowd, J.B., and Aiello, A.E., Socioeconomic and race/ethnic patterns in persistent infection burden among U.S. adults, *Journal of Gerontology A: Biological Sciences and Medical Sciences*, 64A(2), 272–279, 2009.
- Zheng, H. and Little, R.J.A., Inference for the population total from probability-proportional-to-size samples based on predictions from a penalized spline non-parametric model, *Journal of Official Statistics*, 21(1), 1–20, 2005.