

SURVEY SAMPLING

PARIMAL MUKHOPADHYAY



Alpha
Science

0212.2
M953

SURVEY SAMPLING

Parimal Mukhopadhyay



E2009003566

Alpha Science International Ltd.
Oxford, U.K.

Parimal Mukhopadhyay

Professor Indian Statistical Institute (Retd.)
203 B.T. Road, Kolkata, India

Copyright © 2007

Alpha Science International Ltd.
7200 The Quorum, Oxford Business Park North
Garsington Road, Oxford OX4 2JZ, U.K.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, without prior written permission of the publisher.

Printed from the camera-ready copy provided by the Author.

ISBN-13: 978-1-84265-374-6

ISBN-10: 1-84265-374-1

Printed in India

SURVEY SAMPLING

Dedicated to the memory of my uncle
Late Tripureshwar Mukhopadhyay

Preface

This monograph makes a comprehensive review of some of the developments in the theory of survey sampling in the superpopulation model-based approach, starting from an assessment of the situation in the classical fixed-population area.

Chapters 1 and 2 are based on the model of fixed population. Chapter 1 introduces the preliminary concepts, sampling designs, estimators, sampling strategies, various classical estimators, etc. Chapter 2 addresses some inferential problems, e.g., uniformly minimum variance unbiased estimation, admissibility, sufficiency, minimax estimation in survey sampling. The remaining chapters have developed from the assumption of some superpopulation models depicting the survey population.

Chapter 3 considers model-dependent optimal strategies in the prediction-theoretic approach. Chapter 4 deals with the robustness of these strategies under specific model-failures; Chapter 5, the class of strategies which combine randomization both due to sampling designs and the superpopulation models. Chapter 6 addresses the asymptotic properties of these strategies and their robustness.

The following chapter examines the robustness of model-dependent and modelbased strategies in the asymptotic sense. This chapter also identifies regression superpopulation models for which the prediction-estimators become robust. Biasrobust estimation including non-parametric calibration are also discussed. In the design-based conditional approach of Chapter 8, inference is restricted to a part of the sample space which satisfies certain properties. Conditionally optimum estimators are studied in this light. The next chapter addresses design-based calibration estimation of finite population parameters under different distance functions. The concepts of calibration with restricted weights, extended calibration, mitigated calibration are examined. Model-based calibration is discussed in the following chapter and its optimality investigated. The concluding chapter considers yet another approach to estimation, empirical maximum likelihood estimation in finite population sampling. The performance of this approach vis-a-vis calibration approach is examined.

As has been noted in Chapter 2, the arguments based on a fixed-population model do not lead to any optimality result in general. In a broader perspective, survey population can be looked upon as a realization of a superpopulation and many

decisive results can be attained in this set-up. The thrust of the book is, therefore, on superpopulation model-based inference. One important finding is that the generalized regression estimator (*greg*) plays a prominent part, specially, in large scale sample surveys.

This book does not cover, among others, Bayesian approaches in survey sampling, model-based variance estimation and an important application, small area estimation. These have been covered in some details in Mukhopadhyay (2000a, 1996, 1998b) respectively.

The book, to some extent, may be considered as an up-to-date version of, but not restricted to, works in Mukhopadhyay (1996). However, many useful results of the earlier book have not been revisited here.

In writing this book I have attempted to arrange and reconcile the results systematically in a lucid manner and indicate new research areas. Various examples and supplementary exercises have been added to clarify the ideas. We have assumed that the reader has a basic degree in Statistics and is acquainted with the developments in survey sampling at the level of Cassel, *et al.* and Mukhopadhyay (1998a). The book can not be a stand-alone text book for the fixed-population part, but may serve as a self-contained study material for the model-based part, generally of use to the researchers.

The book was partially written during my assignment in the University of South Africa, Pretoria. My family helped me a lot by silent inspiration. An acknowledgement is due to my daughters-in-law Jayita and Shilpi who assisted me in arranging the manuscripts.

Kolkata, India

Parimal Mukhopadhyay

Contents

Preface

vii

1. The Preliminaries	1
1.1 Introduction	1
1.2 The Basic Model	1
1.3 Different Types of Sampling Designs	5
1.4 The Estimators	7
1.5 Exercises and Complements	15
2. Some Inferential Problems under Fixed Population Set-up	17
2.1 Introduction	17
2.2 The <i>pdf</i> of Data and Likelihood Function of y	18
2.3 Sufficiency, Rao-Blackwellization	20
2.4 Uniformly Minimum Variance Unbiased Estimation	22
2.5 Admissibility of Estimators	24
2.6 Minimax Strategies	29
2.7 Average Variance of a Strategy under a Superpopulation Model	33
2.7.1 Superpopulation model	33
2.7.2 Average variance under ξ	34
2.8 Some Asymptotic Results	38
2.9 Exercises and Complements	41
3. Model-Dependent Optimal Strategies	45
3.1 Introduction	45
3.2 Principles of Inference Based on Theory of Prediction	46
3.3 Prediction under Polynomial Regression Models	49
3.4 Prediction under Multiple Regression Models	52
3.5 Predicting a Superpopulation Mean	56
3.6 Reconciling \hat{T}^* and $\hat{\hat{T}}^*$	58
3.7 Exercises and Complements	60

4. Robustness of Model-Dependent Optimal Strategies	63
4.1 Introduction	63
4.2 Bias and MSE under Alternative Models	63
4.3 Bias of \hat{T}_g^*	64
4.4 Mean Square Errors of \hat{T}_g^*	68
4.5 Balanced Samples and $\hat{T}^*(0, 1; v(x))$	69
4.6 Efficiency of a Balanced Sample	72
4.7 Approximately Balanced Samples	74
4.8 A Post-Sample Robust Predictor	75
4.8.1 Empirical Studies	76
4.9 Robustness of $\hat{T}^*(X, v)$ under Alternative Models	77
4.10 Exercises and Complements	79
5. Model-Assisted Sampling Strategies	81
5.1 Introduction	81
5.2 Different Choices of $\hat{\beta}$ under $\xi(X, v)$	82
5.3 Generalized Regression Predictor for Different Choices of $\hat{\beta}$	83
5.4 Simple Variance Estimators of \hat{T}_{gr}	86
5.4.1 WSR form for stratified samples	88
5.5 (p, Q, R) strategies	92
5.5.1 Prediction-error and variance of $\hat{\theta}(Q, R)$	94
5.5.2 Bias of $\hat{T}(p, Q, R)$ under model mis-specification	95
5.6 R-Generalized Strategies	95
5.7 Extended (p, Q, R) Sampling Strategies	96
5.7.1 Optimum EQR predictor	97
5.8 Exercises and Complements	98
6. Asymptotically Optimum Sampling Strategies	101
6.1 Introduction	101
6.2 Brewer's ADU Sampling Strategies	102
6.3 ADU-Subclass of Wright's (p, Q, R) Sampling Strategies	106
6.4 Asymptotic MSE of Wright's Strategies	109
6.5 Asymptotic Efficiency of a (p, Q, R) ADU Strategy	111
6.6 ADU-ness under the Set-up of Continuously Increasing Nested Populations	111
6.7 ADU-Subclass of Extended (Q, R) Predictors	115
6.7.1 ADU $\hat{T}_E(Q, R)$ predictors and \hat{T}_{egr} predictors	117

6.8	Discussion	118
6.9	Exercises and Complements	119
7.	Robust Strategies	123
7.1	Introduction	123
7.2	Robustness of Prediction-Estimators	123
7.2.1	Prediction-estimators based on regression models	124
7.2.2	Comparison with greg	127
7.3	Internally Bias-Calibrated Estimators	128
7.3.1	Generalized linear models	130
7.3.2	Stratified models	132
7.4	Bias-Robust Estimation using Nonparametric Calibration	133
7.4.1	Nonparametric Predictor	134
7.4.2	Bias-adjusted nonparametric estimator	136
7.4.3	Choice of a kernel function and the bandwidth	138
7.4.4	Nonparametric estimation of a finite population distribution function	139
7.5	Outlier Robust Estimators	141
7.5.1	An outlier-robust alternative to the BLUE	141
7.5.2	Outlier-robust estimation of distribution function	143
7.6	Exercises and Complements	148
8.	Design-Based Conditional Approach	151
8.1	Introduction	151
8.2	Conditional Design Unbiasedness	152
8.3	Calculation of the SCW-Estimator $\hat{y}_{\pi \eta}$	156
8.4	Multivariate Difference Estimator	159
8.4.1	The optimal difference estimator	160
8.5	Linear Estimators Based on Multivariate Normality Assumption	163
8.6	Comparison with the Generalized Regression Estimator	164
8.7	Extended Generalized Regression Estimator and the Optimum Difference Estimators	167
8.8	Discussion	169
8.9	Exercises and Complements	170
9.	The Design-Based Calibration Approach	173
9.1	Introduction	173
9.2	Calibration Estimators	174
9.3	Calibration Estimators Based on a Functional Form Approach	177
9.4	Calibration Estimators with Restricted Weights	180

9.5	Robust Calibration Estimators	183
9.6	Extended Calibration Estimators	184
9.6.1	Extended calibration estimators with restricted weights	186
9.6.2	Mitigated calibration	187
9.7	Cosmetic and Calibration Estimator	189
9.8	Exercises and Complements	191
10.	Model Based Calibration Approach	197
10.1	Introduction	197
10.2	The Motivation	198
10.3	The Model	199
10.4	Estimation of \bar{y}	200
10.4.1	Optimality of model-calibrated estimation	206
10.5	Estimator of a Finite Population Distribution Function	207
10.6	Estimation of Quadratic Finite Population Function	208
10.6.1	Estimating finite population variance and covariance	210
10.7	Exercises and Complements	212
11.	Empirical Likelihood Approach	215
11.1	Introduction	215
11.2	Empirical Likelihood	215
11.3	Pseudo-Empirical likelihood Function for SRSWOR	217
11.4	Pseudo-Empirical Likelihood Estimation in Complex Surveys	218
11.5	Asymptotic Properties of PEMLE	220
11.6	Empirical likelihood in Stratified Random Sampling	222
11.7	Asymptotic Properties of the Estimators in Stratified Random Sampling	223
11.8	Model-Calibrated Empirical Likelihood Approach	225
11.8.1	Model calibrated PEMLE of \bar{y}	225
11.8.2	Optimality of model-calibrated PEMLE approach	226
11.8.3	Model calibrated PEMLE of $F_N(t)$	226
11.9	Model-Calibrated PEMLE of Quadratic Functions	227
11.9.1	Estimating variances and covariances under a linear model	228
11.10	Exercises and Complements	229
	References	231
	Author Index	249
	Subject Index	253

Chapter 1

The Preliminaries

1.1 Introduction

Sample survey, finite population sampling or survey sampling is a method of drawing inference about the characteristic of a finite population by observing only a part of the population. Different statistical techniques have been developed to cater to these needs during the last few decades.

In this chapter we examine a model for a fixed finite population which formed the basis of the earlier exposition of the theory. The set-up established in this chapter would be fundamental to our discussion mainly up to Chapter 2.

1.2 The Basic Model

We assume that we have a finite population of distinct units, with a known population size and a variable of interest taking real values on these units. In an enumerative survey the primary task of the survey statistician is to estimate some descriptive characteristics of the population, e.g., population total, mean, variance by suitably choosing a subset (sample) of the population and observing the values of this variable only on the units in the selected subset. (In analytic surveys we consider estimation of superpopulation parameters and these will be considered in Chapter 3 onwards. For the time being we consider the fixed population model and the associated enumerative surveys).

To formulate the basic fixed population model precisely let us consider a few definitions.

DEFINITION 1.2.1: A finite (survey) population \mathcal{P} is a collection of a known number N

N of identifiable units labelled $1, \dots, i, \dots, N$, $\mathcal{P} = \{1, \dots, i, \dots, N\}$, where i stands for the physical unit labeled i .

The above definition excludes from its coverage the populations of the following types: batches of industrial products of the same specification coming out from a production process as the units are not distinguishable individually; population of fishes in a lake as the population size is unknown. Collection of households in an area, industrial units in an urban complex, agricultural fields in a village are examples of survey populations.

Let y be a study variable having value y_i on $i = 1, \dots, N$. Associated with \mathcal{P} we have a vector $\mathbf{y} = (y_1, \dots, y_N)$ which constitutes the parameter for the model of a survey population, $\mathbf{y} \in \mathcal{R}^N$, the parameter space. One is often interested in estimating a parameter function $\theta(\mathbf{y})$, e.g., population total, $T = \sum_{i=1}^N y_i$, population mean, $\bar{y} = \sum_{i=1}^N y_i / N = T/N$, population variance $S^2 = (N-1)^{-1} \sum_{i=1}^N (y_i - \bar{y})^2$ by choosing a sample (a part of the population, defined below) from \mathcal{P} and observing the values of y only on the units in the sample.

DEFINITION 1.2.2: A sample is a part of the population.

A sample may be selected with replacement (*wr*) or without replacement (*wor*) of the units already selected to the original population.

A sample when selected by a *wr*-sampling procedure may be written as a sequence,

$$S = \{i_1, \dots, i_n\}, \quad 1 \leq i_t \leq N, \quad t = 1, \dots, n, \quad (1.2.1)$$

where i_t denotes the label of the unit selected at the t th draw and is not necessarily equal to $i_{t'}$ for $t \neq t'$ ($t = 1, \dots, n$). For a without replacement sampling procedure, a sample when written as a sequence, is

$$S = \{i_1, \dots, i_n\}, \quad 1 \leq i_t \leq N, i_t \neq i_{t'} \text{ for } t \neq t' (t = 1, \dots, n) \quad (1.2.2)$$

since repetition of units in S is not possible. Arranging the units in the sample S in an increasing (decreasing) order of magnitudes of labels and considering only the distinct units, a sample may also be written as a set s . For a *wr*-sampling of n draws, a sample written as a set is, therefore,

$$s = (j_1, \dots, j_{\nu(S)}), \quad 1 \leq j_1 < \dots < j_{\nu(S)} \leq N \quad (1.2.3)$$

where $\nu(S)$ is the number of distinct units in S . In a *wor*-sampling procedure, a sample of n draws, written as a set is,

$$s = (j_1, \dots, j_n), \quad 1 \leq j_1 < \dots < j_n \leq N. \quad (1.2.4)$$

Thus, if in a *wr*-sampling, $S = \{2, 5, 2, 1\}$, the corresponding s is $s = (1, 2, 5)$ with $\nu(S) = 3$. Similarly, if for a *wor*-sampling procedure, $S = \{3, 7, 1\}$, the corresponding s is $s = (1, 3, 7)$. Clearly, information on the order of selection and repetition of units in the sample S is not available in s .

DEFINITION 1.2.3: Number of distinct units in a sample is its effective sample size. In (1.2.3), $\nu(S)$ is the effective sample size, $1 \leq \nu(S) \leq n$. For a *wor*-sample of n draws, $\nu(S) = \nu(s) = n$.

Note that a sample is a sequence or set of some units from the population and does not include their y -values.

DEFINITION 1.2.4: The sample space is the collection of all possible samples and is often denoted as \mathcal{S} . Thus $\mathcal{S} = \{S\}$ or $\{s\}$ according as we are interested in S or s .

In a simple random sample with replacement (*srswr*) of n draws \mathcal{S} contains N^n samples S . In a simple random sample without replacement (*srswor*) of n draws \mathcal{S} contains $(N)_n$ samples S and $\binom{N}{n}$ samples s where $(a)_b = a(a-1)\dots(a-b+1)$, $a > b$. If the samples s of all possible sizes are considered in a *wor*-sampling procedure, there are 2^N samples in \mathcal{S} .

DEFINITION 1.2.5: Let \mathcal{A} be the minimal σ -field over \mathcal{S} and p a probability measure defined over \mathcal{A} such that $p(s)$ [or $p(S)$] denotes the probability of selecting s [or S], satisfying

$$\begin{aligned} p(s) [p(S)] &\geq 0 \\ \sum_{s \in \mathcal{S}} p(s) [\sum_{S \in \mathcal{S}} p(S)] &= 1. \end{aligned} \quad (1.2.5)$$

One of the main tasks of the survey statistician is to find a suitable $p(s)$ or $p(S)$. The collection (\mathcal{S}, p) is called a sampling design, often denoted as $D(\mathcal{S}, p)$ or simply p . The triplet $(\mathcal{S}, \mathcal{A}, p)$ is the probability space for the model of the finite population.

The expected effective sample size of a sampling design p is

$$E\{\nu(S)\} = \sum_{S \in \mathcal{S}} \nu(S) p(S) = \sum_{\mu=1}^N \mu P[\nu(S) = \mu] = \nu. \quad (1.2.6)$$

We shall denote by ρ_n the class of all fixed effective size $[FES(\nu)]$ designs, i.e.

$$\rho_n = \{p : p(s) > 0 \Rightarrow \nu(S) = \nu\}. \quad (1.2.7)$$

A sampling design p is said to be noninformative if $p(s)[p(S)]$ does not depend on the y -values. In this treatise, unless stated otherwise, we shall consider non-informative designs only. Informative designs have been considered by Basu (1969), Zacks (1969), Liao and Sedransk (1975), Stenger (1977), Bethlehem and Schuerhoff (1984), among others.

Basu (1958), Basu and Ghosh (1967) proved that all the information relevant to making inference about the population characteristic is contained in the set sample s and the corresponding y -values (Theorem 2.3.1). As such, unless otherwise stated, we shall consider samples as sets s only.

The quantities

$$\begin{aligned}\pi_i &= \sum_{s \ni i} p(s), \quad \pi_{ij} = \sum_{s \ni (i,j)} p(s) \\ \pi_{i_1 \dots i_k} &= \sum_{s \ni (i_1, \dots, i_k)} p(s)\end{aligned}\tag{1.2.8}$$

are, respectively, the first order, second order, ..., k th order inclusion-probabilities of units in a sampling design p . The following lemma depicts some relations among inclusion-probabilities and expected effective sample size of a sampling design.

Lemma 1.2.1: For any sampling design p ,

(i)

$$\pi_i + \pi_j - 1 \leq \pi_{ij} \leq \min(\pi_i, \pi_j)$$

(ii)

$$\sum_{i=1}^N \pi_i = \sum_{s \in S} \nu(s) p(s) = \nu$$

(iii)

$$\sum_{i \neq j=1}^N \pi_{ij} = \nu(\nu - 1) + V\{\nu(s)\}.$$

If $p \in \rho_n$,

(iv)

$$\sum_{j(\neq i)=1}^N \pi_{ij} = (\nu - 1)\pi_i$$

(v)

$$\sum_{i \neq j=1}^N \pi_{ij} = \nu(\nu - 1).$$

Result (i) is obvious. Results (ii), (iii) and (iv), (v) are, respectively, due to Godambe (1955), Hanurav (1962), Yates and Grundy (1953).

Further, for any sampling design p ,

$$\theta(1 - \theta) \leq V\{\nu(s)\} \leq (N - \nu)(\nu - 1)\tag{1.2.9}$$

where $\nu = [\nu] + \theta$, $0 \leq \theta < 1$, θ being the fractional part of ν . The lower bound in (1.2.9) is attained by a sampling design for which

$$P[\nu(S) = [\nu]] = 1 - \theta \quad \text{and} \quad P[\nu(S) = \nu + 1] = \theta.$$

Mukhopadhyay (1975) derived a sampling design with fixed nominal sample size $n(> \nu)$, $[p(S) > 0 \Rightarrow n(S) = n \forall S]$ such that $V\{\nu(S)\} = \theta(1 - \theta/(n - [\nu]))$, which is very

close to the lower bound in (1.2.9). Here, $n(S)$ is the nominal sample size (number of draws in) S .

We shall denote by

$p_r(i_r)$ = probability of selecting i_r at the r th draw

$p_r(i_r|i_1, \dots, i_{r-1})$ = conditional probability of selecting i_r at the r th draw given that i_1, \dots, i_{r-1} are selected at the first, \dots , $(r-1)$ th draws respectively.

Suppose the values x_1, \dots, x_N of a closely related (to y) auxiliary variable x on units $1, 2, \dots, N$ respectively, are available. As an example, in an agricultural survey, x may be the area of a plot under a specified crop and y the yield of crop on that plot. The quantities $p_i = x_i/X$, $X = \sum_{i=1}^N x_i$ is called the size-measure of unit i ($i = 1, \dots, N$) and is often used in selection of samples.

1.3 Different Types of Sampling Designs

Sampling designs proposed in the literature can be generally grouped in the following categories:

- (a) Simple random sampling with replacement (*srswr*)
- (b) Simple random sampling without replacement (*srswor*)
- (c) Probability proportional to size with replacement (*ppswr*) sampling: a unit i is selected with probability p_i at the r th draw and a unit once selected is returned to the population before the next draw ($r = 1, 2, \dots$).
- (d) Unequal probability without replacement (*upwor*) sampling: A unit i is selected at the r th draw with probability proportional to $p_i^{(r)}$ and a unit once selected is removed from the population. Here

$$p_1(i) = p_i^{(1)}$$

$$p_r^{(i_r)} = \frac{p_{i_r}^{(r)}}{1 - p_{i_1}^{(r)} - p_{i_2}^{(r)} - \dots - p_{i_{r-1}}^{(r)}}, r = 1, 2, \dots, n. \quad (1.3.1)$$

The quantities $\{p_i^{(r)}\}$ are generally functions of p_i and the p_i -values of the units already selected. In particular, if $p_i^{(r)} = p_i \forall i = 1, \dots, N$, the procedure may be called *probability proportional to size without replacement (ppswor)* sampling procedure. For $n = 2$, for this procedure

$$\pi_i = p_i \left[1 + A - \frac{p_i}{1 - p_i} \right]$$

$$\pi_{ij} = p_i p_j \left[\frac{1}{1 - p_i} + \frac{1}{1 - p_j} \right], \text{ where } A = \sum_{k=1}^N \frac{p_k}{1 - p_k}.$$

The sampling design may also be attained by an inverse sampling procedure where units are drawn *wr*, with probability $p_i^{(r)}$ at the *r*th draw, until for the first time *n* distinct units occur. The *n* distinct units each taken only once constitute the sample.

- (e) Rejective sampling: Drwas are made *wr* and with probability $\{p_i^{(r)}\}$ at the *r*th draw. If all the units turn out distinct, the solution is taken as a sample; otherwise, the whole sample is rejected and fresh dras are made. In some situations $p_i^{(r)} = p_i \forall i$.
- (f) Systematic sampling with varying probability (including equal probability).
- (g) Sampling from groups: The population is divided into *L* groups either at random or following some suitable procedures and a sample of size n_h is drwan from the *h*th group by using any of the above-mentioned sampling dsigns such that the desired sample size $n = \sum_{h=1}^L n_h$ is attained. An example is the Rao-Hartley-Cochran (1962) sampling procedure.

Based on the above methods, there are many uni-stage or multi-stage stratified sampling procedures.

A *FES*(*n*)-sampling design with π_i proportional to p_i is often used for estimating a population total. This is, because an important estimator, the Horviyz-Thompson estimator (HTE) has zero variance if y_i is proportional to p_i . Such a design is called a *pps* design or *IPPS* (inclusion-probability proportional to size) design. Since $\pi_i \leq 1$, it is required that $x_i \leq X/n \forall i$ for such a design.

Many (exceeding sixty) unequal probability without replacement sampling designs have been suggested in the literature, mostly for use along with the HTE. For many of these designs sample size is a variable. Again, some of these designs are sequential in nature (e.g., Chao (1982), Sunter (1977)). Mukhopadhyay (1972), Sinha (1973), Herzel (1986) considered the problem of realizing a sampling design with pre-assigned sets of inclusion-probabilities of first two orders.

In a sample survey, all the possible samples are not generally equally preferable from the view-point of practical advantages. In agricultural surveys, for example, the investigators tend to avoid grids which are located further away from the cell camps, are located in marshy land, inaccessible places, etc. In such cases, the sampler would like to use only a fraction of totality of all possible samples, allotting only a small probability to the non-preferred units. Such designs are called *Controlled Sampling Designs* and have been considered by several authors (e.g., Chakravorty (1963), Srivastava and Saleh (1985), Rao and Nigam (1989, 1990), Mukhopadhyay and Vijayan (1996)).