# SAMPLING METHODOLOGIES
## WITH APPLICATIONS

## Poduri S.R.S. Rao

*Professor of Statistics*
*University of Rochester*
*Rochester, New York*

## CHAPMAN & HALL/CRC

# Contents

# List of tables

# List of figures

*To Durgi, Ann, and Pal*

# Preface

This book contains the methodologies and applications of commonly used procedures of sampling from finite populations. Easy access of the different topics to students and their quick availability to practitioners have been the objectives in preparing this book. With these objectives, the main results along with their logical explanations are presented in the text, but the related derivations are placed in the appendixes to the chapters.

The introductory chapter describes the differences between probability and nonprobability sampling, and presents illustrations of several types of national and international surveys. Properties of simple random sampling along with the estimation of the population mean, total, and variance are considered in the second chapter. Unbiasedness of the sample mean, its variance, and standard error are defined and illustrated. Sample sizes needed to estimate the population mean or total with specified criteria are also examined in this chapter.

The third chapter contains a variety of topics of importance in sample surveys and for statistical analysis in general. Definitions of the bias, variance, and mean square error of an estimator and their relevance to its precision and accuracy are presented and illustrated first. Further topics include the covariance and correlation between two random variables and their sample means, tests of hypotheses related to the means, and the comparison of simple random and systematic random sampling procedures. The effect of the bias arising from the nonresponse in a survey is briefly examined in this chapter, and this topic is continued in the exercises of some of the following chapters and finally in detail in Chapter 11. The appendix to this chapter contains the basic results on conditional and unconditional expectation and variance, and other topics of importance.

Chapters 4 through 11 contain detailed presentations of the estimation of percentages and counts, stratification, subpopulations, single and two-stage cluster sampling with equal and unequal probabilities,

ratio and regression methods of estimation, and the problem of nonresponse and its remedies. The final chapter presents in detail linearization, jackknife, bootstrap, and balanced repeated replication procedures. Major results on small-area estimation and complex surveys appear in this chapter.

Each topic is presented with illustrations, followed by examples and exercises. They are constructed from data on everyday practical situations covering a wide variety of subjects ranging from scholastic aptitude tests to health care expenditures and presidential elections. The examples and exercises are interwoven throughout different chapters.

For the sake of good comprehension of the results on unbiasedness, standard error, and mean square error of the estimators, some of the exercises are constructed as projects that require the selection of all the possible samples of a specified size from a finite population. For some of the exercises and projects, Minitab, Excel, and similar computer programs are quite adequate. The methodological type of questions at the end of the exercises along with the derivations in the appendixes should be of particular interest to advanced students. Solutions to the odd-numbered exercises are presented at the end of the book. Detailed solutions to all the exercises appear in the *Teacher's Manual*.

This book can be recommended as a text for a one-semester or two-quarter course for students in statistics and also in business, political and social sciences, and other areas. One or two courses in basic theoretical and applied statistical concepts would provide the required preparation. This book can also serve as a reference guide for survey practitioners, political and public pollsters, and researchers in some industries.

**Poduri S.R.S. Rao**
University of Rochester, New York